

MIND DESIGN

Philosophy
Psychology
Artificial Intelligence

edited by

JOHN HAUGELAND

A Bradford Book

The MIT Press
Cambridge, Massachusetts
London, England

10

Minds, Brains, and Programs

JOHN R. SEARLE

Abstract: I distinguish between strong and weak AI. According to strong AI, appropriately programmed computers literally have cognitive states, and therefore the programs are psychological theories. I argue that strong AI must be false, since a human agent could instantiate the program and still not have the appropriate mental states. I examine some arguments against this claim, and I explore some consequences of the fact that human and animal brains are the causal bases of existing mental phenomena.

WHAT PSYCHOLOGICAL and philosophical significance should we attach to recent efforts at computer simulations of human cognitive capacities? In answering this question I find it useful to distinguish what I will call "strong" AI from "weak" or "cautious" AI. According to weak AI, the principal value of the computer in the study of the mind is that it gives us a very powerful tool. For example, it enables us to formulate and test hypotheses in a more rigorous and precise fashion than before. But according to strong AI the computer is not merely a tool in the study of the mind; rather the appropriately programmed computer really is a mind in the sense that computers given the right programs can be literally said to *understand* and have other cognitive states. And, according to strong AI, because the programmed computer has cognitive states, the programs are not mere tools that enable us to test psychological explanations; rather, the programs are themselves

the explanations. I have no objection to the claims to weak AI, at least as far as this article is concerned. My discussion here will be directed to the claims I have defined as strong AI, specifically the claim that the appropriately programmed computer literally has cognitive states and that the programs thereby explain human cognition. When I refer to AI, it is the strong version as expressed by these two claims which I have in mind.

I will consider the work of Roger Schank and his colleagues at Yale (cf. Schank and Abelson, 1977), because I am more familiar with it than I am with any similar claims, and because it provides a clear example of the sort of work I wish to examine. But nothing that follows depends upon the details of Schank's programs. The same arguments would apply to Winograd's (1972) SHRDLU, Weizenbaum's (1965) ELIZA, and, indeed, any Turing machine simulation of human mental phenomena.

Briefly and leaving out the various details, one can describe Schank's program as follows: the aim of the program is to simulate the human ability to understand stories. It is characteristic of the abilities of human beings to understand stories that they can answer questions about the story even though the information they give was not explicitly stated in the story. Thus, for example, suppose you are given the following story: "A man went into a restaurant and ordered a hamburger. When the hamburger arrived, it was burned to a crisp, and the man stormed out of the restaurant angrily without paying for the hamburger or leaving a tip." Now, if you are given the question "Did the man eat the hamburger?", you will presumably answer, "No, he did not." Similarly if you are given the following story: "A man went into a restaurant, and ordered a hamburger; when the hamburger came, he was very pleased with it; and as he left the restaurant he gave the waitress a large tip before paying his bill," and you are asked the question "Did the man eat the hamburger?", you will presumably answer, "Yes, he ate the hamburger." Now Schank's machines can similarly answer questions about restaurants in this fashion. In order to do so, they have a "representation" of the sort of information that human beings have about restaurants which enables them to answer such questions as those above, given these sorts of stories. When the machine is given the story and then asked the question, the machine will print out answers of the sort that we

would expect human beings to give if told similar stories. Partisans of strong AI claim that in this question-and-answer sequence, not only is the machine simulating a human ability but also:

- (a) the machine can literally be said to *understand* the story and provide answers to questions; and
- (b) what the machine and its program do *explains* the human ability to understand the story and answer questions about it.

Claims (a) and (b) seem to me totally unsupported by Schank's work, as I will attempt to show in what follows.¹

A way to test any theory of the mind is to ask oneself what it would be like if one's own mind actually worked on the principles that the theory says all minds work on. Let us apply this test to the Schank program with the following *Gedankenexperiment*. Suppose that I'm locked in a room and suppose that I'm given a large batch of Chinese writing. Suppose furthermore, as is indeed the case, that I know no Chinese either written or spoken, and that I'm not even confident that I could recognize Chinese writing as Chinese writing distinct from, say, Japanese writing or meaningless squiggles. Now suppose further that after this first batch of Chinese writing, I am given a second batch of Chinese script together with a set of rules for correlating the second batch with the first batch. The rules are in English and I understand these rules as well as any other native speaker of English. They enable me to correlate one set of formal symbols with another set of formal symbols, and all that "formal" means here is that I can identify the symbols entirely by their shapes. Now suppose also that I am given a third batch of Chinese symbols together with some instructions, again in English, that enable me to correlate elements of this third batch with the first two batches, and these rules instruct me how I am to give back certain Chinese symbols with certain sorts of shapes in response to certain sorts of shapes given me in the third batch. Unknown to me, the people who are giving me all of these symbols call the first batch 'a script,' they call the second batch a 'story,' and they call the third batch 'questions.'

