

# THE LIMITS OF CORRECTNESS

Brian Cantwell Smith

October, 1960

## 1 Introduction

On October 5, 1960<sup>1</sup>, the American Ballistic Missile Early–Warning System station at Thule, Greenland, indicated a large contingent of Soviet missiles headed towards the United States<sup>2</sup>. Fortunately, common sense prevailed at the informal threat–assessment conference that was immediately convened: international tensions weren’t particularly high at the time. The system had only recently been installed. Krushchev was in New York, and all in all a massive Soviet attack seemed very unlikely. As a result no devastating counter–attack was launched. What was the problem? The moon had risen, and was reflecting radar signals back to earth. Needless to say, this lunar reflection hadn’t been predicted by the system’s designers.

Over the last ten years, the Defense Department has spent many millions of dollars on a new computer technology called “program verification” a branch of computer science whose business, in its own terms, is to “prove programs correct”. Program verification has been studied in theoretical computer science departments since a few seminal papers in the 1960s<sup>3</sup>, but it has only recently started to gain in public visibility, and to be applied to real world problems. General Electric, to consider just one example, has initiated verification projects in their own laboratories: they would like to prove that the programs used in their latest computer–controlled washing machines won’t have any “bugs” (even one serious one can destroy their profit margin)<sup>4</sup>. Although it used to be that only the simplest programs could be “proven correct” — programs to put simple lists into order, to compute simple arithmetic functions — slow but steady progress has been made in extending the range of verification techniques. Recent papers have reported correctness proofs for somewhat more complex programs, including small operating systems, compilers, and other material of modern system design.<sup>5</sup>

What, we do well to ask, does this new technology mean? How good are we at it? For example, if the 1960 warning system had been proven correct (which it was not), could we have avoided the problem with the moon? If it were possible to prove that the programs being written to control automatic launch–on–warning systems were correct, would that mean there could not be a catastrophic accident? In systems now being proposed computers will make launching decisions in a matter of seconds, with no time for any human intervention (let alone for musings about

---

<sup>1</sup>Prepared for the Symposium on Unintentional Nuclear War, Fifth Congress of the International Physicians for the Prevention of Nuclear War, Budapest, Hungary, June 28 – July 1, 1985. Copyright (c) 1985 by Brian Cantwell Smith, Xerox Corporation: Intelligent Systems Laboratory, Palo Alto Research Center, 3333 Coyote Hill Road, Palo Alto, California 94304 USA ; Stanford University: Department of Philosophy, and Center for the Study of Language and Information, Ventura Hall, Stanford, California 94305 USA ; and Computer Professionals for Social Responsibility: P. O. Box 717, Palo Alto, California 94301 USA.

<sup>2</sup>Edmund Berkeley, *The Computer Revolution*, Doubleday, 1972, pp. 175–177, citing newspaper stories in the *Manchester Guardian Weekly* of Dec. 1, 1960, a UPI dispatch published in the *Boston Traveller* of Dec. 13, 1960, and an AP dispatch published in the *New York Times* on Dec. 23, 1960.

<sup>3</sup>McCarthy, John, “A Basis for a Mathematical Theory of Computation”, 1963, in P. Braffort and D. Hirschbert, eds., *Computer Programming and Formal Systems*, Amsterdam: North–Holland, 1967, pp. 33–70. Floyd, Robert, “Assigning Meaning to Programs”, *Proceedings of Symposia in Applied Mathematics* 19, 1967 (also in F. T. Schwartz, ed, *Mathematical Aspects of Computer Science*, Providence: American Mathematical Society, 1967). Naur, P., “Proof of Algorithms by General Snapshots”, *BIT* Vol. 6 No. 4, pp. 310–316, 1966.

<sup>4</sup>Al Stevens, BBN Inc., personal communication.

<sup>5</sup>See for example R. S. Boyer, and Moore, J. S., eds., *The Correctness Problem in Computer Science*, London: Academic Press, 1981.

Kruschev's being in New York). Do the techniques of program verification hold enough promise so that, if these new systems could all be proven correct, we could all sleep more easily at night? These are the questions I want to look at today. And my answer, to give away the punch-line, is no. For fundamental reasons – reasons that anyone can understand – there are inherent limitations to what can be proven about computers and computer programs. Although program verification is an important new technology, useful, like so many other things, in its particular time and place, it should definitely not be called verification. Just because a program is “proven correct”, in other words, you cannot be sure that it will do what you intend. First some background.

