

Robert Epstein
Gary Roberts
Grace Beber
Editors

Parsing the Turing Test

*Philosophical and Methodological
Issues in the Quest for the
Thinking Computer*

Robert Epstein • Gary Roberts • Grace Beber
Editors

Parsing the Turing Test

Philosophical and Methodological Issues
in the Quest for the Thinking Computer

 Springer

Contents

Foreword	vii
Acknowledgments	ix
Introduction	xi
About the Editors	xxiii
Part I Setting the Stage	
Chapter 1 The Quest for the Thinking Computer	3
Robert Epstein	
Chapter 2 Alan Turing and the Turing Test	13
Andrew Hodges	
Chapter 3 Computing Machinery and Intelligence	23
Alan M. Turing	
Chapter 4 Commentary on Turing’s “Computing Machinery and Intelligence”	67
John Lucas	
Part II The Ongoing Philosophical Debate	
Chapter 5 The Turing Test: Mapping and Navigating the Debate	73
Robert E. Horn	
Chapter 6 If I Were Judge	89
Selmer Bringsjord	
Chapter 7 Turing on the “Imitation Game”	103
Noam Chomsky	

Chapter 8	On the Nature of Intelligence: Turing, Church, von Neumann, and the Brain	107
	Paul M. Churchland	
Chapter 9	Turing's Test: A Philosophical and Historical Guide	119
	Jack Copeland and Diane Proudfoot	
Chapter 10	The Turing Test: 55 Years Later	139
	John R. Searle	
Chapter 11	Doing Justice to the Imitation Game: A Farewell to Formalism	151
	Jean Lassègue	
Part III The New Methodological Debates		
Chapter 12	How to Hold a Turing Test Contest	173
	Hugh Loebner	
Chapter 13	The Anatomy of A.L.I.C.E.	181
	Richard S. Wallace	
Chapter 14	The Social Embedding of Intelligence: Towards Producing a Machine that Could Pass the Turing Test	211
	Bruce Edmonds	
Chapter 15	How My Program Passed the Turing Test	237
	Mark Humphrys	
Chapter 16	Building a Machine Smart Enough to Pass the Turing Test: Could We, Should We, Will We?	261
	Douglas B. Lenat	
Chapter 17	Mind as Space: Toward the Automatic Discovery of a Universal Human Semantic-affective Hyperspace – A Possible Subcognitive Foundation of a Computer Program Able to Pass the Turing Test	283
	Chris McKinstry	
Chapter 18	Can People Think? Or Machines? A Unified Protocol for Turing Testing	301
	Stuart Watt	

Contents	xxi
Chapter 19 The Turing Hub as a Standard for Turing Test Interfaces	319
Robby Garner	
Chapter 20 Conversation Simulation and Sensible Surprises	325
Jason L. Hutchens	
Chapter 21 A Computational Behaviorist Takes Turing's Test	343
Thomas E. Whalen	
Chapter 22 Bringing AI to Life: Putting Today's Tools and Resources to Work	359
Kevin L. Copple	
Chapter 23 Laplace, Turing and the "Imitation Game" Impossible Geometry: Randomness, Determinism and Programs in Turing's Test	377
Giuseppe Longo	
Chapter 24 Going Under Cover: Passing as Human; Artificial Interest: A Step on the Road to AI	413
Michael L. Mauldin	
Chapter 25 How not to Imitate a Human Being: An Essay on Passing the Turing Test	431
Luke Pellen	
Chapter 26 Who Fools Whom? The Great Mystification, or Methodological Issues on Making Fools of Human Beings	447
Eugene Demchenko and Vladimir Veselov	
Part IV Afterthoughts on Thinking Machines	
Chapter 27 A Wager on the Turing Test	463
Ray Kurzweil and Mitchell Kapor	
Chapter 28 The Gnirut Test	479
Charles Platt	
Chapter 29 The Artilect Debate: Why Build Superhuman Machines, and Why Not?	487
Hugo De Garis and Sam Halioris	
Name Index	511

Chapter 6

If I Were Judge

Selmer Bringsjord

Abstract: I have spent a lot of time through the years attacking the Turing Test and its variants (e.g., Harnad's Total Turing Test). As far as I am concerned, my attacks have been lethal, but of course not everyone agrees. At any rate, in the present paper I shift gears: I pretend that the Turing Test is valid, put on the table a proposition designed to capture this validity, and then slip into the shoes of the judge, determined to deliver a correct verdict as to which contestant is the machine, and which the woman. My strategies for separating mind from machine may well reveal some dizzying new-millennium challenges for Artificial Intelligence.