1. I am not saying, of course, that Schank himself is committed to these claims.

Furthermore, they call the symbols I give them back in response to the third batch 'answers to the questions,' and the set of rules in English that they gave me they call 'the program.' To complicate the story a little bit, imagine that these people also give me stories in English which I understand, and they then ask me questions in English about these stories, and I give them back answers in English. Suppose also that after a while I get so good at following the instructions for manipulating the Chinese symbols and the programmers get so good at writing the programs that from the external point of view—that is, from the point of view of somebody outside the room in which I am locked—my answers to the questions are indistinguishable from those of native Chinese speakers. Nobody looking at my answers can tell that I don't speak a word of Chinese. Let us also suppose that my answers to the English questions are, as they no doubt would be, indistinguishable from those of other native English speakers, for the simple reason that I am a native speaker of English. From the external point of view, from the point of view of someone reading my 'answers,' the answers to the Chinese questions and the English questions are equally good. But in the Chinese case, unlike the English case, I produce the answers by manipulating uninterpreted formal symbols. As far as the Chinese is concerned, I simply behave like a computer; I perform computational operations on formally specified elements. For the purposes of the Chinese, I am simply an instantiation of the computer program.

Now the claims made by strong AI are that the programmed computer understands the stories and that the program in some sense explains human understanding. But we are now in a position to examine these claims in light of our thought experiment.

(a) As regards the first claim it seems to me obvious in the example that I do not understand a word of the Chinese stories. I have inputs and outputs that are indistinguishable from those of the native Chinese speaker, and I can have any formal program you like, but I still understand nothing. Schank's computer for the same reasons understands nothing of any stories whether in Chinese, English, or whatever, since in the Chinese case the computer is me; and in cases where the computer is not me, the computer has nothing more than I have in the case where I understand nothing.

(b) As regards the second claim—that the program explains human understanding—we can see that the computer and its program do not provide sufficient conditions of understanding, since the computer and the program are functioning and there is no understanding. But does it even provide a necessary condition or a significant contribution to understanding? One of the claims made by the supporters of strong AI is this: when I understand a story in English, what I am doing is exactly the same—or perhaps more of the same—as what I was doing in the case of manipulating the Chinese symbols. It is simply more formal symbol manipulation which distinguishes the case in English, where I do understand, from the case in Chinese, where I don't. I have not demonstrated that this claim is false, but it would certainly appear an incredible claim in the example. Such plausibility as the claim has derives from the supposition that we can construct a program that will have the same inputs and outputs as native speakers, and in addition we assume that speakers have some level of description where they are also instantiations of a program. On the basis of these two assumptions, we assume that even if Schank's program isn't the whole story about understanding, maybe it is part of the story. That is, I suppose, an empirical possibility, but not the slightest reason has so far been given to suppose it is true, since what is suggested—though certainly not demonstrated—by the example is that the computer program is irrelevant to my understanding of the story. In the Chinese case I have everything that artificial intelligence can put into me by way of a program, and I understand nothing; in the English case I understand everything, and there is so far no reason at all to suppose that my understanding has anything to do with computer programs—i.e., with computational operations on purely formally specified elements. As long as the program is defined in terms of computational operations on purely formally defined elements, what the example suggests is that these by themselves have no interesting connection with understanding. They are certainly not sufficient conditions, and not the slightest reason has been given to suppose that they are necessary conditions or even that they make a significant contribution to understanding. Notice that the force of the argument is not simply that different machines can have the same input and output while operating on different formal principles—that is not

the point at all—but rather that whatever purely formal principles you put into the computer will not be sufficient for understanding, since a human will be able to follow the formal principles without understanding anything, and no reason has been offered to suppose they are necessary or even contributory, since no reason has been given to suppose that when I understand English, I am operating with any formal program at all.

What is it, then, that I have in the case of the English sentences which I do not have in the case of the Chinese sentences? The obvious answer is that I know what the former mean but haven't the faintest idea what the latter mean. In what does this consist, and why couldn't we give it to a machine, whatever it is? Why couldn't the machine be given whatever it is about me that makes it the case that I know what English sentences mean? I will return to these questions after developing my example a little more.

I have had occasions to present this example to several workers in artificial intelligence and, interestingly, they do not seem to agree on what the proper reply to it is. I get a surprising variety of replies, and in what follows I will consider the most common of these (specified along with their geographical origins). First I want to block out some common misunderstandings about "understanding." In many of these discussions one finds fancy footwork about the word "understanding." My critics point out that there are different degrees of understanding, that "understanding" is not a simple two-place predicate, that there are even different kinds and levels of understanding, and often the law of the excluded middle doesn't even apply in a straightforward way to statements of the form 'x understands y,' that in many cases it is a matter for decision and not a simple matter of fact whether x understands y. And so on. To all these points I want to say of course, of course; but they have nothing to do with the points at issue. There are clear cases where "understanding" applies and clear cases where it does not apply; and such cases are all I need for this argument.² I understand stories in English; to a lesser

2. Also, "understanding" implies both the possession of mental (intentional) states and the truth (validity, success) of these states. For the purposes of this discussion, we are concerned only with the possession of the states.

degree I can understand stories in French; to a still lesser degree, stories in German; and in Chinese, not at all. My car and my adding machine, on the other hand, understand nothing; they are not in that line of business. We often attribute "understanding" and other cognitive predicates by metaphor and analogy to cars, adding machines, and other artifacts, but nothing is proved by such attributions. We say, "The door *knows* when to open because of its photoelectric cell," "the adding machine *knows how* (understands how, is able) to do addition and subtraction but not division," and "the thermostat *perceives* changes in the temperature." The reason we make these attributions is interesting and has to do with the fact that in artifacts we extend our own intentionality;³ our tools are extensions of our purposes, and so we find it natural to make metaphorical attributions of intentionality to them; but I take it no philosophical ice is cut by such examples. The sense in which an automatic door "understands instructions" from its photoelectric cell is not at all the sense in which I understand English. If the sense in which Schank's programmed computers understand stories is supposed to be the metaphorical sense in which the door understands, and not the sense in which I understand English, the issue would not be worth discussing. Newell and Simon write that the sense of "understanding" they claim for computers is exactly the same as for human beings. I like the straightforwardness of this claim, and it is the sort of claim I will be considering. I will argue that in the literal sense the programmed computer understands what the car and the adding machine understand, viz. exactly nothing. The computer understanding is not just (like my understanding of German) partial or incomplete; it is zero.