## 1.1 General Issues in Program Verification.

Computation is by now the most important enabling technology of nuclear weapons systems: it underlies virtually every aspect of the defense system, from the early warning systems, battle management and simulation systems, and systems for communication and control, to the intricate guidance systems that direct the missiles to their targets. It is difficult, in assessing the chances of an accidental nuclear war, to imagine a more important question to ask than whether these pervasive computer systems will or do work correctly. Because the subject is so large, however, I want to focus on just one aspect of computers relevant to their correctness: the use of models in the construction, use, and analysis of computer systems. I have chosen to look at modelling because I think it exerts the most profound and, in the end, most important influence on the systems we build. But it is only one of an enormous number of important questions. First, therefore, – in order to unsettle you a little – let me just hint at some of the equally important issues I will not address:

1. Complexity. At the current state of the art, only very simple programs can be proven correct. Although it is terribly misleading to assume that either the complexity or power of a computer program is a linear function of length, some rough numbers are illustrative. The simplest possible arithmetic programs are measured in tens of lines; the current state of the verification art extends only to programs of up to several hundred. It is estimated that the systems proposed in the Strategic Defense Initiative (Star Wars), in contrast, will require at least 10,000,000 lines of code<sup>6</sup>. By analogy, compare the difference between resolving a two-person dispute and settling the political problems of the Middle East. There's no a priori reason to believe that strategies successful at one level will scale to the other.
2. Human interaction. Not much can be “proven”, let alone specified formally, about actual human behaviour. The sorts of programs that have so far been proven correct, therefore, do not include much substantial human interaction. On the other hand, as the moon-rise example indicates, it is often crucial to allow enough human intervention to enable people to override system mistakes. System designers, therefore, are faced with a very real dilemma: should they rule out substantive human intervention, in order to develop more confidence in how their systems will perform, or should they include it, so that costly errors can be avoided or at least repaired? The ThreeMile Island incident is a trenchant example of just how serious this trade-off can get: the system design provided for considerable human intervention, but then the operators failed to act “appropriately”. Which strategy leads to the more important kind of correctness?

A standard way out of this dilemma is to specify the behaviour of the system relative to the actions of its operators. But this, as we will see below, pressures the designers to specify

---

<sup>6</sup>Fletcher, James, study chairman, and McMillan, Brockway, panel chairman, Report of the Study on Eliminating the Threat Posed by Nuclear Ballistic Missiles (U), Vol. 5, Battle Management, Communications, and Data Processing (U), U. S. Department of Defense, February 1984.

the system totally in terms of internal actions, not external effects. So you end up proving only that the system will behave in the way that it will behave (i.e., it will raise this line level 3 volts), not do what you want it to do (i.e., launch a missile only if the attack is real). Unfortunately, the latter is clearly what is important. Systems comprising computers and people must function properly as integrated systems; nothing is gained by showing that one cog in a misshapen wheel is a very nice cog indeed.

Furthermore, large computer systems are dynamic, constantly changing, embedded in complex social settings. Another famous “mistake” in the American defense system happened when a human operator mistakenly mounted a training tape, containing a simulation of a full-scale Soviet attack onto a computer that, just by chance, was automatically pulled into service when the primary machine ran into a problem. For some tense moments the simulation data were taken to be the real thing.<sup>7</sup> What does it mean to install a “correct” module into a complex social flux?

3. Levels of failure. Complex computer systems must work at many different levels. It follows that they can fail at many different levels too. By analogy, consider the many different ways a hospital could fail. First, the beams used to frame it might collapse. Or they might perform flawlessly, but the operating room door might be too small to let in a hospital bed (in which case you would blame the architects, not the lumber or steel company.) Or the operating room might be fine, but the hospital might be located in the middle of the woods, where no one could get to it (in which case you would blame the planners). Or, to take a different example, consider how a letter could fail. It might be so torn or soiled that it could not be read. Or it might look beautiful, but be full of spelling mistakes. Or it might have perfect grammar, but disastrous contents.

