

Lessons from a Restricted Turing Test

he English logician and mathematician Alan Turing, in an attempt to develop a working definition of intelligence free of the difficulties and philosophical pitfalls of defining exactly what constitutes the mental process of intelligent reasoning, devised a test, instead, of intelligent behavior. The idea, codified in his celebrated 1950 paper "Computing Machinery and Intelligence" [28], was specified as an "imitation game" in which a judge attempts to distinguish which of two agents is a human and which a computer imitating human responses by engaging each in a wide-ranging conversation of any topic and tenor. Turing's reasoning was, presuming that intelligence was only practically determinable behaviorally, that any agent that was indistinguishable in behavior from an intelligent agent was, for all intents and purposes, intelligent.

It is presumably uncontroversial that humans are intelligent as evidenced by their conversational behavior. Thus, any agent that can be mistaken by virtue of its conversational behavior with a human must be intelligent. As Turing himself noted, this syllogism argues that the criterion provides a sufficient, but not necessary, condition for intelligent behavior. The game has since become known as the "Turing Test," a term that has eclipsed even his eponymous machine in Turing's terminological legacy. Turing predicted that by the year 2000, computers would be able to pass the Turing Test at a reasonably sophisticated level, in particular, that the average interrogator would not be able to identify the computer correctly more than 70% of the time after a five-minute conversation.

On November 8, 1991, an eclectic group including academics, business people, press, and passers-by filled two floors of Boston's Computer Museum for a tournament billed as the first actual administration of the Turing Test. The tournament was the first attempt on the recently constituted Loebner Prize established by New York theater equipment manufacturer Hugh Loebner and organized by Robert Epstein, president emeritus of the Cambridge Center for Behavioral Studies, a research center specializing in behaviorist psychology. The Loebner Prize was administered by an illustrious committee headed by Daniel Dennett, Distinguished Professor of Arts and Sciences and director for Cognitive Studies, Tufts University, and included: Robert Epstein; Harry Lewis, Gordon McKay Professor of Computer Science, Harvard University; H. McIlvaine Parsons, senior research scientist, HumRRO; Willard van Orman Quine, Edgar Pierce Professor of Philosophy Emeritus, Harvard University; and Joseph Weizenbaum, professor of computer science emeritus, Massachusetts Institute of Technology.

The prize committee spent almost two years in planning the structure of the tournament. Because this was to be a real competition, rather than a thought experiment, there would be several computer contestants, and therefore several confederates would be needed as well.1 It was decided that there would be 10 agents all together. In the event, 6 were computer programs. Ten judges would converse with the agents and score them. The judges and confederates were both selected from the general public on the basis of a newspaper employment advertisement that required little beyond typing ability, and then screened by interview with the prize committee. They were chosen to have "no special expertise in computer science."

The committee realized early that, given the current state of the art, there was no chance that Turing's test, as originally defined, had the slightest chance of being passed by a

¹We follow the prize committee's terminology in using the terms "confederate," "contestant," and "judge" for the computer program entrants, the humans being compared against, and the human interrogators performing the evaluation, respectively. We use the term "agent" for both confederates and contestants.

computer program. Consequently, they attempted to adjust both the structure of the test and the scoring mechanism, to allow the computers a fighting chance. In particular, the following two rules were added to dramatically restrict Turing's test.

• *Limiting the topic:* In order to limit the amount of area the contestant programs must be able to cope with, the topic of the conversation was to be strictly limited, both for the contestants and the confederates. The judges were required to stay on the subject in their conversations with the agents.

• Limiting the tenor: Further, only behavior evinced during the course of a natural conversation on the single specified topic would be required to be duplicated faithfully by the contestants. The operative rule precluded the use of "trickery or guile. Judges should respond naturally as they would in a conversation with another person." (The method of choosing judges served as a further measure against excessive judicial sophistication.)

As will be seen, these two rules limiting the topic and tenor of the discussion—were quite problematic.