Keywords Artificial Intelligence, Turing Test

6.1 Introduction

I have spent a lot of time through the years attacking the Turing Test and its variants. For example, the Turing Test and *many* variants (e.g., the *Total Turing Test* and the *Total Total Turing Test*) are overthrown in "Could, how could we tell if, and why should-androids have inner lives?" (Bringsjord, 1995). I have also proposed complete replacements for the Turing Test, in "Creativity, the Turing Test, and the (better) Lovelace test" (Bringsjord et al., 2001). As another example, I have recently carefully refined, extended, and defended Searle's (1980) Chinese Room Argument against the Turing Test (Bringsjord, 1992; Bringsjord and Noel, 2002). As far as I am concerned, these attacks have been lethal, but not everyone agrees (at least not yet). At any rate, in the present paper I shift gears: I pretend that the Turing Test is valid, put on the table a proposition designed to capture this validity, and then slip into the shoes of the judge, determined to deliver a correct verdict as to which contestant is the machine, and which is the woman. My strategies for

Rensselaer Polytechnic Institute

separating mind from machine may well reveal some dizzying new-millennium challenges for Artificial Intelligence (AI).

6.2 Validity of the Turing Test in Declarative Form

The basic architecture of the Turing Test will be familiar to all readers; it is given in Turing's (1950) famous *Mind* paper. A judge must attempt to determine which of two sequestered agents is a machine, and which is a woman. The judge can interact with the agents only via (to modernize things a bit) typed e-mail. Many, many variations on the Turing Test have been suggested. In fact, rather long ago, in my *What Robots Can and Can't Be* (1992), I defined the Turing Test Sequence, which assumes tests ranging from those less demanding than Turing Test (e.g., judges cannot know anything about AI, and can only ask questions about a small, determinate domain), to those that – like Kugel's (1990) – require contestants to have a capacity for infinitary processing. (When I introduced the sequence, I also asserted that, sooner or later, a machine can be built by us to pass any and every test in it.) Many of the variations in this sequence have been offered to supplant the original Turing Test, which even to fans of “Strong” AI seems to be afflicted by a certain myopia (e.g., the Turing Test ignores sensorimotor behavior in favor of the purely linguistic variety). However, I now assume for the present paper that the original Turing Test is in fact valid.

Now let us get a little bit more precise about what it means to say that the Turing Test is valid. One possibility is:

(TT.) For every computer c , if c passes Turing Test, then c is conscious.

It may strike you as odd, if not flatly wrong, that we have turned the focus upon consciousness. Someone might object, specifically, as follows: ‘Turing presents his test as a test for “intelligence” or “thought.” Your interpretation of Turing Test cheats over toward precisely what Turing sought to dodge: phenomenal consciousness and qualia.’ But a careful reading of Turing's (1950) paper supports my interpretation. Specifically, I draw your attention to the section therein entitled ‘(4) The Argument from Consciousness.’ The argument from consciousness is one Turing takes to be well-expressed in Professor Jefferson's Lister Oration of 1949:

Not until a machine can write a sonnet or compose a concerto because of thoughts and emotions felt, and not by the chance fall of symbols, could we agree that machine equals brain – that is, not only write it, but know that it had written it. No mechanism could feel (and not merely artificially signal, an easy contrivance) pleasure at its successes, grief when its valves fuse, be warmed by flattery, be made miserable by mistakes, be charmed by sex, be angry or depressed when it cannot get what it wants (Turing, 1950).

What is Turing's response? I am afraid he gives an absolutely dreadful rejoinder; it amounts to:

If one refuses to agree that passing Turing Test (when, say, sonnets are requested by the judge) implies consciousness, one must accept solipsism. Since solipsism is false, one cannot refuse in the manner indicated by Jefferson.

Though you and I no doubt reject solipsism, I cannot for the life of me imagine where Turing gets his first premise. To take a guess, perhaps the underlying idea is that if one rejects the notion that *c*'s passing implies that *c* is conscious, one will have no solution (or candidate solution) to the problem of other minds – save for solipsism. But since the literature is filled with proposed solutions to the problem of other minds that make no appeal to the Turing Test, if this is Turing's underlying idea, it is a fatally flawed one. At any rate, Turing's argument is beside the point we have by now made, which is that Turing clearly holds that if a computer (or robot) passes his test (in part by producing sonnets, etc.), it follows that it is conscious.

However, I propose in the present paper to assume a version of Turing's claim for the Turing Test that insulates him from all the sorts of attacks that naturally arise if (TT₁) is the target. (The thesis claiming that passing Turing Test ensures consciousness is vulnerable to attacks based on thought-experiments in which passing occurs, but subjective awareness is absent. Searle's (1980) CRA is an example of such an attack.) Specifically, I propose something like:

(TT₂) For every computer *c*, if *c* passes Turing Test, then *c* is as intelligent as human persons.

