

## How minds can be computational systems

WILLIAM J. RAPAPORT

*Department of Computer Science, Department of Philosophy and Center for Cognitive Science, State University of New York at Buffalo, Buffalo, NY 14260, USA*

email: rapaport@cs.buffalo.edu

<http://www.cs.buffalo.edu/pub/WWW/faculty/rapaport/>

*Abstract.* The proper treatment of computationalism, as the thesis that cognition is computable, is presented and defended. Some arguments of James H. Fetzer against computationalism are examined and found wanting, and his positive theory of minds as semiotic systems is shown to be consistent with computationalism. An objection is raised to an argument of Selmer Bringsjord against one strand of computationalism, namely, that Turing-Test-passing artifacts are persons, it is argued that, whether or not this objection holds, such artifacts will inevitably be persons.

*Keywords:* cognition, computation, computationalism, incorrigibilism, personhood, semantics, semiotic systems, syntax, Turing machines, Turing test.

### 1. The proper treatment of computationalism

Computationalism is—or ought to be—the thesis that cognition is computable. But what does this mean? What is meant by ‘cognition’ and by ‘computable’?

I take the vague term ‘cognition’ to cover the phenomena that others have called, equally vaguely, ‘mental states and processes’, ‘thinking’, ‘intelligence’, ‘mentality’, or simply ‘the mind’. More specifically, this includes such things as language use, reasoning, conceiving, perceiving, planning, and so on—the topics of such cognitive disciplines as linguistics, cognitive psychology, and the philosophy of mind, among others—in general, of cognitive science. Perhaps this is still vague, but it will be good enough for present purposes.

A slightly more precise story has to be told about ‘computable’. Note, first, that I have said that computationalism is the thesis that cognition is *computable*, not that it is *computation* (as Pylyshyn (1985, p. xiii) characterizes it). There is a subtle but, I think, important difference, as we shall see. The kind of thing that can be computable is a function, i.e. a set of ordered pairs—‘input–output’ pairs, to use computer jargon—such that no two pairs have the same first element but different second elements. Roughly, a *function* is computable if and only if there is an ‘algorithm’ that computes it, i.e. an algorithm that takes as input the first elements of the function’s ordered pairs, manipulates them (in certain constrained ways), and returns the appropriate second elements. To say that it is an *algorithm* that does this is to say that there is an explicitly given, ‘effective procedure’ for converting the input into the output. But what does *this* mean? Following my colleague Stuart C. Shapiro, I tell my introductory computer science students that an algorithm is a

procedure for solving a problem such that: (a) it is unambiguous for the computer or human who will execute it, i.e. all steps of the procedure must be clear and well-defined for the executor, and (b) it is effective, i.e. it must eventually halt and it must be correct. ((a) and (b) are the constraints alluded to above.) However, to be mathematically precise, an algorithm is—following Church’s Thesis—a Turing-machine program (or a recursive function, or a lambda expression, or any one of the other, logically equivalent models of computation, for an excellent historical discussion of the interrelationships of these concepts, see Soare (1996)). We’ll return to this notion in more detail, below.

Thus, to say that cognition is computable is to say that there is an algorithm—more likely, a collection of interrelated algorithms—that computes it. So, what does it mean to say that something ‘computes cognition’? If mental (or psychological) behaviour can be characterized in input–output terms as, perhaps, stimulus–response pairs, then—assuming the set of such pairs is a function (or several functions)—cognition is computable if and only if there is an algorithm (or a collection of algorithms) that computes this function (or functions). Another way to say this is to say that cognition is a recursive function (or a set of recursive functions) (cf. Shapiro 1992, p. 54, 1995, p. 517).

But note that it is quite possible, according to this characterization of computationalism, for cognition itself—for mental states and processes, or brain states and processes—not to *be* a computation, i.e. not to *be* the execution of an algorithm. After all, a computable function need not be given as a computation: it might just be a set of input–output pairs. Arguably, such a set is a trivial computation: a table look-up. But there are other phenomena that are, or may be, computable but not computations. One standard kind of example is illustrated by the solar system, which, arguably, computes Kepler’s laws. However, it could also be said that it is Kepler’s laws that are computable and that describe the behaviour of the solar system, yet the solar system does not compute them, i.e. the behaviour of the solar system is not a computation, even though its behaviour is computable. In theory, a device could convert input into output in some mysterious, or causal but not computational, way, yet be such that the function that is its input–output description is computable. Perhaps the mind works this way: mental states and processes—i.e. cognition—may be *computable* but not *computed*. And similarly for the brain. That is, possibly there are minds—i.e. systems that have the input–output behaviour of cognition—that accomplish it *non-algorithmically* (and maybe human minds are like this.) Personally, I do not believe that this is the case, but it is possible—and it is consistent with computationalism properly treated.

However, if cognition *is* thus *computable*, then any (physical) device that *did* perform cognitive computations would exhibit cognition. It *would* think. And *it* would do so even if *we* (human) cognitive agents did *not* perform cognitive *computations*.

So, is computationalism true? I do not *know*, but I believe so. For one thing, I have not seen any good arguments against it—people such as Dreyfus, Penrose, and Searle notwithstanding. It is not that these people are necessarily wrong. It is just that I do not think that this is the sort of issue that is ready to be refuted by an in-principle argument. It may well be the case that some—maybe even all—aspects of cognition are not computable. But I take the goal of computational cognitive science to be the investigation of the extent to which—and the ways in which—cognition *is* computable. And I take computationalism—understood now as the thesis that *all* cognition is computable—to be its working hypothesis. In fact, much of cognition

is known to be computable: most reasoning is—including first-order, default, and non-monotonic logics, belief revision, as well as searching and game playing. Much of language is. Significant aspects of visual perception, planning, and acting are. What we should do is to see how much of cognition can be computed. The history of computational cognitive science is much too short for us to give up on it yet.

Is computationalism as I have characterized it trivial? Quite the contrary, for it needs to be *demonstrated* that the cognitive functions are computable, and to do this, we need to devise appropriate algorithms. This is by no means a trivial task, as the history of AI has shown. The central question of cognitive science, as I and many other computational cognitive scientists see it, is not (merely): how is human cognition accomplished? Human cognition may, for all we know, not be accomplished computationally. Rather, the central question is the more Kantian one: how is cognition possible? Computationalism properly treated is the hypothesis that cognition can be accomplished computationally, i.e. that it is computable.

## 2. Fetzer's treatment of computationalism

In 'Mental Algorithms', Fetzer (1994) asks whether minds are computational systems, and there, as well as in 'People Are Not Computers' (Fetzer 1998),<sup>1</sup> he answers in the negative. I believe that there are (or will be, in the golden age of cognitive science) some computational systems that are minds, and I believe that mentality (cognition) is *computable*. But, as we have seen, the computationalist should also be willing to believe that it is possible that not all minds are computational systems *if* by that is meant that possibly not all minds *behave* computationally.