Computer systems are the same: they can be “correct” at one level – say, in terms of hardware – but fail at another (i. e. , the systems built on top of the hardware can do the wrong thing even if the chips are fine). Sometimes, when people talk about computers failing, they seem to think only the hardware needs to work. And hardware does from time to time fail, causing the machine to come to a halt, or yielding errant behaviour (as for example when a faulty chip in another American early warning system sputtered random digits’ into a signal of how many Soviet missiles had been sighted, again’ causing a false alert<sup>7</sup> ). And the connections between the computers and the world can break: ’ when the moon-rise problem was first recognized, an attempt to override it failed because an iceberg had acidentally cut an undersea telephone cable<sup>8</sup>. But the more important point is that, in order to be reliable, a system has’ to be correct at every relevant level: the hardware is just the starting place (and by far the easiest, at that). Unfortunately, however, we don’t even know what all the relevant levels are. So-called “fault-tolerant” computers, for example, are particularly good at coping with hardware failures, but the soft- ware that runs on them is not thereby i proved.<sup>9</sup>

4. Correctness and Intention. What does correct mean, anyway? Suppose the people want peace, and the President thinks that means having a strong defense, and the Defense department thinks that means having nuclear weapons systems, and the weapons designers request control

---

<sup>7</sup>See, for example, the Hart-Goldwater report to the Committee on Armed Services of the U. S. Senate: “Recent False Alerts from the Nation’s Missile Attack Warning System” (Washington, D.C.: U.S. Government Printing Office, Oct. 9, 1980); Physicians for Social Responsibility, Newsletter, “Accidental Nuclear War”, (Winter 1982), P. 1.

<sup>8</sup>Berkeley, op. cit. See also Daniel Ford’s two part article “The Button”, New Yorker, April 1, 1985, p. 43, and April 8, 1985, p. 49, excerpted from Ford, Daniel, The Button, New York: Simon and Schuster, 1985.

<sup>9</sup>In point of fact, developing software for fault-tolerant systems is an extremely tricky business.

systems to monitor radar signals, and the computer companies are asked to respond to six particular kinds of radar pattern, and the engineers are told to build signal amplifiers with certain circuit characteristics, and the technician is told to write a program to respond to the difference between a two-volt and a four-volt signal on a particular incoming wire. If being correct means doing what was intended, whose intent matters? The technician's? Or what, with twenty years of historical detachment, we would say should have been intended?

With a little thought any of you could extend this list yourself. And none of these issues even touch on the intricate technical problems that arise in actually building the mathematical models of software and systems used in the so-called "correctness" proofs. But, as I said, I want to focus on what I take to be the most important issue underlying all of these concerns: the pervasive use of models. Models are ubiquitous not only in computer science but also in human thinking and language; their very familiarity makes them hard to appreciate. So we'll start simply, looking at modelling on its own, and come back to correctness in a moment.

## 1.2 The Permeating Use of Models

When you design and build a computer system, you first formulate a model of the problem you want it to solve, and then construct the computer program in its terms. For example, if you were to design a medical system to administer drug therapy, you would need to model a variety of things: the patient, the drug, the absorption rate, the desired balance between therapy and toxicity, and so on and so forth. The absorption rate might be modelled as a number proportional to the patient's weight, or proportional to body surface area, or as some more complex function of weight, age, and sex.

Similarly, computers that control traffic lights are based on some model of traffic of how long it takes to drive across the intersection, of how much metal cars contain (the signal change mechanisms are triggered by metal detectors buried under each street). Bicyclists, as it happens, often have problems with automatic traffic lights, because bicycles don't exactly fit the model: they don't contain enough iron to trigger the metal-detectors. I also once saw a tractor get into trouble because it couldn't move as fast as the system "thought" it would: the cross-light went green when the tractor was only half-way through the intersection.

To build a model is to conceive of the world in a certain delimited way. To some extent you must build models before building any artifact at all, including televisions and toasters, but computers have a special dependence on these models: you write an explicit description of the model down inside the computer, in the form of a set of rules or what are called representations essentially linguistic formulae encoding, in the terms of the model, the facts and data thought to be relevant to the system's behaviour. It is with respect to these representations that computer systems work. In fact that's really what computers are (and how they differ from other machines): they run by manipulating representations, and representations are always formulated in terms of models. This can all be summarized in a slogan: no computation without representation.

The models, on which the representations are based, come in all shapes and sizes. Balsa models of cars and airplanes, for example, are used to study air friction and lift. Blueprints can be viewed as models of buildings: musical scores as models of a symphony. But models can also be abstract. Mathematical models, in particular, are so widely used that it is hard to think of anything that they haven't been used for: from whole social and economic systems, to personality traits in teen-agers, to genetic structures, to the mass and charge of sub-atomic particles. These models, furthermore, permeate all discussion and communication. Every expression of language can be viewed as resting implicitly on some model of the world.