Notice how charitable a move toward such a cashing out of "the Turing Test is valid" is. (TT₂) makes no reference (at least no overt reference) to invisible mental properties, which would surely gladden Turing's empiricist heart. On the other hand, there is a severe defect in (TT₂): It classifies a Turing Test-passing computer as human-level intelligent, but it does not stipulate that the human contestant is smart! To concretize this point, suppose that a machine competes alongside a 2-year old in the confines of the Turing Test. We can probably safely assume that some computer today, or at least in the near future, could leave me in the dark as to which room houses little, just-learning-to-talk Johnny. What this reveals is that propositions like (TT₂) leave concealed a fact that needs to be uncovered: viz., that there are *a lot* of candidate human contestants.

I see this point as another version of one made by the thinker who proposed Turing Test before (!) Turing: Descartes.¹ Here is the relevant passage:

If there were machines which bore a resemblance to our body and imitated our actions as far as it was morally possible to do so, we should always have two very certain tests by which to recognize that, for all that, they were not real men. The first is that they could never use speech or other signs as we do when placing our thoughts on record for the benefit of others. For we can easily understand a machine's being constituted so that it can utter

¹ Actually, Descartes proposed a test that is much more demanding than the Turing Test (Descartes, 1911), but I do not explain and defend this herein. In a nutshell, if you read the passage very carefully, you will see that Descartes' test is passed only if the computer has the capacity to answer arbitrary questions. A machine which has a set of stored chunks of text that happen to perfectly fit the queries given it during a Turing test would not pass Descartes' test – even though it would pass Turing's.

words, and even emit some responses to action on it of a corporeal kind, which brings about a change in its organs; for instance, if it is touched in a particular part it may ask what we wish to say to it; if in another part it may exclaim that it is being hurt, and so on. But it never happens that it arranges its speech in various ways, in order to reply appropriately to everything that may be said in its presence, as even the lowest type of man can do. And the second difference is, that although machines can perform certain things as well as or perhaps better than any of us can do, they infallibly fall short in others, by which means we may discover that they did not act from knowledge, but only for the disposition of their organs. For while reason is a universal instrument which can serve for all contingencies, these organs have need of some special adaptation for every particular action. From this it follows that it is morally impossible that there should be sufficient diversity in any machine to allow it to act in all the events of life in the same way as our reason causes us to act. (Descartes 1911)

To operationalize Descartes' point about the diversity of human cognition, and the challenge this poses to a machine, I will say that passing the Turing Test is actually at least a ternary relation *Pass* taking as arguments a computer c , a human player h , and a judge j . By hypothesis, I am the judge; let us denote me by b_j . The proposition itself then is:

$$(TT_3) \forall c \forall h (Pass(c, h, b_j) \rightarrow \text{then } c \text{ is as intelligent as human persons})$$

Quantification over contestants has considerable logical payoff. I say this because (TT_2) is vacuously true, and surely that's an unwanted consequence. (TT_2) is vacuously true because for every computer c today, it is simply false that c can pass the Turing Test when stacked against all people. Therefore (TT_3) 's antecedent is false. But by the standard semantics of first-order logic, this makes the entire proposition true. I have discussed this situation (Bringsjord, 1995). And yet there is a problem with (TT_3) : I should not be allowed to pick a *particular* individual. We really should be talking about picking from a *class* of humans. It should be obvious why this is so. If I could pick an individual, and if the competing computer does not know this individual inside and out, I have only to ask a few "private" questions to prevail. Turing did not envisage a situation wherein a machine would be trying to impersonate a specific human being. Here is the solution where C ranges over classes of human persons:

$$(TT_2) \forall c \forall h (h \in C \rightarrow (Pass(c, h, b_j) \rightarrow \text{then } c \text{ is as intelligent as human persons}))$$

(TT_2) needs further refinement. The problem is that it makes no reference to how long the test is to last. Following what I did in Bringsjord (1995), let us let τ denote the length of time the Turing Test is played for; we can follow this notation here. We can expand the key relation to take four arguments, and we can quantify over intervals. So we have:

$$(TT_1) \forall c \forall h \forall \tau (h \in C \rightarrow (Pass(c, h, b_j, \tau) \rightarrow \text{then } c \text{ is as intelligent as human persons}))$$

Now (TT_3) can be used to anchor an adversarial relationship between me as judge and computer as competitor. That is, it is now under my control as to what length of time to opt for, and what sort of human to select (= what class C to pick the human from) as my other interlocutor. For reasons already cited, it would not be a particularly clever strategy for me to request that h be instantiated to a 2-year-old, and that τ be set to one minute. So, what *would* be an intelligent pick? To that we now turn.

6.3 My Strategies

To qualify as a Turing Test-passer, I will require that the class C of humans against which c is matched be the *union* of a number of classes. Each of the following four strategies (standardized tests of mental ability, tests for “irrationality”, requests that certain paradoxes be solved, and tests for literary creativity) is associated with at least one subclass within this union.