Now to the replies:

1. *The Systems Reply* (Berkeley): "While it is true that the individual person who is locked in the room does not understand the story, the fact is that he is merely part of a whole system

3. Intentionality is by definition that feature of certain mental states by which they are directed at or about objects and states of affairs in the world. Thus beliefs, desires, and intentions are intentional states; undirected forms of anxiety and depression are not. For further discussion, see Searle (1979).

and the system does understand the story. The person has a large ledger in front of him in which are written the rules, he has a lot of scratch paper and pencils for doing calculations, he has "data banks" of sets of Chinese symbols. Now, understanding is not being ascribed to the mere individual, rather it is being ascribed to this whole system of which he is a part."

My response to the systems theory is simple: Let the individual internalize all of these elements of the system. He memorizes the rules in the ledger and the data banks of Chinese symbols, and he does all the calculations in his head. The individual then incorporates the entire system. There isn't anything at all to the system which he does not encompass. We can even get rid of the room and suppose he works outdoors. All the same, he understands nothing of the Chinese, and *a fortiori* neither does the system, because there isn't anything in the system which isn't in him. If he doesn't understand, then there is no way the system could understand because the system is just a part of him.

Actually I feel somewhat embarrassed even to give this answer to the systems theory because the theory seems to me so implausible to start with. The idea is that while a person doesn't understand Chinese, somehow the *conjunction* of that person and bits of paper might understand Chinese. It is not easy for me to imagine how someone who was not in the grip of an ideology would find the idea at all plausible. Still, I think many people who are committed to the ideology of strong AI will in the end be inclined to say something very much like this; so let us pursue it a bit further. According to one version of this view, while the man in the internalized systems example doesn't understand Chinese in the sense that a native Chinese speaker does (because, for example, he doesn't know that the story refers to restaurants and hamburgers, etc.), still "the man as formal symbol manipulation system" *really does understand Chinese*. The subsystem of the man which is the formal symbol manipulation system for Chinese should not be confused with the subsystem for English.

So there are really two subsystems in the man; one understands English, the other Chinese, and "it's just that the two systems have little to do with each other." But, I want to reply, not only do they have little to do with each other, they are not even

remotely alike. The subsystem that understands English (assuming we allow ourselves to talk in this jargon of "subsystems" for a moment) knows that the stories are about restaurants and eating hamburgers, etc.; he knows that he is being asked questions about restaurants and that he is answering questions as best as he can by making various inferences from the content of the story, and so on. But the Chinese system knows none of this; whereas the English subsystem knows that "hamburgers" refers to hamburgers, the Chinese subsystem knows only that "squiggle-squiggle" is followed by "squoggle-squoggle." All he knows is that various formal symbols are being introduced at one end and are manipulated according to rules written in English, and that other symbols are going out at the other end. The whole point of the original example was to argue that such symbol manipulation by itself couldn't be sufficient for understanding Chinese in any literal sense because the man could write "squoggle-squoggle" after "squiggle-squiggle" without understanding anything in Chinese. And it doesn't meet that argument to postulate subsystems within the man, because the subsystems are no better off than the man was in the first place; they still don't have anything even remotely like what the English-speaking man (or subsystem) has. Indeed, in the case as described, the Chinese subsystem is simply a part of the English subsystem, a part that engages in meaningless symbol manipulation according to rules in English.

Let us ask ourselves what is supposed to motivate the systems reply in the first place—that is, what *independent* grounds are there supposed to be for saying that the agent must have a subsystem within him which literally understands stories in Chinese? As far as I can tell, the only grounds are that in the example I have the same input and output as native Chinese speakers, and a program that goes from one to the other. But the point of the example has been to show that that couldn't be sufficient for understanding, in the sense in which I understand stories in English, because a person, hence the set of systems that go to make up a person, could have the right combination of input, output, and program and still not understand anything in the relevant literal sense in which I understand English. The only motivation for saying there *must* be a subsystem in me which understands Chinese is that I have a program and I can pass the Turing test;

I can fool native Chinese speakers (cf. Turing, 1950). But precisely one of the points at issue is the adequacy of the Turing test. The example shows that there could be two "systems" both of which pass the Turing test but only one of which understands; and it is no argument against this point to say that since they both pass the Turing test, they must both understand, since this claim fails to meet the argument that the system in me which understands English has a great deal more than the system which merely processes Chinese. In short the systems reply simply begs the question by insisting without argument that the system must understand Chinese.