The prize committee specified that there be independent referees stationed in several locations: several in the rooms with the judges and confederates to answer questions concerning interpretation of the preceding rules, and one in the auditorium to serve as a sort of roving ombudsman. I was a referee in the confederates' room, and can assure that my colleagues' efforts and mine there were hardly needed; the confederates performed admirably. Reports from the other referees indicated the same for the judges.²

Hugh Loebner placed only two restrictions on the setting up of the competition by the prize committee: that a competition be held each year, and that a prize be awarded at each competition. The prize at this first competition was a nominal \$1,500, although Loebner has reportedly earmarked \$100,000 for the first computer program to pass the full Turing Test at some later running of the competition. (Costs for the running of the competition itself were paid for by grants from the National Science Foundation and the Sloan Foundation.)

To determine the prize winner, an ingenious scoring mechanism was devised. The Turing Test involves a single binary decision, which is either right or wrong. But to determine a winner, the contestants had to be ranked, so each judge was required to place all of the agents in order from the apparently least human to most human. This alone induced the ordering on the basis of which the prize would be awarded. The contestant with the highest average rank would be deemed the winner of the tournament. However, this does not allow a direct reconstruction of the results of the 100 implicit binary decisions that might be made: which of the agents were humans, and which computers? To allow for this to be deduced as well, each judge was requested to place a single line separating the ranked agents into two groups. Those to the right of the line were claimed by that judge to be humans, those to the left computers. (See Figure 1.) The judges were told that at least two of the agents were human confederates, and at least two were computer contestants, thus limiting the number of places that the line could be (rationally) placed. The binary decisions could then be read from the rankings by noting on which side of the line each agent fell. This

demarcation process was not used in the awarding of the prize, but was carried out for its informational value alone.

The Event

The tournament was to begin at 1:00 P.M. on the scheduled Friday. One room of the computer museum was set up with 10 terminals for the judges, each labeled with a code letter and the specified topic for conversing with the associated agent. In a back room, hidden from the publicly accessible part of the museum for obvious reasons, five computers had been set up to serve the four confederates. (One terminal was intended to be a backup, and in case it was not needed, to be connected to a publicly accessible terminal so that press and the public could interact with it in a separate Turing Test.) In a large auditorium, the 10 conversations were projected, each on its own screen around the perimeter of the room, and A. K. Dewdney provided running commentary.

Unfortunately, there were serious technical difficulties with the rented computer equipment that had been set up for the confederates. None of the three IBM computers could be made to appropriately interact over the prepared lines with their companion terminal in the judges' room. (The two DEC workstations seemed to work fine.) After almost two hours of unsuccessful last-minute engineering, the prize committee decided to begin the competition with only two confederates in place (just the number the judges had been told was the minimum), reducing the number of agents to eight. The time that each judge had to converse with each agent was shortened from approximately 15 minutes to approximately

Rank order of the terminals									
Least human-like	1	2	з	4	5	6	7	8	Most human-like
	${\mathcal B}$	я	\mathcal{E}	\mathcal{D}	С	\mathcal{F}	${\mathcal H}$	G	

Figure 1. Mock-up of the form used to implement the scoring method for the first Loebner competition. The judge writes the letters corresponding to the terminals in order from least to most human-like, and draws a line purporting to separate the computer contestants from the human confederates. In this case, the line has been drawn so that three of the terminals (F, H, and G) were deemed to be connected to humans.

²The confederate room referees, in addition to me, were Susan Cole Dranoff, an attorney at the firm of Ropes and Gray, and Burton Dreben, an Edgar Pierce Professor of Philosophy Emeritus, Harvard University. The judge room referees were Ned Block, professor of philosophy, Massachusetts Institute of Technology, Robert W. Furlong, patent attorney, and Robert Harford, professor of radiology, Thomas Jefferson University. Thomas Sheridan, professor of engineering, MIT, served in the auditorium.

7 in order to accommodate the press's deadlines.

The topics chosen by the 6 contestants were of the sort appropriate for a cocktail party (burgundy wines, dry martinis, small talk, whimsical conversation, dissatisfactions in relationships) or perhaps a child's birthday party (second-grade school topics). The two participating confederates chose to converse on Shakespeare and women's clothing. In the end, and perhaps unsurprisingly, the average rankings placed the two human confederates as "more human-like" than the six contestants. The highestranked contestant, Joseph Weintraub's program (topic: whimsical conversation) was awarded the \$1,500 prize.