What is important, for our purposes, is that every model deals with its subject matter at some particular level of abstraction, paying attention to certain details, throwing away others, grouping together similar aspects into common categories, and so forth. So the drug model mentioned above would probably pay attention to the patients' weight, but ignore their tastes in music. Mathematical models of traffic typically ignore the temperments of taxi drivers. Sometimes what is ignored is at too "low" a level; sometimes too "high": it depends on the purposes for which the model is being used. So a hospital blueprint: would pay attention to the structure and connection of its beams, but not to the arrangements. of proteins in the wood the beams are made of, nor to the efficacy of the resulting operating room.

Models *have to* ignore things exactly because they view the world at a level of abstraction ('abstraction' is from the Latin *abstrahere*, 'to pull or draw away'). And it is good that they do: otherwise they would drown in the infinite richness of the embedding world. Though this isn't the place for metaphysics, it would not be too much to say that every act of conceptualization, analysis, categorization, does a certain amount of violence to its subject matter, in order to get at the underlying regularities that group things together. If you don't commit that act of violence – don't ignore some of what's going on you would become so hypersensitive and so overcome with complexity that you would be unable to act.

To capture all this in a word, we will say that models are inherently partial. All thinking, and all computation, are similarly partial. Furthermore – and this is the important point – thinking and computation have to be partial: that's how they are able to work.

### 1.3 Full-blooded Action

Something that is not partial, however, is action. When you reach out your hand and grasp a plow, it is the real field you are digging up, not your model of it. Models, in other words, may be abstract, and thinking may be abstract, and some aspects of computation may be abstract, but action is not. To actually build a hospital, to clench the steering wheel and drive through the intersection, or to inject a drug into a person's body, is to act in the full-blooded world, not in a partial or distilled model of it.

This difference between action and modelling is extraordinarily important. Even if your every thought is formulated in the terms of some model, to act is to take leave of the model and participate in the whole, rich, infinitely variegated world. For this reason, among others, action plays a crucial role, especially in the human case, in grounding the more abstract processes of modelling or conceptualization. One form that grounding can take, which computer systems can already take advantage of, is to provide feedback on how well the modelling is going. For example, if an industrial robot develops an internal three-dimensional representation of a wheel assembly passing by on a conveyor belt, and then guides its arm towards that object and tries to pick it up, it can use video systems or force sensors to see how well the model corresponded to what was actually the case. The world doesn't care about the model: the claws will settle on the wheel just in case the actualities mesh.

Feedback is a special case of a very general phenomenon: you often learn, when you do act, just how good or bad your conceptual model was. You learn, that is, if you have adequate sensory apparatus, the capacity to assess the sensed experience, the inner resources to revise and reconceptualize, and the luxury of recovering from minor mistakes and failures.

## 1.4 Computers and Models

What does all this have to do with computers and with correctness? The point is that computers, like us, participate in the real world: they take real actions. One of the most important facts about computers, to put this another way, is that we plug them in. They are not, as some theoreticians seem to suppose, pure mathematical abstractions, living in a pure detached heaven. They land real planes at real airports; administer real drugs; and – as those of you here today know all too well – control real radars, missiles and command systems. Like us, in other words, although they base their actions on models, they have consequence in a world that inevitably transcends the partiality of those enabling models. Like us, in other words, and unlike the objects of mathematics, they are challenged by the inexorable conflict between the partial but tractable model, and the actual but infinite world.

And, to make the only too obvious point: we in general have no guarantee that the models are right – indeed we have no guarantee about much of anything about the relationship between model and world. As we will see, current notions of “correctness” don’t even address this fundamental question.

In philosophy and logic, as it happens, there is a very precise mathematical theory called “model theory”. You might think that it would be a theory about what models are, what they are good for, how they correspond to the worlds they are models of, and so forth. You might even hope this was true, for the following reason: a great deal of theoretical computer science, and all of the work in program verification and correctness, historically derives from this model–theoretic tradition, and depends on its techniques. Unfortunately, however, model theory doesn’t address the model–world relationship at all. Rather, what model theory does is to tell you how your descriptions, representations, and programs correspond to your model.

The situation, in other words, is roughly as depicted in Figure 1, below. You are to imagine a description, program, computer system (or even a thought – they are all similar in this regard) in the left hand box, and the very real world in the right. Mediating between the two is the inevitable model, serving as an idealized or preconceptualized simulacrum of the world, in terms of which the description or program or whatever can be understood. One way to understand the model is as the glasses through which the program or computer looks at the world: it is the world, that is, as the system sees it (though not, of course, as it necessarily is).

The technical subject of “model theory”, as I have already said, is a study of the relationship on the left. What about the relationship on the right? The answer, and one of the main points I hope you will take away from this discussion, is that, at this point in intellectual history, we have no theory of this right–hand side relationship.