Furthermore, the systems reply would appear to lead to consequences that are independently absurd. If we are to conclude that there must be cognition in me on the grounds that I have a certain sort of input and output and a program in between, then it looks as though all sorts of noncognitive subsystems are going to turn out to be cognitive. For example, my stomach has a level of description where it does information processing, and it instantiates any number of computer programs, but I take it we do not want to say that it has any understanding. Yet if we accept the systems reply, it is hard to see how we avoid saying that stomach, heart, liver, etc. are all understanding subsystems, since there is no principled way to distinguish the motivation for saying the Chinese subsystem understands from saying that the stomach understands. (It is, by the way, not an answer to this point to say that the Chinese system has information as input and output and the stomach has food and food products as input and output, since from the point of view of the agent, from my point of view, there is no information in either the food or the Chinese; the Chinese is just so many meaningless squiggles. The information in the Chinese case is solely in the eyes of the programmers and the interpreters, and there is nothing to prevent them from treating the input and output of my digestive organs as information if they so desire.)

This last point bears on some independent problems in strong AI, and it is worth digressing for a moment to explain it. If strong AI is to be a branch of psychology, it must be able to distinguish systems which are genuinely mental from those which are not. It must be able to distinguish the principles on which the mind

works from those on which nonmental systems work; otherwise it will offer us no explanations of what is specifically mental about the mental. And the mental-nonmental distinction cannot be just in the eye of the beholder—it must be intrinsic to the systems, for otherwise it would be up to any beholder to treat people as nonmental and, e.g., hurricanes as mental, if he likes. But quite often in the AI literature the distinction is blurred in ways which would in the long run prove disastrous to the claim that AI is a cognitive inquiry. McCarthy, for example, writes: “Machines as simple as thermostats can be said to have beliefs, and having beliefs seems to be a characteristic of most machines capable of problem solving performance” (McCarthy, 1979). Anyone who thinks strong AI has a chance as a theory of the mind ought to ponder the implications of that remark. We are asked to accept it as a discovery of strong AI that the hunk of metal on the wall which we use to regulate the temperature has beliefs in exactly the same sense that we, our spouses, and our children have beliefs, and furthermore that “most” of the other machines in the room—telephone, tape recorder, adding machine, electric light switch, etc.—also have beliefs in this literal sense. It is not the aim of this article to argue against McCarthy’s point, so I will simply assert the following without argument. The study of the mind starts with such facts as that humans have beliefs and thermostats, telephones, and adding machines don’t. If you get a theory that denies this point, you have produced a counter-example to the theory, and the theory is false. One gets the impression that people in AI who write this sort of thing think they can get away with it because they don’t really take it seriously and they don’t think anyone else will either. I propose, for a moment at least, to take it seriously. Think hard for one minute about what would be necessary to establish that that hunk of metal on the wall over there has real beliefs, beliefs with direction of fit, propositional content, and conditions of satisfaction; beliefs that have the possibility of being strong beliefs or weak beliefs; nervous, anxious or secure beliefs; dogmatic, rational, or superstitious beliefs; blind faiths or hesitant cogitations; any kind of beliefs. The thermostat is not a candidate. Neither are stomach, liver, adding machine, or telephone. However, since we are taking the idea seriously, notice that its truth would be fatal to the claim of strong AI to be a

science of the mind, for now the mind is everywhere. What we wanted to know is what distinguishes the mind from thermostats, livers, etc. And if McCarthy were right, strong AI hasn’t a hope of telling us that.

II. *The Robot Reply* (Yale): “Suppose we wrote a different kind of program from Schank’s program. Suppose we put a computer inside a robot, and this computer would not just take in formal symbols as input and give out formal symbols as output, but rather it would actually operate the robot in such a way that the robot does something very much like perceiving, walking, moving about, hammering nails, eating, drinking—anything you like. The robot would, for example, have a television camera attached to it that enabled it to see, it would have arms and legs that enabled it to act, and all of this would be controlled by its computer brain. Such a robot would, unlike Schank’s computer, have genuine understanding and other mental states.”

The first thing to notice about the robot reply is that it tacitly concedes that cognition is not solely a matter of formal symbol manipulation, since this reply adds a set of causal relations with the outside world. But the answer to the robot reply is that the addition of such “perceptual” and “motor” capacities adds nothing by way of understanding, in particular, or intentionality, in general, to Schank’s original program; and to see this, notice that the same thought experiment applies to the robot case. Suppose that instead of the computer inside the robot, you put me inside the room and you give me again, as in the original Chinese case, more Chinese symbols with more instructions in English for matching Chinese symbols to Chinese symbols and feeding back Chinese symbols to the outside. Suppose unknown to me, some of the Chinese symbols that come to me come from a television camera attached to the robot, and other Chinese symbols that I am giving out serve to make the motors inside the robot move the robot’s legs or arms. It is important to emphasize that all I am doing is manipulating formal symbols: I know none of these other facts. I am receiving “information” from the robot’s “perceptual” apparatus, and I am giving out “instructions” to its motor apparatus without knowing either of these facts. I am the

robot's homunculus, but unlike the traditional homunculus, I don't know what's going on. I don't understand anything except the rules for symbol manipulation. Now in this case I want to say that the robot has no intentional states at all; it is simply moving about as a result of its electrical wiring and its program. And furthermore, by instantiating the program, I have no intentional states of the relevant type. All I do is follow formal instructions about manipulating formal symbols.