6.3.1 Standardized Tests as a First Hurdle

My first strategy would be to require the computer c to excel on all established, standardized tests – tests of intelligence, creativity, spatial reasoning, and so on. This strategy relates to a form of AI invented by me and Bettina Schimanski; we refer to it as Psychometric AI (Bringsjord and Schimanski, 2003; Bringsjord and Zenzen, 2003). Together, the two of us are attempting to build a robot capable of reaching high performance on *all* standardized tests; this robot is PERI, who “lives” in the Rensselaer AI and Reasoning Lab. PERI is shown in Fig. 6.1. I have a pretty good idea of how demanding it is to pass the hurdle of reaching such performance because of these efforts. In the present essay, I say but a few words about this challenge.

I would first insist the c take a “broad” IQ test, on which it would need to score as high as the best-performing humans in order to remain a candidate for passing Turing Test. What could possibly be a more obvious strategy? After all, by (TT_3), the Turing Test is about *intelligence*, and we have many established broad tests of intelligence. But what is meant by a “broad” IQ test? This question is best answered by turning to an example of a “narrow” intelligence test. The example of a narrow IQ test that I give here is one that Bettina Schimanski and I have built for an artificial agent to crack: Raven’s (1962) Progressive Matrices (RPM). An example of the type of problem that appears in RPM is shown in Fig. 6.2, which is taken from Carpenter et al. (1990). The query in each RPM problem is implicit, and invariant: viz., pick the option that preserves the vertical and horizontal patterns. Obviously, the items in question do not relate to general (declarative) knowledge, common sense, or the ability to communicate in natural language.

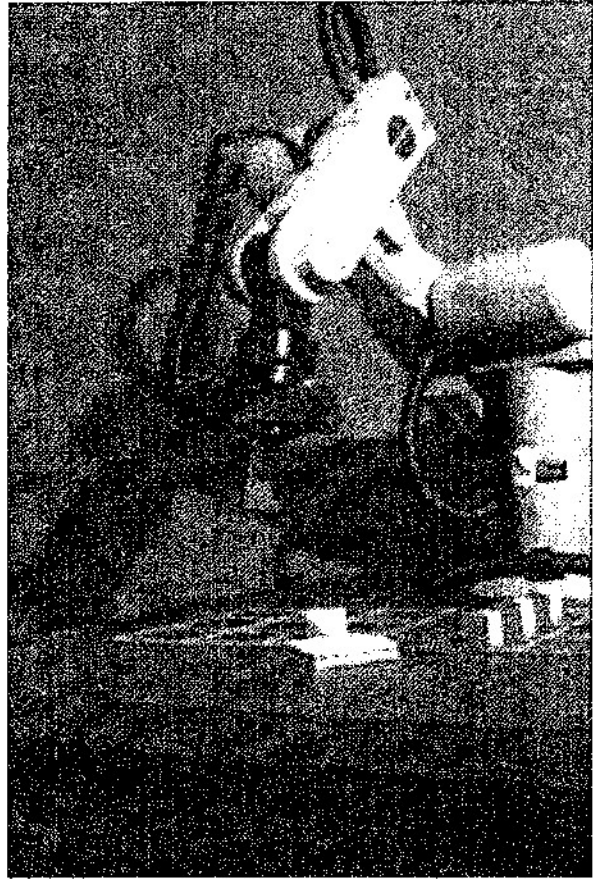


Fig. 6.1 PERI working on the block design puzzle

Sample (& Simple) RPM Problem

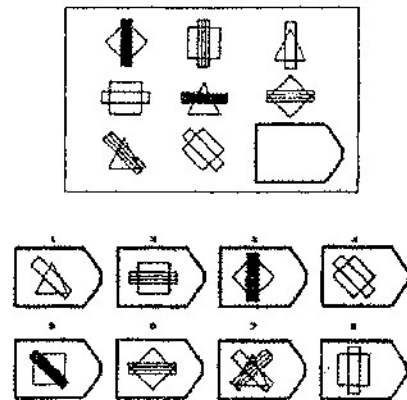


Fig. 6.2 A simple RPM problem "cracked" by a Bringsjord-created intelligent agent

An example of a broad intelligence test is the WAIS (Wechsler Adult Intelligence Scale, available from Psychological Corporation). One of the subtests on the WAIS is block design, in which test-takers must assemble cubes whose sides have different colored patterns to make a larger pattern given as a goal state. (Figure 6.1 shows PERI successfully completing a block design puzzle.) But the WAIS also contains some harder subtests. For example, there is a subtest in which, in conversation, subjects are asked questions designed to determine whether or not they can reason in “common-sense” fashion. For example, subjects might be asked to explain why the tires on automobiles are made of rubber, instead of, say, wood. Another subtest on the WAIS, picture completion, requires that coherent stories be assembled from snapshots of various situations. To our knowledge, no present-day AI system can correctly answer arbitrary questions of this sort, or solve these kinds of narrative-related problems. Hopefully this gives you a tolerably clear sense of the distinction between “narrow” and “broad” in this context, and perhaps you also appreciate that it would be no small feat for a computer to solve these sorts of problems.