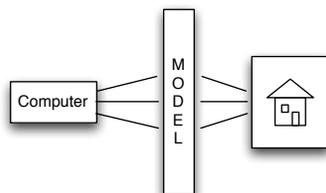


Figure 1: Computers, Models, and the Embedding World

There are lots of reasons for this, some very complex. For one thing, most of our currently accepted formal techniques were developed, during the first half of this century, to deal with mathematics and physics. Mathematics is unique, with respect to models, because (at least to a first level of approximation) its subject matter is the world of models and abstract structures, and therefore

the model–world relationship is relatively unproblematic. The situation in physics is more complex, of course, as is the relationship between mathematics and physics. How apparently pure mathematical structures could be used to model the material substrate of the universe is a question that has exercised physical scientists for centuries. But the point is that, whether or not one believes that the best physical models do more justice and therefore less violence to the world than do models in so-called “higher–level” disciplines like sociology or economics, formal techniques don’t themselves address the question of adequacy.

Another reason we don’t have a theory of the right–hand side is that there is very little agreement on what such a theory would look like. In fact all kinds of questions arise, when one studies the model–world relationship explicitly, about whether it can be treated formally at all, about whether it can be treated rigorously, even if not formally (and what the relationship is between those two), about whether any theory will be more than usually infected with prejudices and preconceptions of the theorist, and so forth. The investigation quickly leads to foundational questions in mathematics, philosophy, and language, as well as computer science. But none of what one learns in any way lessens its ultimate importance. In the end, any adequate theory of action, and, consequently, any adequate theory of correctness, will have to take the model–world relationship into account.

## 2 Correctness and Relative Consistency

Let’s get back, then, to computers, and to correctness. As I mentioned earlier, the word ‘correct’ is already problematic, especially as it relates to underlying intention. Is a program correct when it does what we have instructed it to do? or what we wanted it to do? or what history would dispassionately say it should have done? Analysing what correctness should mean is too complex a topic to take up directly. What I want to do, in the time remaining, is to describe what sorts of correctness we are presently capable of analysing.

In order to understand this, we need to understand one more thing about building computer systems. I have already said, when you design a computer system, that you first develop a model of the world, as indicated in the diagram. But you don’t, in general, ever get to hold the model in your hand: computer systems, in general, are based on models that are purely abstract. Rather, if you are interested in proving your program “correct”, you develop two concrete things, structured in terms of the abstract underlying model (although these are listed here in logical order, the program is very often written first):

1. A specification: a formal description in some standard formal language, specified in terms of the model, in which the desired behaviour is described; and
2. The program: a set of instructions and representations, also formulated in the terms of the model, which the computer uses as the basis for its actions.

How do these two differ? In various ways, of which one is particularly important. The program has to say how the behaviour is to be achieved, typically in a step by step fashion (and often in excruciating detail). The specification, however, is less constrained: all it has to do is to specify what proper behaviour would be, independent of how it is accomplished. For example, a specification for a milk–delivery system might simply be: “Make one milk delivery at each store, driving the shortest possible distance in total.” That’s just a description of what has to happen. The program, on the other hand, would have the much more difficult job of saying how this was to be accomplished. It might be phrased as follows: “drive four blocks north, turn right, stop at Gregory’s Grocery Store on the corner, drop off the milk, then drive 17 blocks north–east, . . .”. Specifications, to use

use some of the jargon of the field, are essentially declarative; they are like indicative sentences or claims. Programs, on the other hand, are procedural: they must contain instructions that lead to a determinate sequence of actions.

What, then, is a *proof of correctness*? It is a proof that ***any system that obeys the program will satisfy the specification.***

There are, as is probably quite evident, two kinds of problems here. The first, often acknowledged, is that the correctness proof is in reality only a proof that two characterizations of something are compatible. When the two differ – i. e. , when you try to prove correctness and fail – there is no more reason to believe that the first (the specification) is any more correct than the second (the program). As a matter of technical practice, specifications tend to be extraordinarily complex formal descriptions, just as subject to bugs and design errors and so forth as programs. In fact they are very much like programs, as this introduction should suggest. So what almost always happens, when you write a specification and a program, and try to show that they are compatible, is that you have to adjust both of them in order to get them to converge.

For example, suppose you write a program to factor a number  $C$ , producing two answers  $A$  and  $B$ . Your specification might be:

Given a number  $C$ , produce numbers  $A$  and  $B$  such that  $A \times B = C$ .