III. *The Brain Simulator Reply* (Berkeley and M.I.T.): "Suppose we design a program that doesn't represent information that we have about the world, such as the information in Schank's scripts, but simulates the actual sequence of neuron firings at the synapses of the brain of a native Chinese speaker when he understands stories in Chinese and gives answers to them. The machine takes in Chinese stories and questions about them as input, it simulates the formal structure of actual Chinese brains in processing these stories, and it gives out Chinese answers as outputs. We can even imagine that the machine operates not with a single serial program but with a whole set of programs operating in parallel, in the manner that actual human brains presumably operate when they process natural language. Now surely in such a case we would have to say that the machine understood the stories; and if we refuse to say that, wouldn't we also have to deny that native Chinese speakers understood the stories? At the level of the synapses what would or could be different about the program of the computer and the program of the Chinese brain?"

Before addressing this reply, I want to digress to note that it is an odd reply for any partisan of artificial intelligence (functionalism, etc.) to make. I thought the whole idea of strong artificial intelligence is that we don't need to know how the brain works to know how the mind works. The basic hypothesis, or so I had supposed, was that there is a level of mental operations that consists in computational processes over formal elements which constitute the essence of the mental and can be realized in all sorts of different brain processes in the same way that any computer program can be realized in different computer hardware: on the assumptions of strong AI, the mind is to the brain as the program is to

the hardware, and thus we can understand the mind without doing neurophysiology. If we had to know how the brain worked in order to do AI, we wouldn't bother with AI. However, even getting this close to the operation of the brain is still not sufficient to produce understanding. To see that this is so, imagine that instead of a monolingual man in a room shuffling symbols we have the man operate an elaborate set of water pipes with valves connecting them. When the man receives the Chinese symbols he looks up in the program, written in English, which valves he has to turn on and off. Each water connection corresponds to a synapse in the Chinese brain, and the whole system is rigged up so that after doing all the right firings—that is, after turning on all the right faucets—the Chinese answers pop out at the output end of the series of pipes.

Now where is the understanding in this system? It takes Chinese as input, it simulates the formal structure of the synapses of the Chinese brain, and it gives Chinese as output. But the man certainly doesn't understand Chinese, and neither do the water pipes, and if we are tempted to adopt what I think is the absurd view that somehow the *conjunction* of man and water pipes understands, remember that in principle the man can internalize the formal structure of the water pipes and do all the "neuron firings" in his imagination. The problem with the brain simulator is that it is simulating the wrong things about the brain. As long as it simulates only the formal structure of the sequence of neuron firings at the synapses, it won't have simulated what matters about the brain, namely its causal properties, its ability to produce intentional states. And that the formal properties are not sufficient for the causal properties is shown by the water pipe example: we can have all the formal properties carved off from the relevant neurobiological causal properties.

IV. *The Combination Reply* (Berkeley and Stanford): "While each of the previous three replies might not be completely convincing by itself as a refutation of the Chinese room counter-example, if you take all three together they are collectively much more convincing and even decisive. Imagine a robot with a brain-shaped computer lodged in its cranial cavity; imagine the computer programmed with all the synapses of a human brain;

imagine that the whole behavior of the robot is indistinguishable from human behavior; and now think of the whole thing as a unified system and not just as a computer with inputs and outputs. Surely in such a case we would have to ascribe intentionality to the system."

I entirely agree that in such a case we would find it rational and indeed irresistible to accept the hypothesis that the robot had intentionality, as long as we knew nothing more about it. Indeed, besides appearance and behavior the other elements of the combination are really irrelevant. If we could build a robot whose behavior was indistinguishable over a large range from human behavior, we would attribute intentionality to it, pending some reason not to. We wouldn't need to know in advance that its computer brain was a formal analogue of the human brain.

But I really don't see that this is any help to the claims of strong AI, and here is why: According to strong AI, instantiating a formal program with the right input and output is a sufficient condition of, indeed is constitutive of, intentionality. As Newell (1980) puts it, the essence of the mental is the operation of a physical symbol system. But the attributions of intentionality that we make to the robot in this example have nothing to do with formal programs. They are simply based on the assumption that if the robot looks and behaves sufficiently like us, we would suppose until proven otherwise that it must have mental states like ours which cause and are expressed by its behavior, and it must have an inner mechanism capable of producing such mental states. If we knew independently how to account for its behavior without such assumptions, we would not attribute intentionality to it, especially if we knew it had a formal program. And this is the point of my earlier reply to objection II.

Suppose we knew that the robot's behavior was entirely accounted for by the fact that a man inside it was receiving uninterpreted formal symbols from the robot's sensory receptors and sending out uninterpreted formal symbols to its motor mechanisms, and the man was doing this symbol manipulation in accordance with a bunch of rules. Furthermore, suppose the man knows none of these facts about the robot; all he knows is which operations to perform on which meaningless symbols. In such a case we

would regard the robot as an ingenious mechanical dummy. The hypothesis that the dummy has a mind would now be unwarranted and unnecessary, for there is now no longer any reason to ascribe intentionality to the robot or to the system of which it is a part (except of course for the man's intentionality in manipulating the symbols). The formal symbol manipulations go on, the input and output are correctly matched, but the only real locus of intentionality is the man, and he doesn't know any of the relevant intentional states; he doesn't, for example, *see* what comes into the robot's eyes, he doesn't *intend* to move the robot's arm, and he doesn't *understand* any of the remarks made to or by the robot. Nor, for the reasons stated earlier, does the system of which man and robot are a part.