I would give the would-be Turing Test-passer not just intelligence tests, but, as I said, *all* established, standardized tests of mental ability. For example, I would subject *c* to tests of creativity. These tests, in the context of building artificial agents, are discussed by Bringsjord and Ferrucci (2000). On the assumption that a would-be Turing Test-passer clears the first hurdle by matching brilliant humans on all tests of mental ability, I move to the second.

6.3.2 Irrationality

In the next hurdle, I would pick as my class of humans those with college educations, but no advanced training in formal reasoning. I would give the computer *c* problems like this one²:

Problem 1

1. If there is a king in the hand, then there is an ace, or else if there is not a king in the hand, then there is an ace.
2. There is a king in the hand.

Given these premises, what can one infer?

Almost certainly your own verdict is this: One can infer that there is an ace in the hand. Your verdict seems correct, even perhaps *obviously* correct, and yet a little

²Problem 1 (or, actually, Illusion 1) is from “How to make the impossible seem probable” (Johnson-Laird and Savary, 1995). Variations are presented and discussed in “Rules and illusions: a critical study of Rips’s” (Johnson-Laird, 1997a).

logic suffices to show that not only are you wrong, but that in fact what you can infer is that there *isn't* an ace in the hand!³

To see this, note that “or else” is to be understood as exclusive disjunction,⁴ so (using obvious symbolization) the two premises become

$$(1') ((K \rightarrow A) \vee (\neg K \rightarrow A)) \wedge \neg((K \rightarrow A) \wedge (\neg K \rightarrow A))$$

$$(2') K$$

Figure 6.3 shows a proof in the standard first-order Fitch-style system \mathcal{F} , constructed in HYPERPROOF (Barwise and Etchemendy, 1994), that demonstrates that from these two given one can correctly conclude $\neg A$.

• $((K \rightarrow A) \vee (\neg K \rightarrow A)) \wedge \neg((K \rightarrow A) \wedge (\neg K \rightarrow A))$	✓ Given
• K	✓ Given
• $\neg((K \rightarrow A) \wedge (\neg K \rightarrow A))$	✓ \wedge Elim
• $(K \rightarrow A) \vee (\neg K \rightarrow A)$	✓ Taut Con
• $(K \rightarrow A)$	✓ Assume
• $K \wedge A$	✓ Taut Con
• $\neg A$	✓ \wedge Elim
• $\neg(K \rightarrow A)$	✓ Assume
• $\neg K \wedge \neg A$	✓ Taut Con
• $\neg A$	✓ \wedge Elim
• $\neg A$	✓ \vee Elim

Fig. 6.3 A proof that there is no ace in the hand in \mathcal{F}

³You should not feel bad about succumbing to Illusion 1; after all, you have a lot of company. Johnson-Laird has recently reported that

“Only one person among the many distinguished cognitive scientists to whom we have given [Illusion 1] got the right answer; and we have observed it in public lectures – several hundred individuals from Stockholm to Seattle have drawn it, and no one has ever offered any other conclusion” (Johnson-Laird, 1997b).

Time and time again, in public lectures, I have replicated Johnson-Laird’s numbers – presented in (Johnson-Laird and Savary, 1995) – among those not formally trained in logic.

⁴Even when you make the exclusive disjunction explicit, the results are the same. For example, you still have an illusion if you use

Illusion 1’

(1’) If there is a king in the hand then there is an ace, or if there is not a king in the hand then there is an ace, but not both.

(2’) There is a king in the hand.

Given these premises, what can you infer?

While psychologists of reasoning create problems like this in order to carry out experiments on humans, AI researchers might be interested in programming computers that *generate* such illusions. Indeed, this is exactly my interest, and I have elsewhere discussed how it is that I came to work on an algorithm able to generate this problem:

Problem 2

3. The following three assertions are either all true or all false:
 - If Billy is happy, Doreen is happy
 - If Doreen is happy, Frank is as well
 - If Frank is happy, so is Emma
4. The following assertion is definitely true: Billy is happy.

Can it be inferred from (3) and (4) that Emma is happy?