Thus, it is misleading for Fetzer to say that 'the idea that human thought *requires* the execution of mental algorithms appears to provide a foundation for research programs in cognitive science' (Fetzer 1994, p. 1, my italics). Rather, what provides the foundation is the idea that thought in general (and human thought in particular) can be explained (or described, or accounted for) in terms of algorithms (call them 'mental algorithms' if you will). As I noted, this is, or can be seen as, an elaboration of behaviourism: *behaviourism* was concerned with describing human thought in stimulus–response terms (i.e. input–output terms) and only those. *Cognitivism* posits processes that mediate the stimulus (input) with the response (output). And *computational* cognitivism posits that those processes are computable (cf. Rapaport 1993).

Let's consider the notion of an algorithm in a bit more detail. It is difficult to talk about what an algorithm is, since the notion is an informal one. What licenses its use is Church's Thesis, the fact that all formal explications of the informal notion have turned out to be logically equivalent, thus giving support to its being more than just an intuitively plausible idea. I offered one informal specification above, due to Shapiro. Fetzer offers two variants: (1) 'algorithms . . . are definite (you always get an answer), reliable (you always get a correct answer), and completable (you always get a correct answer in a finite interval of time)' (Fetzer 1994, p. 4)<sup>2</sup> and (2) 'algorithms are . . . completely reliable . . . procedures . . . that can be carried out in a finite number of steps to solve a problem' (Fetzer 1998, p. 375). In all three cases, the informal notion of algorithm is a relative one: an algorithm is an algorithm *for* a problem or question. Thus, if an alleged algorithm for some problem *P* fails to solve *P* correctly, it fails to be an algorithm *for P*. However, it does not

necessarily thereby fail to be an *algorithm*, for it may be an algorithm for some other problem  $P'$ .

This is where the notion of a 'heuristic' enters. Consider an AI program that uses heuristic game-playing techniques to play chess. It may not always make the 'best' move, but if it makes a move that is 'good enough', that will suffice for the purposes of playing chess. Yet it is an *algorithm* that does this, not the algorithm you thought it was, perhaps, but an algorithm nonetheless. A *heuristic for problem P* can be defined as an *algorithm* for some problem  $P'$ , where the solution to  $P'$  is 'good enough' as a solution to  $P$  (cf. Korf 1992, Shapiro 1992, pp. 54–55, Findler 1993, Herman 1993, and Korfhage 1993.) Thus, to argue, as Fetzer (1994, p. 6) does, that computationalism is false to the extent that it relies on heuristics 'rather' than algorithms is to set up a false dichotomy: the heuristics that AI researchers and computational cognitive scientists use *are* algorithms. In fact, in AI we do not need guaranteed solutions at all, just algorithmic processes that are cognitive. It is best to stick to the Turing-machine analysis of 'algorithm' (i.e. of computation) and omit any reference to the problem for which an algorithm is designed. What is important for computationalism properly treated is whether cognitive processes are algorithmic (i.e. computable) in the Turing-machine sense.

Are there functions that are intuitively algorithmic but are not recursive or Turing-computable? Fetzer cites Cleland (1993), who 'appeals to a class of "mundane functions", which includes recipes for cooking and directions for assembling devices, as examples of effective procedures that are not Turing computable', because they manipulate 'things rather than numerals' (Fetzer 1994, p. 15). However, this takes the standard introduction-to-computer-science analogy for an algorithm and tries to make more of it than is there. The computationalist's point is not that cooking, for example, is algorithmic in the sense that the recipe is an algorithm to be 'followed' (and, incidentally, computers do not 'follow' algorithms, they *execute* them) but that cooking is the *result* of algorithmic processes: I can write a very complex program for a robot who will plan and execute the preparation of dinner in an algorithmic fashion, even using 'pinches' of salt, but the recipe is not an algorithm. A suitably-programmed Turing machine *can* cook (or so the working hypothesis of computational cognitive science would have it).

But is *all* of cognition algorithmic? Fetzer makes two claims that reveal, I think, a serious misunderstanding of the computational cognitive enterprise:

The strongest possible version of the computational conception would therefore appear to incorporate the following claims: that all thinking is reasoning, that all reasoning is reckoning, that all reckoning is computation, and that the boundaries of computability are the boundaries of thought. Thus understood, the thesis is elegant and precise, but it also appears to suffer from at least one fatal flaw: it is (fairly obviously, I think) untrue! The boundaries of thought are vastly broader than those of reasoning, as the exercise of imagination and conjecture demonstrates. Dreams and daydreams are conspicuous examples of non-computational thought processes. (Fetzer 1994, p. 5, cf. 1998, p. 381)

The first claim, that all thinking is reasoning, is a red herring. I agree with Fetzer that 'the boundaries of thought are vastly broader than those of *reasoning*': when I recall a pleasant afternoon spent playing with my son, I am not *reasoning* (there are no premises or conclusions), but I *am* thinking, and—the computationalist maintains—my thinking is computable. In this case, my recollection (my reverie, if you will) is the result of a computable (i.e. algorithmic) mental process.

However, a slight amendment to Fetzer's slippery slope makes me willing to slide down it: although *not* all thinking is *reasoning*, all reasoning *is* reckoning, all

thinking (including reckoning) is computable, and the boundaries of computability are the boundaries of thought.

This last point is also disputed by Fetzer, who says that ‘computability does not define the boundaries of thought. The execution of mental algorithms appears to be no more than one special kind of thinking’ (Fetzer 1994, p. 2). On the working hypothesis that all thinking is computable, computability *does* define the boundaries of thought, and even if mental algorithms are *just* ‘a special kind of thinking’, they are an *important* kind, because systems other than humans can execute them, and these systems can thus be said to think. But the computationalist should, as noted, be willing to recognize the possibility that actual thinkings, even though computable, might not take place by computation.<sup>3</sup> In *that* sense, I can agree with Fetzer.<sup>4</sup> But this is merely to say that some acts of thinking might be computable but not carried out computationally.

The notion of ‘boundaries’, however, is a slippery one. An analogy might clarify this. The distance between locations *A* and *B* might be ‘drivable’—i.e. one can get from *A* to *B* by driving. But it might also be ‘walkable’—i.e. one can get from *A* to *B* by walking (cf. McCarthy 1968 (see reference for reprint version 1985, p. 300)). If any two locations are drivable *but also walkable*, then in one sense drivability does *not* define the boundaries of getting from any *A* to any *B*, because it is *also* walkable. Yet, in another sense, drivability *does* define the boundaries: any two locations *are* drivable.