To see the point contrast this case with cases where we find it completely natural to ascribe intentionality to members of certain other primate species, such as apes and monkeys, and to domestic animals, such as dogs. The reasons we find it natural are, roughly, two: we can't make sense of the animal's behavior without the ascription of intentionality, and we can see that the beasts are made of stuff similar to our own—an eye, a nose, its skin, etc. Given the coherence of the animal's behavior and the assumption of the same causal stuff underlying it, we assume both that the animal must have mental states underlying its behavior, and the mental states must be produced by mechanisms made out of the stuff that is like our stuff. We would certainly make similar assumptions about the robot unless we had some reason not to, but as soon as we knew that the behavior was the result of a formal program, and that the actual causal properties of the physical substance were irrelevant, we would abandon the assumption of intentionality.

There are two other responses to my example which come up frequently (and so are worth discussing) but really miss the point.

V. *The Other Minds Reply* (Yale): "How do you know that other people understand Chinese or anything else? Only by their behavior. Now the computer can pass the behavioral tests as well as they can (in principle), so if you are going to attribute cognition to other people, you must in principle also attribute it to computers."

The objection is worth only a short reply. The problem in this discussion is not about how I know that other people have cognitive states, but rather what it is that I am attributing to them when I attribute cognitive states to them. The thrust of the argument is that it couldn't be just computational processes and their output because there can be computational processes and their output without the cognitive state. It is no answer to this argument to feign anesthesia. In "cognitive sciences" one presupposes the reality and knowability of the mental in the same way that in physical sciences one has to presuppose the reality and knowability of physical objects.

VI. *The Many Mansions Reply* (Berkeley): "Your whole argument presupposes that AI is only about analogue and digital computers. But that just happens to be the present state of technology. Whatever these causal processes are that you say are essential for intentionality (assuming you are right), eventually we will be able to build devices that have these causal processes and that will be artificial intelligence. So your arguments are in no way directed at the ability of artificial intelligence to produce and explain cognition."

I have no objection to this reply except to say that it in effect trivializes the project of strong artificial intelligence by redefining it as whatever artificially produces and explains cognition. The interest of the original claims made on behalf of artificial intelligence is that it was a precise, well defined thesis: mental processes are computational processes over formally defined elements. I have been concerned to challenge that thesis. If the claim is redefined so that it is no longer that thesis, my objections no longer apply, because there is no longer a testable hypothesis for them to apply to.

Let us now return to the questions I promised I would try to answer: Granted that in my original example I understand the English and I do not understand the Chinese, and granted therefore that the machine doesn't understand either English or Chinese; still there must be something about me that makes it the case that I understand English and a corresponding something lacking in me which makes it the case that I fail to understand Chinese. Now why couldn't we give those somethings, whatever they are, to a machine?

I see no reason in principle why we couldn't give a machine the capacity to understand English or Chinese, since in an important sense our bodies with our brains are precisely such machines. But I do see very strong arguments for saying that we could not give such a thing to a machine where the operation of the machine is defined solely in terms of computational processes over formally defined elements; that is, where the operation of the machine is defined as an instantiation of a computer program. It is not because I am the instantiation of a computer program that I am able to understand English and have other forms of intentionality (I am, I suppose, the instantiation of any number of computer programs), but as far as we know it is because I am a certain sort of organism with a certain biological (i.e., chemical and physical) structure, and this structure under certain conditions is causally capable of producing perception, action, understanding, learning, and other intentional phenomena. And part of the point of the present argument is that only something that had those causal powers could have that intentionality. Perhaps other physical and chemical processes could produce exactly these effects; perhaps, for example, Martians also have intentionality, but their brains are made of different stuff. That is an empirical question, rather like the question whether photosynthesis can be done by something with a chemistry different from that of chlorophyll.

But the main point of the present argument is that no purely formal model will ever be by itself sufficient for intentionality, because the formal properties are not by themselves constitutive of intentionality, and they have by themselves no causal powers except the power, when instantiated, to produce the next stage of the formalism when the machine is running. And any other causal properties which particular realizations of the formal model have are irrelevant to the formal model because we can always put the same formal model in a different realization where those causal properties are obviously absent. Even if by some miracle Chinese speakers exactly realize Schank's program, we can put the same program in English speakers, water pipes, or computers, none of which understand Chinese, the program notwithstanding.

What matters about brain operation is not the formal shadow cast by the sequence of synapses but rather the actual properties

of the sequences. All the arguments for the strong version of artificial intelligence that I have seen insist on drawing an outline around the shadows cast by cognition and then claiming that the shadows are the real thing.

By way of concluding I want to state some of the general philosophical points implicit in the argument. For clarity I will try to do it in a question-and-answer fashion, and I begin with that old chestnut:

"Could a machine think?"

The answer is, obviously, yes. We are precisely such machines.

"Yes, but could an artifact, a man-made machine, think?"

Assuming it is possible to produce artificially a machine with a nervous system, neurons with axons and dendrites, and all the rest of it, sufficiently like ours, again the answer to the question seems to be obviously 'yes'. If you can exactly duplicate the causes, you could duplicate the effects. And indeed it might be possible to produce consciousness, intentionality and all the rest of it using chemical principles different from those human beings use. It is, as I said, an empirical question.

"OK, but could a digital computer think?"