The surprises were in the bipartite classifications. Five judges ranked the top contestant as human, and there were eight instances of such misclassifications of computers as humans overall. (For the record, the contestant discussing relationship problems, the entry of Kenneth Colby, well known for his PARRY program, which simulates paranoid responses, was twice rated human. The simulated seven-year old accounted for the final misclassification.)

Cynthia Clay, the Shakespeare aficionado, was thrice misclassified as a computer. At least one of the judges made her classification on the premise that "[no] human would have that amount of knowledge about Shakespeare" [26]. Lisette Gozo was honored as the most human of the agents for her discussion of women's clothing, although one judge rated two computer programs above her.

Perhaps the most conspicuous characteristic of the six computer programs was their poor performance. It was widely recognized that computer experts could readily distinguish the contestants from the confederates. Indeed, many of the techniques being used by the programs were easily spotted by those familiar with the ELIZA program that prize committee member Weizenbaum developed in 1965. The repetition of previous statements verbatim (subject only to pronominal adjustments, sometimes wrong), answers transparently keyed to trigger words, and similar tricks of the ELIZA trade were common. For example, the following example from the whimsical conversation program is illustrative of the regurgitation technique:³

01:11:42 Judge 10

It has been a long day! Now I think I'm ready for some fun. The weather isn't very pleasant today. Do you have any jokes that could lift the clouds?

01:13:01 Terminal 5

Everybody talks about the weather, but nobody seems to do much about it. Why do you need to tell me now you think you're ready for some fun? Groannnnn.....

Note that the first sentence in the reply, which concerns the weather, seems to be triggered by the keywords "weather" or "clouds" in the judge's contribution. The second sentence is comprised of the judge's second sentence, with pronouns transposed, and preceded by the introductory "Why do you need to tell me."

Such techniques are subject to well-known and obvious errors, as can be exemplified in the transcripts of interactions with the winning computer entrant. For instance, a naive application of pronoun transposition interacts poorly with tag questions like "aren't you."

01:17:52 Judge 10

You're trying to make this difficult for me aren't you? I just explained that I was ready for fun!

01:18:53 Terminal 5

Would you pre]fer it if it was not true that IO re trying to make this difficult for you aren't you? I just explained that i was ready for fun! I am, therfore I compute.... What makes you think I am trying to make this difficult for you aren't i?

Similarly, repetition of a compound phrase is conversationally bizarre.

01:35:26 Judge 6

I can't say I understand the question...but the whimsy of a poet is one thing and the whimsy of most of us is another.

01:36:51 Terminal 5

How do you know you can t say you understand the question...but the whimsy of a poet is one thing and the whimsy of most of us is another.

One reporter noted that Weizenbaum himself was "disturbed" by how easily people were fooled by these programs [26], and more than one of the judges reported that they were disappointed in the programs' capabilities after their expectations had been raised by interacting with ELIZA in the interviewing process. Epstein, in a speech after the event, noted he had learned from the day's proceedings that "little progress has been made in the last 25 years," that is, since ELIZA.

Analysis

The obvious question, then, is how to reconcile the apparent success of the programs in fooling judges with their patently low technology. Clearly, part of the answer relies on the phenomenon that P. T. Barnum used to amass a fortune: people are easily fooled, and are especially easily fooled into reading structure into chaos, reading meaning into nonsense. This accounts for the popularity of newspaper horoscopes and roadside psychics. This is not a flaw in the human mental capacity. Sensitivity to subtle patterns in our environment is extremely important to our ability to perceive, learn, and communicate. Clouds look like ships, and Rorschach blots seem like vignettes. How much different is interpreting non sequitur as whimsical conversation?

Ned Block, a professor of philosophy at MIT (and by coincidence a referee at the competition, stationed with the judges) has argued that the Turing Test is a sorely inadequate test of intelligence because it relies solely on the ability to fool people [3].⁴ Certainly, it has been known since Weizenbaum's surprising experiences with ELIZA that a test based on fooling people is confoundingly simple to pass.

⁴This is not the only case in which exception has been taken to the appropriateness of the Turing Test as a barometer of intelligence. See the discussion in the next section.

³All the following excerpts are taken verbatim from electronic transcripts of the competition provided by and copyright 1991 of the Cambridge Center for Behavioral Studies. No changes were made except for the adjustment of line breaks. In particular, spelling errors and extraneous characters were left as is.