The second mistaken claim concerns dreaming. To say that dreams are not computational (Fetzer 1998, p. 379) because they themselves do not compute any (interesting) functions, or because they are not heuristics, or because ‘they have no definite starting point and no definite stopping point’ (Fetzer 1998, p. 379) is beside the point. The point is that (or whether) there are computational processes that can *result* in dreams. In fact, this point holds for all mental states and processes: the question is not whether a particular mental state or process is itself an algorithm that computes something, but whether there are algorithms that result in that mental state or process.

Dreams, in any case, are a bad example, since neuroscientists are not really sure what they are or what purposes they serve. My understanding is that they result from possibly random neuron firings that take place when we sleep and that are interpreted *by us as if* they were due to external causes.<sup>5</sup> Suppose for the sake of argument that this is the case. Then, insofar as our ordinary interpretations of neuron firings in non-dreamlike situations are computable, so are dreams.

What about ‘a certain look, a friendly smile, a familiar scent [that] can trigger enormously varied associations of thoughts under what appear to be the same relevant conditions’ (Fetzer 1994, p. 13) or ‘the associative character of ordinary thought’ (as exemplified in the stream-of-consciousness ‘Cornish Game Clams’ example (Fetzer 1998, pp. 385–387))? Fetzer says ‘that these thought processes do not satisfy the computational conception and therefore properly count against it’ (Fetzer 1998, p. 387). Again, however, what computationalism properly treated says is, not that such thoughts *are algorithms*, but that they *can be* the *results* of algorithms. Programs like Racter arguably behave in this associative way, yet are computational, and spreading-activation theories of associationist thinking can account for this behaviour computationally (cf. Quillian 1967). And, of course, ‘what *appear* to be the same relevant conditions’ may in fact be *different* ones.

Instead of algorithms, some writers on computationalism, such as Haugeland (1985), talk about ‘automatic formal systems’—essentially, syntactic, i.e. symbol-

manipulation, systems (cf. Fetzer 1994, pp. 2–3). Fetzer says that ‘what turns a purely formal system into a cognitive system . . . is the existence of an ‘interpretation’ in relation to which the well-formed formulae, axioms, and theorems of that formal system become meaningful and either true or false’ (Fetzer 1994, pp. 2–3). That is *one* way of turning a syntactic system into a cognitive one. But it is important to see that this is an external, third-person attribution of cognition to the system, for it is an external agent that provides the interpretation (cf. Rapaport 1988). This is one aspect of Dennett’s (1971) ‘intentional stance’.

But another way of turning a syntactic system into a cognitive one—as I have argued in ‘Syntactic Semantics’ and ‘Understanding Understanding’ (Rapaport 1988, 1995)—is to ensure that the formal system is sufficiently rich and has some input–output connections with the external world. (The ‘some’ hedge is to allow for cases of ‘handicapped’ humans, cf. Maloney (1987, 1989, Ch. 5) and Shapiro (1995, pp. 521–522).) Such a rich syntactic system need have no interpretation externally imposed on it. Syntax can give rise to semantics of a holistic, conceptual-role variety. Most likely, some of the formal system’s internal symbols (or terms of its language of thought, or nodes of its semantic network) would be internal representations of external entities, causally produced therefrom by perception. And others of its internal symbols (or terms, or nodes) would be concepts of those perceptually-produced symbols. As I argue in ‘Understanding Understanding’ (Rapaport 1995), these would be the system’s internal interpretations of the other symbols. This would be the system’s first-person ‘interpretation’ of its own symbols.

Fetzer would probably not agree with my analysis of syntax and semantics. Unfortunately, he does not provide arguments for claims such as the following:

When . . . marks [that form the basis for the operation of a causal system without having any meaning for the system] are envisioned as *syntax*, . . . they are viewed as the . . . bearers of meaning, which presupposes a point of view. In this sense, syntax *is* relative to an interpretation, interpreter or mind.

It is the potential to sustain an interpretation that qualifies marks as elements of a formal system . . . (Fetzer 1994, p. 14.)

But *why* is syntax thus relative? *Who* ‘views’ the marks as ‘bearers of meaning’? And why do the marks of a formal system need ‘the potential to sustain an interpretation’? The only way I will grant this without argument is if we allow the system itself to provide its own interpretation, by allowing it to map some marks into others, which are ‘understood’ by the system in some primitive way. This process of self-interpretation, however, turns out to be purely syntactic and computable (see Rapaport 1995 for elaboration).

Let me now turn to Fetzer’s positive theory, that minds are not computational systems but *semiotic* systems (Fetzer 1994, p. 17). A semiotic system, according to Fetzer, is something ‘for which something can stand for something else in some respect or other’ (Fetzer 1994, p. 17). Thus, a semiotic system consists of three entities and three relations: the three entities are a sign, a sign user (e.g. a mind), and what the sign stands for. The sign could be either an ‘icon’ (which resembles what it stands for), an ‘index’ (which causes or is an effect of what it stands for), or else a ‘symbol’ (which is conventionally associated with what it stands for). None are ‘merely syntactical marks’, i.e. ‘symbols’ in the computational sense’ (Fetzer 1994, p. 17). As for the three relations, the sign and what it stands for participate in a ‘grounding’ relation, the sign and its user participate in a ‘causal’ relation, and all three participate in an ‘interpretant’ relation.

If minds are semiotic systems, and semiotic systems are not computational, then neither are minds. Fetzer gives two reasons for the second premise: [1] ‘*the same sign* may be variously viewed as an icon . . . , as an index . . . , or as a symbol . . . , and . . . [2] *inductive reasoning* employing heuristics . . . , which are usually reliable but by no means effective procedures, appears to be fundamental to our survival’ (Fetzer 1994, pp. 17–18). But, in the first case, I fail to see what the ambiguity of signs has to do with not being computational. And, in the second case, heuristics *are*, as we have seen, effective procedures. In addition, computational theories of inductive reasoning are a major research area in AI (cf. Angluin and Smith 1992 and Muggleton and Page (in press)). So what are the details of Fetzer’s arguments?

Consider Fetzer’s figure 1 (Fetzer 1994, p. 19), reproduced here as figure 1. It raises more questions than it answers: what is the causation relation between sign-user  $z$  and sign  $S$ ? Which causes which? What is the grounding relation between sign  $S$  and thing  $x$  that  $S$  stands for? The diagram suggests that sign  $S$  is grounded by thing  $x$  that it stands for, which, in turn, suggests that the two binary relations ‘is grounded by’ and ‘stands for’ are the same. But Fetzer says that sign  $S$  stands for thing  $x$  for sign-user  $z$ , which is a 3-place relation that is not shown or easily showable in figure 1. What is shown instead is a binary ‘interpretant’ relation between sign-user  $z$  and thing  $x$ , yet earlier we were told that the interpretant relation was a 3-place one. I offer a slightly clearer diagram in figure 2.