This is a specification, not a program, because it doesn't tell you how to come up with  $A$  and  $B$ . All it tells you is what properties  $A$  and  $B$  should have. In particular, suppose I say: ok,  $C$  is 5,332,114; what are  $A$  and  $B$ ? Staring at the specification just given won't help you to come up with the answer. Suppose, on the other hand, given this specification, that you then write a program – say, by successively trying pairs of numbers until you find two that work. Suppose further that you then set out to prove that your program meets your specifications. And, finally, suppose that this proof can be constructed (I won't go into details here; I hope you can imagine that such a proof could be constructed). With all three things in hand – program, specification, and proof – you might think you were done.

In fact, however, things are rarely that simple, as even this simple example can show. In particular, suppose, after doing all this work, that you try your program out, confident that it must work because you have a proof of its correctness. You randomly give it 14 as an input, expecting 2 and 7. But in fact it gives you back the answers  $A = 1$  and  $B = 14$ . In fact, you realize upon further examination, it will always give back  $A = 1$  and  $B = C$ . It does this, even though you have a proof of its being correct, because you didn't make your specification meet your intentions. You wanted both  $A$  and  $B$  to be different from  $C$  (and also different from 1), but you forgot to say that. In this case you have to modify both the program and the specification. A plausible new version of the latter would be:

Given a number  $C$ , produce  $A$  and  $B$  such that  $A \neq 1$  and  $B \neq 1$  and  $A \times B = C$ .

And so one and so forth: the point, I take it, is obvious. If the next version of the program given 14, produces  $A = -1$  and  $B = -14$ , you would similarly have met your new specification, but still failed to meet your intention. Writing “good” specifications – which is to say, writing specifications that capture your intention – is hard.

It should be apparent, nonetheless, that developing even straightforward proofs of “correctness” is nonetheless very useful. It typically forces you to delineate, very explicitly and completely, the model on which both program and specification are based. A great many of the simple bugs that occur in programs, of which the problem of producing 1 and 14 was an example, arise from sloppiness and unclarity about the model. Such bugs are not identified by the proof, but they are

often unearthed in the attempt to prove. And of course there is nothing wrong with this practice; anything that helps to eradicate errors and increase confidence is to be applauded. The point, rather, is to show exactly what these proofs consist in.

In particular, as the discussion has shown, when you show that a program meets its specifications, all you have done is to show that two formal descriptions, slightly different in character, are compatible. This is why I think it is somewhere between misleading and immoral for computer scientists to call this “correctness”. What is called a proof of correctness is really a proof of the compatibility or consistency between two formal objects of an extremely similar sort: program and specification. As a community, we computer scientists should call this relative consistency, and drop the word ‘correctness’ completely.

What proofs of relative consistency ignore is the second problem intimidated earlier. Nothing in the so-called program verification process per se deals with the right-hand side relationship: the relationship between the model and the world. But, as is clear, it is over inadequacies on the right hand side – inadequacies, that is, in the models in terms of which the programs and specifications are written – that systems so commonly fail.

The problem with the moon-rise, for example, was a problem of this second sort. The difficulty was not that the program failed, in terms of the model. The problem rather, was that the model was overly simplistic; it didn’t correspond to what was the case in –the world. Or, to put it more carefully, since all models fail to correspond to the world in indefinitely many ways, as we have already said, it didn’t correspond to what was the case in a crucial and relevant way. In other words, to answer one of our original questions, even if a formal specification had been written for the 1960 warning system, and a proof of correctness generated, there is no reason to believe that potential difficulties with the moon would have emerged.

You might think that the designers were sloppy; that they would have thought of the moon if they had been more careful. But it turns out to be extremely difficult to develop realistic models of any but the most artificial situations, and to assess how adequate these models are. To see just how hard it can be, think back on the case of General Electric, and imagine writing appliance specifications, this time for a refrigerator. To give the example some force, imagine that you are contracting the refrigerator out to be built by an independent supplier, and that you want to put a specification into the contract that is sufficiently precise to guarantee that you will be happy with anything that the supplier delivers that meets the contract.

Your first version might be quite simple say, that it should maintain an internal temperature of between 3 and 6 degrees Centigrade; not use more than 200 Watts of electricity; cost less than \$100 to manufacture; have an internal volume of half a cubic meter; and so on and so forth. But of course there are hundreds of other properties that you implicitly rely on: it should, presumably, be structurally sound: you wouldn’t be happy with a deliciously cool plastic bag. It shouldn’t weigh more than a ton, or emit loud noises. And it shouldn’t fling projectiles out at high speed when the ‘door’ is opened. In general, it is impossible, when writing specifications, to include everything that you want: legal contracts, and other humanly’ interpretable specifications, are always stated within a background of common sense, to cover the myriad unstated and unstatable assumptions assumed – to hold in force. (Current computer programs, alas, have no common sense, as the cartoonists know so well.)