Clouds look like ships, and Rorschach blots seem

like vignettes. How much different is interpreting non-sequiter

as whimsical conversation?

People are even more easily fooled when their ability to detect fooling is explicitly vitiated, for instance, by a prohibition against using "trickery or guile."5 When I asked Weintraub during the postcontest press conference how he would have unmasked his program, his response-typing gibberish in to see if the program spat it back verbatim at a later time à la ELIZA-was certainly outside the established rules. In fact, the referees had discussed that very technique the previous night at a meeting with the prize committee to calibrate our collective understanding of the rules. I pointed out to Weintraub that his response fell under the "trickery-andguile" prohibition, and he took another stab at the question. His second attempt to specify a winning strategy against his program succumbed to the same problem. (It involved repeating questions multiple times.)

Weintraub's problem in answering the question points to the craftiness of his solution to the Loebner Prize puzzle. His entry is unfalsifiable independent of its performance and solely on the basis of the choice of topic. As almost everyone familiar with the rules has noted, whimsical conversation is not in fact a *topic* but a *style* of conversation (at least as practiced by Weintraub's program). And whimsical conversation in the mold of Weintraub's program is essentially nonsensical conversation, a series of non sequiturs. Thus, when Weintraub's program is unresponsive, fails to make any sense, or shows a reckless abandonment of linguistic appropriateness, it, unlike its competitor programs, is operating *as advertised*. It is being "whimsical." At those times when, by chance, the program trips over an especially suggestive response, a judge can grab at it as the real article. (The strategy is reminiscent of that used by the program Racter to create "free-verse" poetry, another unfalsifiable genre.) Weintraub's strategy was an artful dodge of the competition rules. He had found a loophole and exploited it elegantly. I for one believe that, in so doing, he heartily deserved to win.

We might call this winning strategy "PARRY's finesse,"6 after Kenneth Colby's previously mentioned PARRY program [4]. PARRY was designed to engage in a dialogue in the role of a paranoid patient. The program was perhaps the first to be subject to an actual controlled experiment modeled on the Turing Test [5], in which psychiatrists were given transcripts of electronically mediated dialogues with PARRY and with actual paranoids and were asked to pick out the simulated patient from the real. The fact that the expert judges, the psychiatrists, did no better than chance, has been credited to the fact that unresponsiveness and non sequitur are typical behaviors of paranoids. Weizenbaum's response to the experiment-in the form of his own model of a deviant mentalityparodies PARRY's finesse succinctly:7

The contribution here reported should lead to a full understanding of one of man's most troublesome disorders: infantile autism.... It responds *exactly* as does an autistic patient—that is, not at all.... This program has the advantage that it can be implemented on a plain typewriter not connected to a computer at all [29].

Post hoc thinking of this sort can go a long way toward rationalizing the various misclassifications of the whimsical-conversation program or, in the same vein, the program that talks at the level of a second-grader. (Who could fail to give a seven-yearold child the benefit of the doubt?) It leads to noting other insidious forms of scoring bias that crept into the competition. One possible source of such bias, for instance, follows from the technical problems that caused two of the confederates to be eliminated. Once the number of confederates had been reduced to the announced minimum, it became impossible for a judge to rationally place the demarcation line between "humans" and "computers" in such a way as to rate a human as a computer without also rating a computer as a human. Of course, the converse was not true.

This might have accounted for one or two more of the errors. Epstein points out in response to this observation that "(1) Two of the ten judges drew the line after just one entry, in spite of our instructions. (2) Three of the five judges who mistook Weintraub's program for a person rated it above one or both confederates. (3) Two judges mistook a confederate for a computer. In fact, in two (and only two) cases could our instructions have forced the judge to mistake a computer for a person" (personal communication to Harry Lewis, 1992). The third point is, of course, irrelevant; the first hardly gratifying; the second accounted for by Weintraub's

⁵Daniel Dennett, the head of the prize committee, has himself argued against placing "tacit restrictions on the lines of questioning of the judges," calling this a "a common *misapplication* of the sort of testing exhibited by the Turing Test that often leads to drastic overestimation of the powers of actually existing computer systems" [6].

⁶Dennett [6] uses the term "parrying" for the Eliza-like technique of randomly generating a canned response as an option of last resort, a key tool for implementers of PARRY's finesse.