So why are semiotic systems and computational systems disjoint?

. . . the marks that are manipulated by means of programs might be meaningful for the *users* of that system . . . but are not therefore meaningful for use **by** that system itself. (Fetzer 1998, p. 375, my boldface, cf. Fetzer 1998, p. 383.)

But why are they not meaningful for the system? I have argued in ‘Syntactic Semantics’ and ‘Understanding Understanding’ that with enough structure, they *can* be meaningful for the system. For example, consider a program that (a) computes the greatest common divisor of two natural numbers, (b) that has a ‘knowledge base’ of information about arithmetic and about what a greatest common divisor is, and (c) that has natural-language competence (i.e. the ability to interact with the user in natural language, cf. Shapiro and Rapaport (1991)). Such an AI program can be asked what it is doing and how it does it, it can answer that it is computing greatest common divisors, and it can explain what they are and how it computes them, in exactly the same sort of way that a human student who has just learned

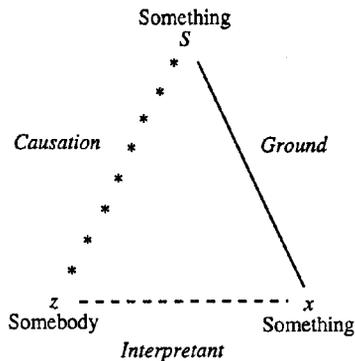


Figure 1. Fetzer’s diagram of a semiotic system (from Fetzer 1998, p. 384).

how to compute them can answer these questions. Not only does the user of such a system ascribe an interpretation *to* it, according to which the *user* says that the system is computing greatest common divisors, but the system *itself* can be said to ‘understand’ what it is doing. It could even turn around and ascribe *to the user* an interpretation of the *user’s* greatest-common-divisor-computing behaviour!

Such a possibility shows that it is incorrect to say, as Fetzer does, that ‘the marks ... are not ... meaningful for use by that system itself’. Fetzer says that this is because the ‘grounding relationship between these marks and that for which they stand’ is absent (Fetzer 1994, p. 18). But why is it absent in, say, my greatest-common-divisor-computing computational agent but not in the human student who computes greatest common divisors? And does Fetzer really think that it is the *grounding* relation that is absent, or rather the relation of sign *S* standing for thing *x* for computer *z* that is absent? It would seem that the point he wants to make is that although sign *S* might stand for thing *x*, it does not do so for the computer *z*, but only for a human user of the computer.

At this point, Fetzer refers to his figure 2 (reproduced here as figure 3), which is supposed to be like figure 1 except that the grounding relation between sign *S* and thing *x* is missing. But there are other differences: sign *S* is now called an ‘input’, thing *x* is called an ‘output’, and sign-user *z* is called a program—presumably, a program for a mental algorithm. But who ever said that ‘the marks by means of which’ digital computers ‘operate’ are only the input? Or that what they stand for are the output? (And, anyway, what does it mean for sign *S* to stand for thing *x* yet not be grounded by *x*?)

Consider Cassie, a computational cognitive agent implemented in the SNePS knowledge-representation and reasoning system by our research group at SUNY Buffalo (Shapiro and Rapaport 1987). As a reader of narratives, Cassie’s input is English sentences (Shapiro and Rapaport 1995). Her output is other English sentences expressing her understanding of the input or answering a question that was input. The marks by means of which Cassie operates are not these sentences, but

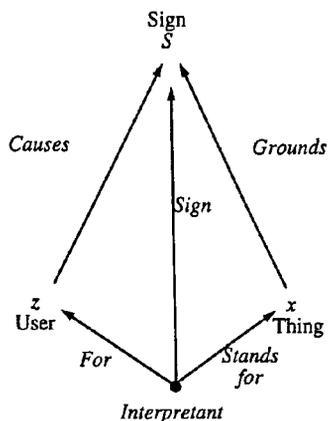


Figure 2. Another diagram of a semiotic system: sign-user *z* causes sign *S*, thing *x* grounds sign *S*, and the 3-place interpretant relation is that sign *S* stands for thing *x* for sign-user *z*.

the nodes of her internal semantic-network ‘mind’. They are not the input (though they may be causally produced by the input).

Do these marks mean anything? Yes: some of them stand for, or are grounded by, whatever in the external world caused them to be built. But Cassie can have no direct access to such external things, so this is not what they mean *for her*. Do they mean anything for Cassie? Again, yes: their meaning for her is their location in Cassie’s entire semantic network. Insofar as some of these marks are directly caused by external objects, and others are concepts of those marks, then those others stand in a grounding relation to the directly caused ones. But all of this is internal to—and meaningful to—Cassie (cf. Ehrlich 1995 and Rapaport 1995). And it is all syntactic symbol manipulation. And it is all computable.

I conclude that computationalism properly treated withstands Fetzer’s objections and that the semiotic approach is consistent with it.

### 3. Bringsjord’s treatment of computationalism

In ‘Computationalism Is Dead, Now What?’ Bringsjord (1998) analyses computationalism into four postulates and a theorem:

**Computationalism** consists of the following four propositions.

**CTT:** A function  $f$  is effectively computable if and only if  $f$  is Turing-computable.

**P=aTM:** Persons are Turing machines.

**TT:** The Turing Test is valid.

**P-BUILD:** Computationalists will succeed in building persons.

**TT-BUILD:** Computationalists will succeed in building Turing Test-passing artifacts. (This proposition is presumably entailed by its predecessor.)

(Bringsjord 1998, p. 395).

As a true believer in computationalism properly treated, I accept all five, with one small caveat about **P=aTM**: I would prefer to say, not that *persons* are Turing machines, but that *minds* are—even better, that cognition is *computable*, in fact, putting it this way links it more closely to Bringsjord’s **CTT**.

Bringsjord, however, denies **P-BUILD**. (Since he accepts **TT-BUILD**, it follows that he must reject **P=aTM**.) He does not, however, offer us a definition of ‘person’. In Chapter IX, ‘Introspection’, of his book, *What Robots Can and Can’t Be*

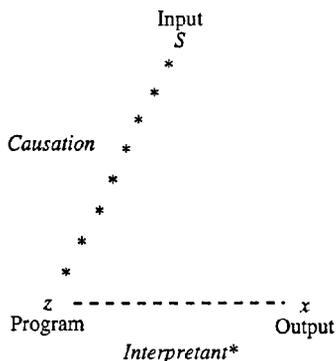


Figure 3. Fetzer’s diagram of a (computational) symbol system (from Fetzer 1998, p. 384).