So it is hard to make sure that everything that meets your specification will really be a refrigerator; it is also hard to make sure that your requirements don’t rule out perfectly good refrigerators. Suppose for example a customer plugs a toaster in, puts it inside the refrigerator and complains that the object he received doesn’t meet the temperature specification and must therefore not be a refrigerator. Or suppose he tries to run it upside down. Or complains that it doesn’t work in outer space, even though you didn’t explicitly specify that it would only work within the earth’s

atmosphere. Or spins it a 10,000 rpm. Or even just unplugs it. In each case you would say that the problem lies not with the refrigerator but with the use. But how is use to be specified? The point is that, as well as modelling the artifact itself, you have to model the relevant part of the world in which it will be embedded. It follows that the model of a refrigerator as a device that always maintains an internal temperature of between 3 and 6 degrees is too “strict” to cover all possible situations. One could try to model what appropriate use would be, though specifications don’t ordinarily, even try to identify all the relevant circumstantial factors. As well as there being a background set of constraints with respect to which a model is formulated, in other words, there is also a background set of assumptions on which a specification is allowed at any point to rely.

### 3 The Limits of Correctness

It’s time to summarize what we’ve said so far. The first challenge to developing a perfectly “correct” computer system stems from the sheer complexity of real-world tasks. We mentioned at the outset various factors that contribute to this complexity: human interaction, unpredictable factors of setting, hardware problems, difficulties in identifying salient levels of abstraction, etc. Nor is this complexity of only theoretical concern. A December 1984 report of the American Defense Science Board Task Force on “Military Applications of New-Generation Computing Technologies” identifies the following gap between current laboratory demonstrations and what will be required for successful military applications – applications they call “Real World; Life or Death”. In their estimation the military now needs (and, so far as one can tell, expects to produce) an increase in the power of computer systems of nine orders of magnitude, accounting for both speed and amount of information to: be processed. That is a 1,000,000,000-fold increase over current research systems, equivalent to the difference between a full century of the entire New York metropolitan area, compared to one day in the life of a hamlet of one hundred people. And remember that even current systems are already several orders of magnitude more complex than those for which we can currently develop proofs of relative consistency.

But sheer complexity has not been our primary subject matter. The second challenge to computational correctness, more serious, comes from the problem of formulating or specifying an appropriate model. Except in the most highly artificial or constrained domains, modelling the embedding situation is an approximate, not a complete, endeavour. It has the best hopes of even partial success in what Winograd has called “systematic domains”: areas where the relevant stock of objects, properties, and relationships are most clearly and regularly pre-defined. Thus bacteremia, or warehouse inventories, or even flight paths of airplanes coming into airports, are relatively systematic domains, at least compared to conflict negotiations, any situations involving intentional human agency, learning and instruction, and so forth. The systems that land airplanes are hybrids – combinations of computers and people – exactly because the unforeseeable happens, and because what happens is in part the result of human action, requiring human interpretation. Although it is impressive how well the phone companies can model telephone connections, lines, and even develop statistical models of telephone use, at a certain level of abstraction, it would nevertheless be impossible to model the content of the telephone conversations themselves.

Third, and finally, is the question of what one does about these first two facts. It is because of the answer to this last question that I have talked, so far, somewhat interchangeably about people and computers. With respect to the ultimate limits of models and conceptualization, both people and computers are restrained by the same truths. If the world is infinitely rich and variegated, no prior conceptualization of it, nor any abstraction, will ever do it full justice. That’s ok – or at least we might as well say that it’s ok, since that’s the world we’ve got. What matters is that we

not forget about that richness – that we not think, with misplaced optimism, that machines might magically have access to a kind of “correctness” to which people cannot even aspire.

It is time, to put this another way, that we change the traditional terms of the debate. The question is not whether machines can do things, as if, in the background, lies the implicit assumption that the object of comparison is people. Plans to build automated systems capable of making a “decision”, in a matter of seconds, to annihilate Europe, say, should make you uneasy; requiring a person to make the same decision in a matter of the same few seconds should make you uneasy too, and for very similar reasons. The problem is that there is simply no way that reasoning of any sort can do justice to the inevitable complexity of the situation, because of what reasoning is. Reasoning is based on partial models. Which means it cannot be guaranteed to be correct. Which means, to suggest just one possible strategy for action, that we might try, in our treaty negotiations, to find mechanisms to slow our weapons systems down.