If by "digital computer" we mean anything at all which has a level of description where it can correctly be described as the instantiation of a computer program, then again the answer is, of course, yes, since we are the instantiations of any number of computer programs and we can think.

"But could something think, understand, etc. *solely* by virtue of being a computer with the right sort of program? Could instantiating a program, the right program of course, by itself be a sufficient condition of understanding?"

This I think is the right question to ask, though it is usually confused with one or more of the earlier questions, and the answer to it is "no."

"Why not?"

Because the formal symbol manipulations by themselves don't have any intentionality: they are meaningless; they aren't even *symbol* manipulations, since the symbols don't symbolize anything. In the linguistic jargon they have only a syntax but no semantics. Such intentionality as computers appear to have is

solely in the minds of those who program them and those who use them, those who send in the input and who interpret the output.

The aim of the Chinese room example was to try to show this by showing that as soon as we put something into the system which really does have intentionality, a man, and we program the man with the formal program, you can see that the formal program carries no additional intentionality. It adds nothing, for example, to a man's ability to understand Chinese.

Precisely that feature of AI which seemed so appealing—the distinction between the program and the realization—proves fatal to the claim that simulation could be duplication. The distinction between the program and its realization in the hardware seems to be parallel to the distinction between the level of mental operations and the level of brain operations. And if we could describe the level of mental operations as a formal program, it seems we could describe what was essential about the mind without doing either introspective psychology or neurophysiology of the brain. But the equation "Mind is to brain as program is to hardware" breaks down at several points, among them the following three:

First, the distinction between program and realization has the consequence that the same program could have all sorts of crazy realizations which had no form of intentionality. Weizenbaum (1976), for example, shows in detail how to construct a computer using a roll of toilet paper and a pile of small stones. Similarly, the Chinese story-understanding program can be programmed into a sequence of water pipes, a set of wind machines, or a monolingual English speaker, none of which thereby acquires an understanding of Chinese. Stones, toilet paper, wind, and water pipes are the wrong kind of stuff to have intentionality in the first place (only something that has the same causal powers as brains can have intentionality), and though the English speaker has the right kind of stuff for intentionality, you can easily see that he doesn't get any extra intentionality by memorizing the program, since memorizing it won't teach him Chinese.

Second, the program is purely formal, but the intentional states are not in that way formal. They are defined in terms of their content, not their form. The belief that it is raining, for example, if defined not as a certain formal shape, but as a certain mental

content, with conditions of satisfaction, a direction of fit (cf. Searle, 1979), etc. Indeed, the belief as such hasn't even got a formal shape in this syntactical sense, since one and the same belief can be given an indefinite number of different syntactical expressions in different linguistic systems.

Third, as I mentioned before, mental states and events are a product of the operation of the brain, but the program is not in that way a product of the computer.

"Well if programs are in no way constitutive of mental processes, then why have so many people believed the converse? That at least needs some explanation."

I don't know the answer to that. The idea that computer simulations could be the real thing ought to have seemed suspicious in the first place because the computer isn't confined to simulating mental operations, by any means. No one supposes that computer simulations of a five-alarm fire will burn the neighborhood down or that a computer simulation of a rainstorm will leave us all drenched. Why on earth would anyone suppose that a computer simulation of understanding actually understood anything? It is sometimes said that it would be frightfully hard to get computers to feel pain or fall in love, but love and pain are neither harder nor easier than cognition or anything else. For simulation, all you need is the right input and output and a program in the middle that transforms the former into the latter. That is all the computer has for anything it does. To confuse simulation with duplication is the same mistake, whether it is pain, love, cognition, fires, or rainstorms.

Still, there are several reasons why AI must have seemed and to many people perhaps still does seem in some way to reproduce and thereby explain mental phenomena, and I believe we will not succeed in removing these illusions until we have fully exposed the reasons that give rise to them.

First, and perhaps most important, is a confusion about the notion of "information processing." Many people in cognitive science believe that the human brain with its mind does something called "information processing," and analogously the computer with its program does information processing, but fires and rainstorms on the other hand don't do information processing at all. Thus though the computer can simulate the formal features of

any process whatever, it stands in a special relation to the mind and brain because when the computer is properly programmed, ideally with the same program as the brain, the information processing is identical in the two cases, and this information processing is really the essence of the mental. But the trouble with this argument is that it rests on an ambiguity in the notion of "information." In the sense in which people "process information" when they reflect, say, on problems in arithmetic or when they read and answer questions about stories, the programmed computer does not do "information processing." Rather, what it does is manipulate formal symbols. The fact that the programmer and the interpreter of the computer output use the symbols to stand for objects in the world is totally beyond the scope of the computer. The computer, to repeat, has a syntax but no semantics. Thus if you type into the computer "2 plus 2 equals?" it will type out "4." But it has no idea that "4" means 4 or that it means anything at all. And the point is not that it lacks some second-order information about the interpretation of its first-order symbols, but rather that its first-order symbols don't have any interpretations as far as the computer is concerned. All the computer has is more symbols. The introduction of the notion of "information processing" therefore produces a dilemma: Either we construe the notion of "information processing" in such a way that it implies intentionality as part of the process or we don't. If the former, then the programmed computer does not do information processing, it only manipulates formal symbols. If the latter, then although the computer does information processing, it is only in the sense in which adding machines, typewriters, stomachs, thermostats, rainstorms, and hurricanes do information processing—namely, they have a level of description where we can describe them as taking information in at one end, transforming it, and producing information as output. But in this case it is up to outside observers to interpret the input and output as information in the ordinary sense. And no similarity is established between the computer and the brain in terms of any similarity of information processing in the two cases.