(Bringsjord 1992), to which he refers us for the defense of his denial of **P-BUILD**, he argues that persons satisfy ‘hyper-weak incorrigibilism’ but that robots cannot. I shall argue that robots *can* satisfy this feature, and I will conclude with some observations on the nature of personhood that suggest a more general way in which **P-BUILD** can be defended.

Hyper-weak incorrigibilism is essentially the following thesis (Bringsjord 1992, pp. 333–335):

- (1) Let  $F$  be a contingent property (i.e. one such that it is possible for something to have it and possible for something (else) to lack it).
- (2) Let  $s$  be a cognitive agent.
- (3) Let  $C'$  be a set of ‘psychological’ or ‘Cartesian’ properties (such as *being sad*, *being in pain*, *being happy*, or *seeming to see a ghost*).
- (4) Let  $C''$  be defined as follows: if  $F' \in C'$ , then *seeming to have*  $F' \in C''$ .
- (5) Suppose  $F \in C''$  (i.e.  $F$  is of the form: *seeming to have*  $F'$ , for some  $F' \in C'$ , thus *seeming to be sad*, *seeming to be in pain*, *seeming to be happy* (and *seeming to seem to see a ghost?*) are all in  $C''$ ).
- (6) Then it is necessary that if  $s$  believes that  $F_s$ , then  $F_s$  (i.e. it is necessary that if you believe that you seem to have  $F'$ , then you seem to have  $F'$ ).

Now, Bringsjord’s argument concerns (1) ‘logicist cognitive engineering’ (as opposed to connectionist cognitive engineering), where cognitive engineering is roughly the interdisciplinary attempt to build a (computational) cognitive agent, and (2) a proposition he calls **AILOGFOUND**, namely:

If [logicist cognitive engineering] is going to succeed, then the robots to be eventually produced by this research effort will be such that

- (i) if there is some significant mental property  $G$  that persons have, these robots must also have  $G$ ,
- (ii) the objects of their ‘beliefs’ (hopes, fears, etc.)—the objects of their *propositional attitudes*—are represented by formulas of some symbol system, and these formulas will be present in these robots’ knowledge bases, and
- (iii) they will be physical instantiations of automata (the physical substrate of which will be something *like* current silicon hardware, but may be something as extravagant as optically-based parallel hardware).

(Bringsjord 1992, pp. 330–331.)

He uses hyper-weak incorrigibilism to refute **AILOGFOUND** as follows: from the claim that hyper-weak incorrigibilism is a significant mental property that persons have (which I will not for the moment deny) and from (i), ‘it follows . . . that the flashy robot (call it ‘ $r$ ’) to be eventually produced by Cognitive Engineering will be able to introspect infallibly with respect to  $C''$  (Bringsjord 1992, p. 341). Therefore, by instantiating hyper-weak incorrigibilism to robot  $r$ , we get:

$$\forall F[(F \text{ is contingent} \wedge F \in C'') \supset \Box(B_r Fr \supset Fr)]; \quad (1)$$

By (ii), this implies:

$$\forall F[(F \text{ is contingent} \wedge F \in C'') \supset \Box(\langle\langle Fr \rangle\rangle \in D(\mathcal{B}) \supset Fr)], \quad (2)$$

where  $\langle\langle Fr \rangle\rangle$  denotes the first-order formula corresponding to the proposition ‘ $Fr$ ’,  $\mathcal{B}$  is a set of first-order formulas that  $r$  ‘believes initially’, and  $D(\mathcal{B}) = \{\alpha \mid \mathcal{B} \vdash \alpha\}$  is the robot’s knowledge base (Bringsjord 1992, p. 340). *The underlying assumption here is that for robot  $r$  to believe a proposition is for a first-order representation of*

the proposition to be an element of robot  $r$ 's knowledge base, or, as it is sometimes called, the robot's 'belief box'.

Next, let  $F^* \in C''$ , Bringsjord suggests taking *seeming to be in pain* as  $F^*$ . Hence,

$$\square[(\ll F^* r \gg \in D(\mathcal{B})) \supset F^* r] \quad (3)$$

(For example, necessarily, if the formula ' $r$  seems to be in pain' is in robot  $r$ 's belief box, then  $r$  seems to be in pain.)

Suppose that  $\neg F^* r$ , for example, that it is *not* the case that the robot seems to be in pain. (Bringsjord does not specify this assumption, but I think that it is needed for what follows.) Now, 'since it's physically possible that the hardware substrate of  $r$  fail, and since, in turn, it's physically possible that this failure be the cause of  $\ll F^* r \gg \in D(\mathcal{B})$ ' (Bringsjord 1992, p. 342), we have:

$$\diamond(\ll F^* r \gg \in D(\mathcal{B}) \wedge \neg F^* r), \quad (4)$$

which contradicts (3). Therefore, **AILOGFOUND** is false. Therefore, logicist cognitive engineering will fail.

There are a number of questions I have about this argument. For one thing, what does it mean for a cognitive agent to *seem* to have a property? It could mean that it seems *to someone else* that the agent has the property, for example, it might seem *to me* that the robot is in pain if it *acts* (as if) in pain. But for *seeming* to have a property to make sense in hyper-weak incorrigibilism, I think it has to mean that it seems *to the agent* that the agent itself has the property.<sup>6</sup> But, then, what does this mean for the robot? Does it mean that there is a 'seeming box' such that if the first-order formula expressing that the robot has the property is in the seeming-box, then it seems to the robot that it has the property? Not only does this make *seeming* to have a property much like *believing* oneself to have it, but I suggest (without argument) that that is just what it is. At any rate, let us suppose so for the sake of argument.<sup>7</sup> Then the brunt of hyper-weak incorrigibilism is this:

$$\square(B_r B_r F' r \supset B_r F' r), \quad (5)$$

where  $F' \in C'$ , e.g. necessarily, if the robot believes that it believes itself to be in pain, then it believes itself to be in pain.

This leads to another problem I have with Bringsjord's argument: the assumption underlying the inference to the belief-box version of hyper-weak incorrigibilism (2) is somewhat simplistic. Sophisticated computational cognitive agents need not equate belief in a proposition  $P$  with the mere presence of a representation of  $P$  in the knowledge base. For example, for reasons independent of the present considerations (see Shapiro and Rapaport 1992 and Rapaport *et al.* 1997), we represent Cassie's believing a proposition in two distinct ways: one way is by 'asserting' the node representing the propositional object of her mental act of believing. This is a 'flag' that *is* roughly equivalent to placing the node in a 'belief box', as in the (simplified) semantic-network representation using the assertion flag of Cassie's belief that John is rich, as shown in figure 4. The other way is to represent explicitly that Cassie herself believes the proposition in question, as in the (simplified) semantic-network representation using an explicit 'I (Cassie) believe that' operator, as shown in figure 5. One reason for allowing these two different representations is that we want to be able to represent that Cassie *believes* that she believes something even if she believes that it is *not* the case. This can be done roughly as shown in figure 6. Here, Cassie

believes that she herself believes  $M1$ , but she does not:  $M2$  represents her believing that she herself believes  $M1$ ,  $M3$  represents her belief that  $\neg M1$ .