⁷Dennett [6] discusses this and other problems with the PARRY tests. Arbib [1] presents a contravening view, rejoined by Weizenbaum [30].

use of PARRY's finesses; and the final comment is exactly my point.

But post hoc rationalization, like telling your boss off, may be enjoyable at the moment, but is, in the long run, ungratifying. The important questions do not involve microanalysis of the particular competition as run, but the larger questions of the purpose, design, and even existence of the Loebner Prize itself.

Why a Loebner Prize?

There is a long history of argumentation in the philosophical literature opposing the appropriateness of the Turing Test as a litmus test of intelligence. Certain arguments against the effectiveness of the test in answering questions about the intelligence of computers or the possibility of human thought center around the behaviorist nature of the test. Intelligence, it may be claimed, is not determinable simply by surface behavior. Variants of this argument have been given by Block [2], Gunderson [15], and Searle [23, 24]. Others have suggested that the Turing Test is not sufficient in that the behaviors under adjudication are too limited [10, 15]. On the basis of such counterarguments, Moor [18] has argued for a drastically limited view of the Turing Test, not as an operational definition of intelligence at all, but rather as a mode for accumulating evidence leading to an inductive argument for the intelligence of the machine. (See the reply by Stalker [25] and a later clarification by Moor [19] for further arguments.) Moor [20] provides a good introduction to these issues. French [11] provides a strong argument that, as a sufficient condition for intelligence, the Turing Test is so difficult as to be uninteresting. Nonetheless, none of these sorts of presumptive counterarguments to the use of a Turing Test are the basis for the discussion in the remainder of this article. The issues of whether an operational definition of intelligence is appropriate, and whether the particular definition codified in the Turing Test is too narrow, though important questions, can be taken as resolved in favor of the Turing Test for the purposes of the present discussion. Thus, we will side with the behaviorist interpretation favored by the organization administering the prize, the Cambridge Center for Behavioral Studies. Nonetheless, these arguments do provide another strong basis on which to question the appropriateness of the Loebner Prize. A full discussion is, unfortunately, well beyond the scope here, but readers are urged to consult the cited literature. Having sided, for this occasion, with the philosophical appropriateness of Turing's design as a test of intelligent behavior, we turn to the question of whether the Loebner Prize competition is itself an appropriate enterprise.

Prizes for technological advances have existed before, and much can be learned by comparison with previous models.8 Just as humankind has dreamed of mimicking the human power of thought, so have we longed to possess the avian power of flight. Human-powered flight entered the mythology of the ancient Chinese and Romans, the designs of da Vinci, yet was only accomplished within the last generation as a direct result of a prize set up for the express purpose of promoting that technology. The Kremer Prize, established in 1959 by British engineer and industrialist Henry Kremer, provided for an award of £5000 for the first humanpowered vehicle to fly a specified half-mile figure-eight course. It was awarded in 1977, less than 20 years later, to a team headed by Paul Macready, Jr., for a flight by Bryan Allen in the Gossamer Condor.

The success of the Kremer Prize depended on two factors.

• *Pursuing a purpose:* The goals of the Kremer Prize were clear. At the time of the institution of the prize, there were no active efforts to build human-powered aircraft. The goal of the prize was to provide an incentive to enter the field of human-powered flight. It was tremendously successful at this goal. By the time that the *Gossamer Condor* made its award-winning flight, Macready's team was in competition with several other teams with planes that were flying substantial distances solely under human power. • *Pushing the envelope:* The basic sci-

⁸In fact, other limited Turing tests have been carried out as well. See the discussion by Moor [18, p. 1129–30] for some examples.

ences underlying human-powered flight were, by 1959, well understood. These included aerodynamics, mechanics, anatomy and physiology, and materials technology. It was even possible for Robert Graham, an expert in the field of human-powered flight and a founding member of the Cranfield Man-Powered Aircraft Committee, to state at that time that "Man could fly, if only someone would put up a prize for it" [14, p. 23]. Overcoming the human difficulties in building a team that had collective mastery of these various fields and the engineering difficulties in creatively combining them were astonishing accomplishments. Nonetheless, as it turned out, no new basic discoveries were required at the time of the founding of the Kremer Prize to win it.⁹ The task was just beyond the edge of the current technology. Unfortunately, since our ability to dream far outstrips our ability to build, the establishment of tests of ridiculous difficulty is not difficult to imagine. At a time when an awardwinning human-powered flight was one of one meter at an altitude of 10 centimeters (the 1912 Prix Peugeot), the Paris newspaper La Justice established a prize for the first nonstop human-powered flight from Paris to Versailles and back. (It was never won.)