Most human subjects answer “Yes”, but get the problem wrong – because their *reasons* for answering with an affirmative are incorrect. They say “Yes” because they notice that since Billy is happy, if the three conditionals are true, one can “chain” through them to arrive at the conclusion that Emma is happy. But this is only part of the story, and the other part has been ignored: viz., that it could be that all three conditionals are false. Some subjects realize that there are two cases to consider (conditionals all true, conditionals all false), and because they believe that when the conditionals are all false one cannot prove that Emma is happy, they respond with “No”. But this response is also wrong. The correct response is “Yes”, because *in both cases it can be proved that Emma is happy*. This can be shown using propositional logic; the proof, once again constructed in HYPERPROOF, is shown in Fig. 6.4. This proof establishes

$$\{\neg(B \rightarrow D), \neg(D \rightarrow F)\} \vdash E$$

Note that the trick is exploiting the inconsistency of the set $\{\neg(B \rightarrow D), \neg(D \rightarrow F)\}$ in order to get a contradiction. Since everything follows from a contradiction, E can then be derived.

What does all this have to do with the Turing Test? I would expect to unmask many would-be Turing Test-passers with problems like these, for the simple reason that a machine, *ceteris paribus*, would not be fooled. That is, the machine might well parse these problems *correctly*, and would then reason them out in accordance with normatively correct structures from symbolic logic; that is, reason them out essentially as shown in Figs. 6.4 and 6.3. In other words, the machine must be smart enough to appear dull. Note that a computer smart enough to meet this challenge would presumably be capable of some form of meta-reasoning. To *really* test for meta-reasoning I would next see if my digital opponent can solve a paradox or two. Suppose that “generic” paradox P is the derivation of contraction $\phi \wedge \neg \phi$ from set Φ of premises. It is unavoidable that, in trying to solve P , one must reason about this reasoning; that is, it is unavoidable that the attempt to solve P involves meta-reasoning. In addition, P may involve propositions that are themselves very expressive, so that modeling them requires very sophisticated modes of representation and

\diamond $\vdash ((H(b) \rightarrow H(d)) \wedge (H(d) \rightarrow H(f)) \wedge (H(f) \rightarrow H(e))) \vee$ $(\neg(H(b) \rightarrow H(d)) \wedge \neg(H(d) \rightarrow H(f)) \wedge \neg(H(f) \rightarrow H(e)))$ $\vdash H(b)$ $\vdash (H(b) \rightarrow H(d)) \wedge (H(d) \rightarrow H(f)) \wedge (H(f) \rightarrow H(e))$ $\vdash H(b) \rightarrow H(d)$ $\vdash H(d)$ $\vdash H(d) \rightarrow H(f)$ $\vdash H(f)$ $\vdash H(f) \rightarrow H(e)$ $\vdash H(e)$ $\vdash (\neg(H(b) \rightarrow H(d)) \wedge \neg(H(d) \rightarrow H(f)) \wedge \neg(H(f) \rightarrow H(e)))$ $\vdash \neg(H(b) \rightarrow H(d))$ $\vdash H(b) \wedge \neg H(d)$ $\vdash \neg(H(d) \rightarrow H(f))$ $\vdash H(d) \wedge \neg H(f)$ $\vdash \neg H(e)$ $\vdash H(d) \wedge \neg H(d)$ $\vdash H(e)$ $\vdash H(e)$	\checkmark Given \checkmark Given \checkmark Given \checkmark Assume \checkmark \wedge Elim \checkmark \rightarrow Elim \checkmark \wedge Elim \checkmark \rightarrow Elim \checkmark \wedge Elim \checkmark \rightarrow Elim \checkmark Assume \checkmark \wedge Elim \checkmark Taut Con \checkmark \wedge Elim \checkmark Taut Con \checkmark Assume \checkmark Taut Con \checkmark \neg Intro \checkmark \vee Elim
--	--

Fig. 6.4 A proof that “Emma is happy” in F

reasoning. Finally, many paradoxes are infinitary in nature, which implies that agents who would solve them must be able to, at least in some sense, “grasp” infinitary reasoning. All of this may seem rather vague to you. Fortunately, I have recently presented a paradox that makes my points concrete: this paradox is an “infinitized” version of Yablo’s (1993) paradox. I present it now as an excerpt from “The mental eye defense of an infinitized version of Yablo’s paradox” (Bringsjord and van Heuveln, 2003), in which a full discussion can be found.