But now, given our interpretation of seeming to be in pain (say) as believing oneself to be in pain, what happens on the simplistic belief-box theory? The consequent of (5) becomes:

$$\langle\langle Fr \rangle\rangle \in D(\mathcal{B}), \quad (6)$$

i.e. a formula representing the proposition that  $r$  is in pain is in  $r$ 's belief box. But what is the antecedent of (5)? The question is: How does one represent nested beliefs in a belief-box theory? Bringsjord says that we need an epistemic logic and that  $\langle\langle B_r Fr \rangle\rangle \in D(\mathcal{B})$ —but then how does  $\langle\langle B_r Fr \rangle\rangle$  being in  $D(\mathcal{B})$  relate to  $\langle\langle Fr \rangle\rangle$  being in  $D(\mathcal{B})$ ? Bringsjord does not say, so let me speculate that, in fact, nested beliefs

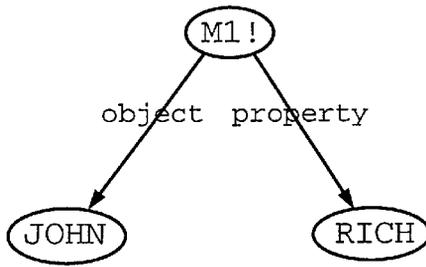


Figure 4.  $M1$  is a simplified SNePS representation of the proposition that John is rich, using an assertion flag (!) to indicate that Cassie believes it.

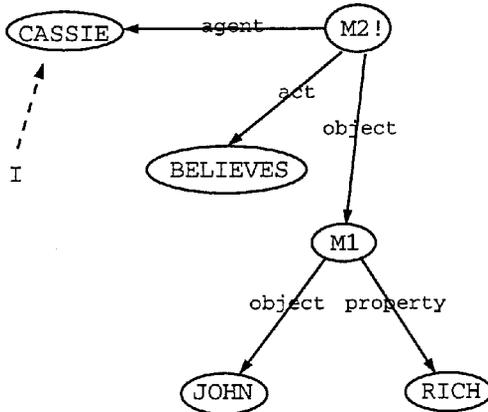


Figure 5.  $M2$  is a simplified SNePS representation of the proposition that Cassie believes  $M1$ . The 'I-pointer' to the node that represents Cassie indicates that the node labeled 'Cassie' is her self-concept. She would express  $M2$  as 'I believe that John is rich'. Since  $M2$  is asserted, Cassie believes that she (herself) believes that John is rich. If  $M1$  were asserted, then Cassie would believe that John is rich. But  $M1$  is not asserted, so it is not the case that Cassie believes that John is rich (even though she believes that she believes it) (see Rapaport *et al.* (1997) for more on the 'I-pointer').

collapse:  $B_r B_r Fr$  just becomes  $\langle Fr \rangle \in D(\mathcal{B})$ . But then (5) is a tautology. Hence, any robot satisfies hyper-weak incorrigibilism.

Let me trade on the kind words of Bringsjord in his footnote 3, in which he praises our work on Cassie, and suggest how Cassie would fare in this case. First, hyper-weak incorrigibilism would have to be part of her ‘unconscious’ processing mechanism, much as her reasoning ability is—i.e. hyper-weak incorrigibilism would not be encoded explicitly as a ‘conscious’ belief but would govern the way she thinks. Second, what it would mean is this: whenever Cassie believed (via the assertion flag) that she believed (represented via the ‘I (Cassie) believe that’ operator) that she is in pain (say), then she would believe (via the assertion flag) that she *was* in pain—i.e. asserting ‘Cassie believes that she herself is in pain’ would automatically assert that Cassie was in pain. In figure 7, hyper-weak incorrigibilism would mean that if M12 were asserted, then M10 would automatically also be asserted. Now suppose that Cassie believes (via assertion) that she is *not* in pain (M11 in figure 7). And suppose that a hardware failure asserts a Cassie-believes-that belief that she herself *is* in pain (M12 in figure 7). Hyper-weak incorrigibilism would then assert that she is in pain (i.e. it would cause an assertion operator to be applied to M10 in figure 7), contradicting her explicit (asserted) belief that she is not in pain, i.e. M10 and M11, both asserted, are inconsistent.

What actually happens with Cassie? *At this point, the SNePS belief-revision system (SNeBR (Martins and Shapiro 1988, Martins and Cravo 1991, Cravo and Martins 1993, Ehrlich 1995)) is invoked, alerting Cassie to the fact that her beliefs are inconsistent, and she would have to give one of them up.* Whichever one she gives up will maintain hyper-weak incorrigibilism, for either she will correct the hardware failure and unassert M12, the Cassie-believes-that belief that she is in pain, or she will decide that she *is* in pain after all.

Thus, I do not think that Bringsjord has made his case that **ALOGFOUND** and therefore logicist cognitive engineering fail, and hence that hyper-weak incorrigibilism rules out Turing-Test-passing artifacts from personhood. Moreover, why could

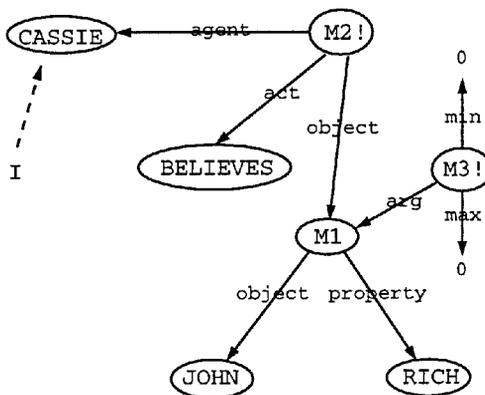


Figure 6. M3 represents the proposition that  $\neg M1$ . The min-max-arg case frame says that at least 0 and at most 0 of the propositions pointed to by the arg arc are true. M3 is asserted, so Cassie believes  $\neg M1$ . Since M2 is also asserted, Cassie believes that she believes M1, but she also believes that  $\neg M1$ . (See Shapiro and Rapaport 1987 and Martins and Shapiro 1988 for details.)

not an analogous neural hardware (i.e. brain) failure cause a similar problem for us? Maybe hyper-weak incorrigibilism is not a necessary feature of persons after all. (Indeed, if hyper-weak incorrigibilism were implemented in Cassie's mind as an explicit 'rule' that she believed, then, when she realized that her beliefs were inconsistent, she could simply reject *it* instead of either M12 or M10!)