It is striking to realise, once the comparison between machines and people is raised explicitly, that we don’t typically expect “correctness” for people in anything like the form that we presume it for computers. In fact quite the opposite, and in a revealing way. Imagine, in some by-gone era, sending a soldier off to war, and giving him (it would surely have been a “him”) final instructions. “Obey your commander, help your fellow soldier”, you might say, “and above all do your country honour”. What is striking about this is that it is considered not just a weakness, but a punishable weakness – a breach of morality – to obey instructions blindly (in fact, and for relevant reasons, you generally can’t follow instructions blindly; they have to be interpreted to the situation at hand). You are subject to court-martial, for example, if you violate fundamental moral principles, such as murdering women and children, even if following strict orders.

In the human case, in other words, our social and moral systems seem to have built in an acceptance of the uncertainties and limitations inherent in the model-world relationship. We know that the assumptions and preconceptions built into instructions will sometimes fail, and we know that instructions are always incomplete; we exactly rely on judgment, responsibility, consciousness, and so forth, to carry someone through those situations all situations, in fact – where model and world part company. In fact we never talk about people, in terms of their overall personality, being correct; we talk about people being reliable, a much more substantive term. It is individual actions, fully situated in a particular setting, that are correct or incorrect, not people in general, or systems. What leads to the highest number of correct human actions is a person’s being reliable, experienced, capable of good judgment, etc.

There are two possible morals here, for computers. The first has to do with the notion of experience. In point of fact, program verification is not the only, or even the most common, method of obtaining assurance that a computer system will do the right thing. Programs are usually judged acceptable, and are typically accepted into use, not because we prove them “correct”, but because they have shown themselves relatively reliable in their destined situations for some substantial period of time. And, as part of this experience, we expect them to fail: there always has to be room for failure. Certainly no one would ever accept a program without this in situ testing: a proof of correctness is at best added insurance, not a replacement, for real life experience. Unfortunately, for the ten million lines of code that is supposed to control and coordinate the Star Wars Defense System, there will never, God willing, be an in situ test.

One answer, of course, if genuine testing is impossible, is to run a simulation of the real situation. But simulation, as our diagram should make clear, –tests onZy the Zeft-hand side reZationship. Simulations are defined in terms of models; they don’t test the relationship between the model and world. That is exactly why simulations and tests can never replace embedding a program in the real world. All the war-games we hear about and hypothetical military scenarios, and electronic battlefield simulators, and so forth, are all based on exactly the kinds of models we have been

talking about all along. In fact the subject of simulation, worthy of a whole analysis on its own, is really just our whole subject welling up all over again.

I said earlier that there were two morals to be drawn, for the computer, from the fact that we ask people to be reliable, not correct. The second moral—is for those who, when confronted with the fact that genuine or adequate experience cannot be had, would say “oh, well, let’s build responsibility and morality into the computers – if people can have it, there’s no reason why machines can’t have it too.” Now I will not argue that this is inherently impossible, in a metaphysical or ultimate philosophical sense, but a few short comments are in order. First, from the fact that humans sometimes are responsible, it does not follow that we know what responsibility is: from tacit skills no explicit model is necessarily forthcoming. We simply do not know what aspects of the human condition underlie the modest levels of responsibility to which we sometimes rise. And second, with respect to the goal of building computers with even human levels of full reliability and responsibility, I can state with surety that the present state of artificial intelligence is about as far from this as mosquitos are from flying to the moon.

But there are deeper morals even than these. The point is that even if we could make computers reliable, they still wouldn’t necessarily always do the correct thing. People aren’t provably “correct”, either: that’s why we hope they are responsible, and it is surely one of the major ethical facts is that correctness and responsibility don’t coincide. Even if, in another 1,000 years, someone were to devise a genuinely responsible computer system, there is no reason to suppose that it would achieve “perfect correctness” either, in the sense of never doing anything wrong. This isn’t a failure in the sense of a performance limitation; it stems from the deeper fact that models must abstract, in order to be useful. The lesson to be learned from the violence inherent in the model–world relationship, in other words, is that there is an inherent conflict between the power of analysis and conceptualization, on the one hand, and sensitivity to the infinite richness, on the other.

But perhaps this is an overly abstract way to put it. Perhaps, instead, we should just remember that there will always be another moon–rise.