6.3.3 A Paradox for a Computer to Tackle

The paradox runs as follows⁵:

Recall the familiar natural numbers $\mathbf{N} = \{0, 1, 2, \dots\}$. With each $n \in \mathbf{N}$ associate a sentence as follows, using a truth predicate, T :

⁵I specify an infinitary version of Yablo’s paradox, expressed in the “background” logic that allows for meta-proofs regarding infinitary logical systems like $\mathcal{L}_{\omega_1, \omega}$. This system is presented in encapsulated form in *Mathematical Logic* (Ebbinghaus et al., 1984), from which the student interested in infinitary logic can move to *Languages with Expressions of Infinite Length* (Karp, 1964), then to *Model Theory for Infinitary Logic* (Keisler, 1971), and then *Large Infinitary Languages* (Dickmann, 1975).

$$s(0) = \forall k(k > 0 \rightarrow \neg T(s(k)))$$

$$s(1) = \forall k(k > 1 \rightarrow \neg T(s(k)))$$

$$s(2) = \forall k(k > 2 \rightarrow \neg T(s(k)))$$

$$s(3) = \forall k(k > 3 \rightarrow \neg T(s(k)))$$

⋮

Expressed with help from the infinitary system $L\omega_1\omega$, we can say that

$$s(0) = \bigwedge_{k>0} \neg T(s(k)), \quad s(1) = \bigwedge_{k>1} \neg T(s(k)), \quad s(2) = \bigwedge_{k>2} \neg T(s(k)) \dots$$

Next, suppose that $T(s(0))$. From this it follows immediately that $\neg T(s(1)) \wedge \neg T(s(2)) \dots$, which in turn implies by conjunction elimination in $\mathcal{L}\omega_1\omega$ that $\neg T(s(1))$. But in addition, if $T(s(0))$ is true, it follows again that $\neg T(s(1)) \wedge \neg T(s(2)) \dots$, and hence that

$$\neg T(s(2)) \wedge \neg T(s(3)) \dots,$$

which implies that $T(s(1))$. By reduction, we can infer $\neg T(s(0))$. The same indirect proof can be given to show

$$\neg T(s(1)), \neg T(s(2)), \neg T(s(3)) \dots$$

Hence we can infer by the ω -rule

$$\frac{\alpha(1), \alpha(2), \dots}{\alpha(n)}$$

that

$$(*) \bigwedge_{k \in \mathbb{N}} \neg T(s(k))$$

Hence $\neg T(s(1)), \neg T(s(2)), \neg T(s(3)) \dots$, that is, $T(s(0))$. But $\neg T(s(0))$ follows from (*) – contradiction.

I have serious doubts that a computer will ever be able to solve a paradox like this one. (Can you solve it? Or is it a true paradox?) But perhaps you are wondering what a solution would consist of. Well, one possible type of solution for a paradox P is to provide a formal theory on which the premises in P , Φ , are true, but the contradiction cannot be derived. This kind of solution for the paradox I have presented would be remarkable, because the formal theory will in some sense subsume a logical system that in and of itself far exceeds what machines can today (in any sense of the word) understand. It is hard to see how even a future machine would achieve this understanding.

If I assume for the sake of argument that some computer does pass the Turing Test with C set to professional logicians and formal philosophers (a group up to the challenge of solving paradoxes), I resort to my final weapon: literary creativity.

6.3.4 Literary Creativity

The idea here is to see if the computer in the Turing Test is capable of producing stories indistinguishable from those produced by accomplished authors of literary fiction. This means that the class C of humans against which the computer c is matched includes the likes of John Updike, Toni Morrison, Mark Helprin, and so on. I have refined this scenario into what I call the “short short story game”, or just S³G for short. The idea is simple; it is summed up in Fig. 6.5. The computer and human both receive one relatively simple sentence from me, say: “Barnes kept the image to himself, kept the horror locked away as best he could.” (For a much better one, see the “loaded” sentence shown in Fig. 6.5.⁶) Both human and machine must now fashion a short short story of no more than 500 words. The machine’s objective,

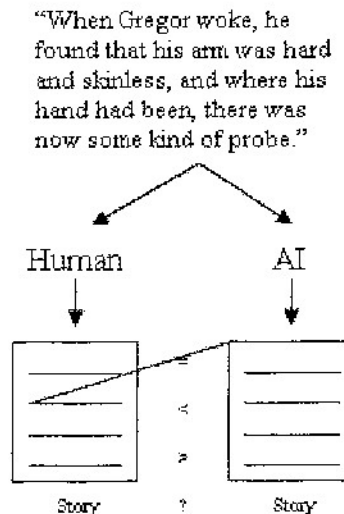


Fig. 6.5 S³G

⁶The actual opening is as follows:

As Gregor Samsa awoke one morning from uneasy dreams he found himself transformed in his bed into a gigantic insect. He was lying on his hard, as it were armor-plated, back and when he lifted his head a little he could see a dome-like brown belly divided into stiff arched segments on top of which the bed quilt could hardly keep in position and was about to slide off completely. His numerous legs, which were pitifully thin compared to the rest of his bulk, waved helplessly before his eyes. (Kafka, 1948)

of course, is to produce narrative that leaves me in the dark as to whether it is authored by mind or machine. For reasons explained in *Artificial Intelligence and Literary Creativity: Inside the Mind of Brutus, a Storytelling Machine* (Bringsjord and Ferrucci, 2000), I think it will be exceedingly difficult for any computer to match the likes of John Updike. Of course, at 500 words, it may be possible. As the ultimate test, I would as judges allow word length to reach that of a full-length novel (which would of course require that I increase τ considerably).