Secondly, in much of AI there is a residual behaviorism or operationalism. Since appropriately programmed computers can have input/output patterns similar to human beings, we are tempted to postulate mental states in the computer similar to human

mental states. But once we see that it is both conceptually and empirically possible for a system to have human capacities in some realm without having any intentionality at all, we should be able to overcome this impulse. My desk adding machine has calculating capacities but no intentionality, and in this paper I have tried to show that a system could have input and output capabilities which duplicated those of a native Chinese speaker and still not understand Chinese, regardless of how it was programmed. The Turing test is typical of the tradition in being unashamedly behavioristic and operationalistic, and I believe that if AI workers totally repudiated behaviorism and operationalism, much of the confusion between simulation and duplication would be eliminated.

Third, this residual operationalism is joined to a residual form of dualism; indeed, strong AI only makes sense given the dualistic assumption that where the mind is concerned the brain doesn't matter. In strong AI (and in functionalism, as well) what matters are programs, and programs are independent of their realization in machines; indeed, as far as AI is concerned, the same program could be realized by an electronic machine, a Cartesian mental substance, or a Hegelian world spirit. The single most surprising discovery that I have made in discussing these issues is that many AI workers are shocked by my idea that actual human mental phenomena might be dependent on actual physical-chemical properties of actual human brains. But I should not have been surprised; for unless you accept some form of dualism, the strong AI project hasn't got a chance. The project is to reproduce and explain the mental by designing programs; but unless the mind is not only conceptually but empirically independent of the brain, you cannot carry out the project, for the program is completely independent of any realization. Unless you believe that the mind is separable from the brain both conceptually and empirically—dualism in a strong form—you cannot hope to reproduce the mental by writing and running programs since programs must be independent of brains or any other particular forms of instantiation. If mental operations consist of computational operations on formal symbols, it follows that they have no interesting connection with the brain, and the only connection would be that the brain just happens to be one of the indefinitely many types of machines

capable of instantiating the program. This form of dualism is not the traditional Cartesian variety that claims there are two sorts of *substances*, but it is Cartesian in the sense that it insists that what is specifically mental about the mind has no intrinsic connection with the actual properties of the brain. This underlying dualism is masked from us by the fact that AI literature contains frequent fulminations against "dualism"; what the authors seem to be unaware of is that their position presupposes a strong version of dualism.

"Could a machine think?" My own view is that *only* a machine could think, and indeed only very special kinds of machines, namely brains and machines that had the *same causal powers* as brains. And that is the main reason why strong AI has had little to tell us about thinking: it has nothing to tell us about machines. By its own definition it is about programs, and programs are not machines. Whatever else intentionality is, it is a biological phenomenon and it is as likely to be as causally dependent on the specific biochemistry of its origins as lactation, photosynthesis, or any other biological phenomena. No one would suppose that we could produce milk and sugar by running a computer simulation of the formal sequences in lactation and photosynthesis; but where the mind is concerned, many people are willing to believe in such a miracle, because of a deep and abiding dualism: the mind they suppose is a matter of formal processes and is independent of specific material causes in the way that milk and sugar are not.

In defense of this dualism, the hope is often expressed that the brain is a digital computer (early computers, by the way, were often called "electronic brains"). But that is no help. Of course the brain is a digital computer. Since everything is a digital computer, brains are too. The point is that the brain's causal capacity to produce intentionality cannot consist in its instantiating a computer program, since for any program you like it is possible for something to instantiate that program and still not have any mental states. Whatever it is that the brain does to produce intentionality, it cannot consist of instantiating a program, since no program by itself is sufficient for intentionality.

Acknowledgments: I am indebted to a rather large number of people for discussion of these matters and for their patient

attempts to overcome my ignorance of artificial intelligence. I would especially like to thank Ned Block, Hubert Dreyfus, John Haugeland, Roger Schank, Robert Wilensky, and Terry Winograd.

11

Methodological Solipsism Considered as a Research Strategy in Cognitive Psychology

JERRY A. FODOR

... to form the idea of an object and to form an idea simply is the same thing; the reference of the idea to an object being an extraneous denomination, of which in itself it bears no mark or character.

—Hume (1888), p. 20

THE PAPER distinguishes two doctrines, both of which inform theory construction in much of modern cognitive psychology: the representational theory of mind (according to which propositional attitudes are relations that organisms bear to mental representations) and the computational theory of mind (according to which mental processes have access only to formal (nonsemantic) properties of the mental representations over which they are defined.

It is argued that the acceptance of some such formality condition is warranted, at least for that part of psychology which concerns itself with the mental causation of behavior. The paper closes with a discussion of the prospects for a "naturalistic" psychology: one which defines its generalizations over relations between mental representations and their environmental causes. Two related arguments are proposed, both leading to the conclusion that no such research strategy is likely to prove fruitful.

Your standard contemporary cognitive psychologist—your thoroughly modern mentalist—is disposed to reason as follows. To think (e.g.,) that Marvin is melancholy is to represent Marvin in a certain way; viz. as being melancholy (and not, for example,