The history of human-powered flight indicates that only when the purpose of the prize is clear and the task is just beyond the edge of current technology is a prize an appropriate incentive. The Kremer Prize is a prime example of a prize that meets these criteria. The Loebner Prize is not.¹⁰

We turn first to the goals of the Loebner Prize. It was, according to the formal statement in the competition application, "established...to further the scientific understanding of complex human behavior." Along these lines Loebner has been quoted as saying "People had been discussing the Turing Test; people had been discussing AI, but nobody was doing anything about it" [17]. (The several

⁹⁶ The flight [of the *Gossamer Condor*] has shown that, with what appears to be a comparatively unsophisticated design, controlled manpowered flight over a reasonable distance is possible" [22, p. 341].

It is questionable whether the notion of a Turing test limited in ways

specificed by the Loebner contest is even a coherent one.

thousand members of the American Association for Artificial Intelligence (AAAI) may be surprised to learn that nobody is doing anything about it.)

Others have argued that the prize will serve to publicize the Turing Test, thereby increasing the public's awareness and understanding of artificial intelligence. Increased public understanding of AI is certainly a laudable goal, especially since the regular appearance of superficial popularizations in the press serves more to mislead the public by alternately raising and dashing expectations than to inform it by cogent coverage of actual results. A flurry of the standard stories in the press such as "Computer fools half of human panel" [13] and "Test a breakthrough in artificial intelligence" [16] was certainly one of the side effects of the Loebner Prize competition, but perhaps not a laudable contribution.

Overselling of AI by the media (and, occasionally, practitioners¹¹) has, in its brief history, been a repeated and persistent problem, and the hubristic claims of the organizers of the Loebner Prize that they are "confident that within 10 to 20 years a system will pass this electronic litmus test" [27] perpetuates the hyperbole. Robert Epstein in his recent ar-

¹¹Dreyfus [8] provides pertinent examples.

ticle describing the event, its genesis, and his speculations as to its importance constructs a standard claim of this sort:

Thinking computers will be a new race, a sentient companion to our own. When a computer finally passes the Turing Test, will we have the right to turn it off? Who should get the prize money—the programmer or the computer? Can we say that such a machine is "self-aware"? Should we give it the right to vote? Should we give it the right to vote? Should it pay taxes? If you doubt the significance of these issues, consider the possibility that someday soon you will have to argue them with a computer [9].

Not surprisingly, the winner of the Loebner Prize has jumped on the publicity bandwagon by taking out an advertisement pushing his program as the "first to pass the Turing Test."¹² Conversely, a prize whose execution mistakenly convinces fellow scientists that little progress has been made in a quarter century does little to promote the field. In summary, there is a difference between publicity and increased public understanding. Events of this sort-and the Loebner competition has been no exception-tend to generate the former rather than the latter.

Dennett has hinted at a completely different goal for the Loebner Prize. "It is useful to have the demonstration of the particular foibles that human beings exhibit in 1991.... We won't learn much about AI from the Loebner Prize for a long time, but we will learn some nonnegligible things about social psychology, perhaps, in the meantime" (Dennett, personal communication). For instance, the competition might be justified "as a proving ground for the environmental conditions necessary to permit the

¹²Dennett has, on behalf of the Loebner Prize committee, demanded that the advertising claim be discontinued, at peril of lawsuit, and Weintraub has apparently complied.

Turing Test to someday occur. In other words, the Loebner competition can tell us what we need to know about how humans behave in computer-mediated interactions" (Dranoff, personal communication). This line of teleology for the Loebner prize, that it serves not as a test of the abilities of the computers but of the psychologies of the various participants, has often been proposed informally. Such a "conspiracy theory" of the prize as a vast psychology experiment executed on unwitting and unconsenting adults is as unlikely as it is disturbing. Of course, there is already an extensive literature on how humans behave in computer-mediated interactions, and the Loebner competition is not likely to contribute to it; it was not designed or executed as a controlled scientific experiment, and that was not its apparent intention, despite the hopes of Dennett and Dranoff that firm conclusions in psychology might be gleaned from it.