Acknowledgment I am indebted to my colleague Yingnui Yang, cocreator of mental metalogic theory.

References

- Barwise, J. and Etchemendy, J., 1994, *Hyperproof*, CSLI, Stanford, CA.
- Bringsjord, S., 1992, *What Robots Can and Can't Be*, Kluwer, Dordrecht, The Netherlands.
- Bringsjord, S., 1995, Could, how could we tell if, and why should-androids have inner lives? in: *Android Epistemology*, K. Ford, C. Glymour, and P. Hayes, eds., MIT Press, Cambridge, MA, pp. 93–122.
- Bringsjord, S. and Ferrucci, D., 2000, *Artificial Intelligence and Literary Creativity: Inside the Mind of Brutus, a Storytelling Machine*, Lawrence Erlbaum, Mahwah, NJ.
- Bringsjord, S. and Noel, R., 2002, Real robots and the missing thought experiment in the Chinese room dialectic, in: *Views into the Chinese Room: New Essays on Searle and Artificial Intelligence*, J. Preston and M. Bishop, eds., Oxford University Press, Oxford, pp. 144–166.
- Bringsjord, S. and Schimanski, S., 2003, What is Artificial Intelligence? · Psychometric AI as an answer, *Proceedings of the 18th International Joint Conference on Artificial Intelligence (IJCAI-03)*, Morgan Kaufmann, San Francisco, CA, pp. 887–893.
- Bringsjord, S. and van Heuveln, B., 2003, The mental eye defense of an infinitized version of Yablo's paradox, *Analysis* 63(1): 61–70.
- Bringsjord, S. and Zenzen, M., 2003, *Superminds: People Harness Hypercomputation, and More*, Kluwer, Dordrecht, The Netherlands.
- Bringsjord, S., Ferrucci, D., and Bello, P., 2001, Creativity, the Turing Test, and the (better) Lovelace test, *Minds and Machines* 11: 3–27.
- Carpenter, P., Just, M., and Shell, P., 1990, What one intelligence test measures: a theoretical account of the processing in the Raven progressive matrices test, *Psychological Review* 97: 404–431.
- Descartes, R., 1911, *The Philosophical Works of Descartes*, Vol. 1, translated by Elizabeth S. Haldane and G. R. T. Ross, Cambridge University Press, Cambridge.
- Dickmann, M. A., 1975, *Large Infinitary Languages*, North-Holland, Amsterdam, The Netherlands.
- Ebbinghaus, H. D., Flum, J., and Thomas, W., 1984, *Mathematical Logic*, Springer, New York.
- Johnson-Laird, P., 1997a, Rules and illusions: a critical study of Rips's, *The Psychology of Proof*, *Minds and Machines* 7(3): 387–407.
- Johnson-Laird, P. N., 1997b, An end to the controversy? A reply to Rips, *Minds and Machines* 7: 425–432.
- Johnson-Laird, P. and Savary, F., 1995, How to make the impossible seem probable, *Proceedings of the 17th Annual Conference of the Cognitive Science Society*, Lawrence Erlbaum, Hillsdale, NJ, pp. 381–384.
- Kafka, F., 1948, The metamorphosis, in: *The Penal Colony*, F. Kafka, T. W. Muir, and E. Muir, eds., Schocken Books, New York.

- Karp, C., 1964, *Languages with Expressions of Infinite Length*, North-Holland, Amsterdam, The Netherlands.
- Keisler, H., 1971, *Model Theory for Infinitary Logic*, North-Holland, Amsterdam, The Netherlands.
- Kugel, P., 1990, Is it time to replace Turing's test? Paper presented at *Artificial Intelligence: Emerging Science or Dying Art Form?*, sponsored by AAAI and the State University of New York's program in Philosophy and Computer and Systems Sciences, the University at Binghamton, New York, June 27.
- Raven, J. C., 1962, *Advanced Progressive Matrices Set II*, H. K. Lewis, London. Distributed in the USA by The Psychological Corporation, San Antonio, TX.
- Searle, J., 1980, Minds, brains and programs, *Behavioral and Brain Sciences* 3: 417–424; <http://members.aol.com/NeoNoetics/MindsBrainsPrograms.html>.
- Turing, A., 1950, Computing machinery and intelligence, *Mind* 59(236): 433–460.
- Yablo, S., 1993, Paradox without self-reference, *Analysis* 53: 251–252.