In any case, I reject clause (ii) of **AI-LOGFOUND** (the clause that posits belief boxes). Replacing it with a more sophisticated belief-representation mechanism can, I believe, save computationalism. At the very least, it shows that **P-BUILD** has not yet been refuted.

But suppose Bringsjord can strengthen his argument to meet my objections. I still think that **P-BUILD** is true, for an entirely different—and more general—reason. Suppose that **TT-BUILD** is true (as both Bringsjord and I hold). I maintain that such Turing-Test-passing artifacts will be such that we will, *in fact* treat them morally as if they were persons—i.e. they will have the moral and perhaps legal status of persons (cf. Asimov (1976) and Rapaport (1988)). That is, our concept of *person* will be broadened to include not only human persons, but non-human computational ones. In much the same way that our concept of *mind* can be broadened to include not only human minds but other animal and computational ones, thus making *Mind* something like an abstract data type implemented in humans, other animals, and computers, so we will come to see an 'abstract data type' *Person* as implemented in both humans and robots. (Legally, corporations already implement it, cf. Willick (1985) and Rapaport (1988, Section 4.2).) I maintain, that is, that it would be morally wrong to harm an entity that, in virtue of passing the Turing Test, we accept as being intelligent, *even if philosophers like Fetzer, Bringsjord, and Searle are right about propositions like P-BUILD*. We already consider it morally wrong to harm

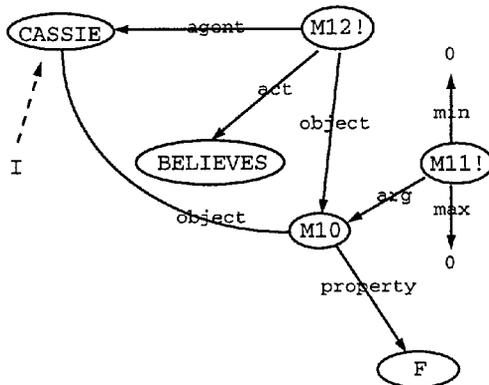


Figure 7. M10 represents that Cassie is F, M12 represents that Cassie believes that she herself is F, M11 represents  $\neg$ M10. By hyper-weak incorrigibilism, if M12 is asserted—i.e. if Cassie believes M12, i.e. if Cassie believes that she herself believes that she herself is F, i.e. if Cassie believes that she herself seems to be F—then M10 must be asserted—i.e. Cassie believes that she herself is F, i.e. Cassie seems to be F. If M11 is asserted, then Cassie believes  $\neg$ M10, i.e. Cassie believes that she herself is not F, i.e. Cassie does not seem to be F. If both M10 and M11 are asserted, then Cassie has inconsistent beliefs, and SNeBR will be invoked (see text).

non-human animals, and the more intelligent (the more human-like?) the animal, the more it seems to be morally wrong. Surely we would owe intelligent robots no less.<sup>8</sup>

## Notes

1. An interesting title given that, at the time that Turing wrote his original paper on what are now called 'Turing machines', a 'computer' was a *human* who did computations for a living!
2. Shapiro, in conversation, has pointed out to me that interactive algorithms, such as those used by our Cassie system for natural-language interaction with a (human) user (Shapiro and Rapaport 1987, 1995, Rapaport 1988, Shapiro 1989), are not reliable or completable, yet they are clearly algorithmic.
3. Although, on Occam's-razor-like grounds, any such non-computational thinking might be so marginal as to be eliminable.
4. I might even be able to agree with Penrose: Penrose's arguments, as I understand them, are of the form: *if* mental phenomena are quantum processes, then they are not algorithmic (Penrose 1989, cf. Fetzer 1994, p. 10). It is important to remember, however, that the antecedent has not been established, not even by Penrose. Even so, if mental phenomena are *computable*, then even if they are not *computed*, perhaps because they are quantum phenomena, computationalism wins.
5. Note the need to posit the brain's (or the mind's) ability to be partitioned into syntax-like neuron firings that are interpreted by semantic-like neuron firings. But it is all just neuron firings, i.e. syntax (see below). Also note the methodologically solipsistic nature of this theory.
6. This use of 'the agent itself' is a 'quasi-indicator', see Castañeda (1966), Rapaport (1986) and Rapaport *et al.* (1997).
7. In discussion, Bringsjord has offered the following counterexample to my identification of seeming to have a property with believing oneself to have it: in the case of optical illusions in which appearance differs from reality, such as the Muller-Lyer illusion, in which two lines of equal length appear to have different lengths, it seems to me that the lines have different lengths even while it is the case that I sincerely believe them to have the same length. I agree that this shows that seeming to have a property can differ from believing oneself to have it. However, arguably this is not a case of what Bringsjord calls 'Cartesian' properties. For Cartesian properties, I think that seeming to have a property *is* believing oneself to have it.
8. An ancestor of this paper was presented as part of an invited symposium, 'Are Minds Computational Systems?', at the 88th Annual Meeting of the Southern Society for Philosophy and Psychology, Nashville, 5 April 1996, with other papers by James H. Fetzer and Selmer Bringsjord. I am grateful to Fetzer for organizing the symposium and to Bringsjord, Fetzer, and members of the SNePS Research Group for comments on earlier versions.

## References

- Angluin, D. and Smith, C. H., 1992, Inductive inference. *Encyclopedia of Artificial Intelligence*, 2nd edition, edited by S. C. Shapiro (New York: John Wiley and Sons), pp. 672–682.
- Asimov, I., 1976, The bicentennial man. *The Bicentennial Man and Other Stories* (Garden City, NY: Doubleday), pp. 135–172.

- Bringsjord, S., 1992, *What Robots Can and Can't Be* (Dordrecht, The Netherlands: Kluwer Academic Publishers).
- Bringsjord, S., 1998, Computationalism is dead, now what? *Journal of Experimental & Theoretical Artificial Intelligence*, this issue, pp. 393–402.
- Castaneda, H.-N., 1966, 'He': a study in the logic of self-consciousness, *Ratio*, **8**: 130–157.
- Cleland, C. E., 1993, Is the Church–Turing thesis true?, *Minds and Machines*, **3**: 283–312.
- Cravo, M. R., and Martins, J. P., 1993, SNePSwD: a newcomer to the SNePS family. *Journal of Experimental and Theoretical Artificial Intelligence*, **5**: 135–148.
- Dennett, D. C., 1971, Intentional systems. *Journal of Philosophy*, **68**: 87–106, reprinted in Dennett, D. C., *Brainstorms* (Montgomery, VT: Bradford Books), pp. 3–22.
- Ehrlich, K., 1995, Automatic vocabulary expansion through narrative context. *Technical Report 95-09* (Buffalo: SUNY Buffalo Department of Computer Science).
- Fetzer, J. H., 1994, Mental algorithms: are minds computational systems? *Pragmatics and Cognition*, **2**: 1–29.
- Fetzer, J. H., 1998, People are not computers: (most) thought processes are not computational procedures. *Journal of Experimental & Theoretical Artificial Intelligence*, this issue, pp. 371–391.
- Findler, N. V., 1993, Heuristic. *Encyclopedia of Computer Science*, 3rd edition, edited by A. Ralston and E. D. Reilly (New York: Van Nostrand Reinhold), pp. 611–612.
- Haugeland, J., 1985, *Artificial Intelligence: The Very Idea* (Cambridge, MA: MIT Press).
- Herman, G. T., 1993, Algorithms, Theory of. *Encyclopedia of Computer Science*, 3rd edition, edited by A. Ralston and E. D. Reilly (New York: Van Nostrand Reinhold), pp. 37–39.
- Korf, R. E., 1992, Heuristics. *Encyclopedia of Artificial Intelligence*, 2nd edition, edited by S. C. Shapiro (New York: John Wiley and Sons), pp. 611–615.
- Korfhage, R. R., 1993, Algorithm. *Encyclopedia of Computer Science*, 3rd edition, edited by A. Ralston and E. D. Reilly (New York: Van Nostrand Reinhold), pp. 27–29.
- Maloney, J. C., 1987, The right stuff: the mundane matter of mind. *Synthese*, **70**: 349–372.
- Maloney, J. C., 1989, *The Mundane Matter of the Mental Language* (Cambridge, UK: Cambridge University Press).
- Martins, J. P. and Cravo M. R., 1991, How to change your mind. *Notas*, **25**: 537–551.
- Martins, J. and Shapiro, S. C., 1988, A model for belief revision. *Artificial Intelligence*, **35**: 25–79.
- McCarthy, J., 1968, Programs with common sense. *Semantic Information Processing*, edited by M. Minsky (Cambridge, MA: MIT Press), pp. 403–418, reprinted in Brachman, R. J. and Levesque, H. J. (eds), 1985, *Readings in Knowledge Representation* (Los Altos, CA: Morgan Kaufmann), pp. 299–307.
- Muggleton, S. and Page, D., 1998, Special issue on inductive logic programming, *Journal of Logic Programming* (in press).
- Penrose, R., 1989, *The Emperor's New Mind: Concerning Computers, Minds, and the Laws of Physics* (Oxford: Oxford University Press).
- Pylshyn, Z., 1985, *Computation and Cognition: Toward a Foundation for Cognitive Science*, 2nd edition (Cambridge, MA: MIT Press).
- Quillian, M. R., 1967, Word concepts: a theory and simulation of some basic semantic capabilities. *Behavioral Science*, **12**: pp. 410–430, reprinted in Brachman, R. J. and Levesque, H. J. (eds), 1985, *Readings in Knowledge Representation* (Los Altos, CA: Morgan Kaufmann), pp. 97–118.
- Rapaport, W. J., 1986, Logical foundations for belief representation. *Cognitive Science*, **10**: 371–422.
- Rapaport, W. J., 1988, Syntactic semantics: foundations of computational natural-language understanding. *Aspects of Artificial Intelligence*, edited by J. H. Fetzer (Dordrecht, Holland: Kluwer Academic Publishers), pp. 81–131, reprinted in Dietrich, E. (ed.), 1994, *Thinking Computers and Virtual Persons: Essays on the Intentionality of Machines* (San Diego: Academic Press), pp. 225–273.
- Rapaport, W. J., 1993, Cognitive science. *Encyclopedia of Computer Science*, 3rd edition, edited by A. Ralston and E. D. Reilly (New York: Van Nostrand Reinhold), pp. 185–189.
- Rapaport, W. J., 1995, Understanding understanding: syntactic semantics and computational cognition. *Philosophical Perspectives*, Vol. 9, *AI, Connectionism, and Philosophical Psychology*, edited by J. E. Tomberlin (Atascadero, CA: Ridgeview), pp. 49–88.
- Rapaport, W. J., Shapiro, S. C. and Wiebe, J. M., 1997, Quasi-indexicals and knowledge reports. *Cognitive Science* **21**: 63–107.
- Shapiro, S. C., 1989, The CASSIE projects: an approach to natural language competence. *Proceedings of the 4th Portuguese Conference on Artificial Intelligence*, Lisbon (Berlin: Springer-Verlag), pp. 362–380.
- Shapiro, S. C., 1992, Artificial Intelligence. *Encyclopedia of Artificial Intelligence*, 2nd edition, edited by S. C. Shapiro (New York: John Wiley and Sons), pp. 54–57, revised version appears in Ralston, A., and Reilly, E. D. (eds), *Encyclopedia of Computer Science*, 3rd edition (New York: Van Nostrand Reinhold), pp. 87–90.
- Shapiro, S. C., 1995, Computationalism. *Minds and Machines*, **5**: 517–524.
- Shapiro, S. C. and Rapaport, W. J., 1987, SNePS considered as a fully intensional propositional semantic network. *The Knowledge Frontier: Essays in the Representation of Knowledge*, edited by N. Cercone and G. McCalla (New York: Springer-Verlag), pp. 262–315.

- Shapiro, S. C. and Rapaport, W. J., 1991, Models and minds: knowledge representation for natural-language competence. *Philosophy and AI: Essays at the Interface*, edited by R. Cummins and J. Pollock (Cambridge, MA: MIT Press), pp. 215–259.
- Shapiro, S. C. and Rapaport, W. J., 1992, The SNePS family. *Computers and Mathematics with Applications*, **23**: 243–275, reprinted in Lehmann, F. (ed.), 1992, *Semantic Networks in Artificial Intelligence* (Oxford: Pergamon Press), pp. 243–275.
- Shapiro, S. C. and Rapaport, W. J., 1995, An introduction to a computational reader of narrative. *Deixis in Narrative: A Cognitive Science Perspective*, edited by J. F. Duchan, G. A. Bruder, and L. E. Hewitt (Hillsdale, NJ: Lawrence Erlbaum Associates), pp. 79–105.
- Soare, R. I., 1996, Computability and recursion. *Bulletin of Symbolic Logic* **2**: 284–321.
- Willick, M. S., 1985, Constitutional law and artificial intelligence: the potential legal recognition of computers as ‘persons’. *Proceedings of the of the 9th International Joint Conference on Artificial Intelligence (IJCAI-85)*, Los Angeles (Los Altos, CA: Morgan Kaufmann), pp. 1271–1273.