Thus, it is difficult to imagine a clear scientific goal that the Loebner Prize might satisfy. Turing's test as originally defined, on the other hand, had a clear goal: to serve as a sufficient condition for demonstrating that a human artifact exhibited intelligent behavior. Even this goal is lost in the Loebner Prize competition. By limiting the test, it no longer serves its original purpose (and arguably no purpose at all), as Turing's syllogism fails.¹³ It is questionable whether the notion of a Turing Test limited in the ways specified by the Loebner Prize

¹⁰Several other factors markedly differentiate the Kremer and Loebner Prizes. First, whereas the committee administering the Kremer Prize primarily consisted of scientists specializing in the engineering of human-powered aircraft, it has been observed that current researchers in artificial intelligence, computational linguistics, and natural language processing were conspicuous by their absence from the Loebner Prize committee. (This problem has since been corrected.) Second, competition for the Kremer Prize was on an as-needed, as opposed to regular, basis, and no prize was awarded until the test was completed in the presence of a judge certified by the committee. Finally, the successful participants in the human-powered flight competitions were uniformly groups with strong backgrounds in the component technologies. In the case of the Loebner Prize, the participants were almost without exception amateurs.

¹³Robert Epstein has claimed that "We have changed the Turing Test as Turing would have if he were alive [27]. But it seems likely that Turing would have appreciated that the limitations imposed on the test by the Loebner committee invalidate it as even a sufficient criterion for intelligent behavior, and would not have sanctioned such gross modifications. A reviewer notes that "none of the conditions assumed by Turing are redundant for a meaningful test not the unlimited domain, not the unlimited time, not the interactive nature of the test, not the interrogator's full awareness that one of the respondents *is* a machine."

committee is even a coherent one. The prize committee spent some time with the referees attempting to explicate the notion of "natural conversation without trickery or guile."

It was suggested that a criterion be used as to whether you might say the utterance in conversation with a stranger seated next to you on an airplane. For instance, what might a competition judge legitimately ask on the topic of Washington, D.C.? Certainly, the question, "Are there any zoos in Washington?" is the kind of thing you might ask a stranger when flying to the capital for the first time, whereas "Is Washington bigger than a breadbasket?" is just as certainly a trick question. What about "Is there much crime in Washington?" Undoubtedly acceptable. "Are there any dogs in Washington?" An odd question for an airplane conversation. "Are there many dogs in Washington?" Sounds better. "Are there many marmosets in Washington?" Odd. "Are there many marmosets in the Washington zoo?" Okay again. The explanation of such examples begins to sound like arguments about angels and sharp objects.

Similar problems accrue to the notion of limiting the topic of discourse. Is the last question about Washington, D.C. or marmosets? (One of the referees in fact thought this and similar questions should be ruled out, since it was not strictly on the topic of the city alone.) How about "Are the buildings in Washington very modern?" Perhaps a question about architecture, as the following question surely is: "Do you know any examples of neo-Georgian architecture in Washington?" Are culinary topics ruled out, as in "What foods is our nation's capital best known for?" Such issues are not idle in the context of the Loebner competition. Cynthia Clay, the Shakespeare expert, was asked why Gov. Mario Cuomo has been referred to recently as "Hamlet on the Hudson." The question caused much consternation among the referees peering over Clay's shoulder. Her response was "His brooding," after which she coolly changed the topic back to Shakespeare. Or had it ever left?

The reason that Turing chose natural language as the behavior definitional of human intelligence is because of its open-ended, freewheeling nature. "The question-andanswer method seems to be suitable for introducing almost any of the fields of human endeavor that we wish to include" [28, p. 435]. In attempting to limit the *task* of the contestants through limiting the *domain alone*, the prize committee succeeded in doing neither.

The distinction between domain and task is crucial. Finance is a domain, but not a task; withdrawing money from a bank account is a task, one that is achievable through both human and computer intermediaries these days; taking dictation of a funds transfer request is a task that only humans can currently undertake with reliability. Had Babbage limited his differential analyzer to multiply only even numbers, the design would have been no more successful. This is a limitation of domain that does not yield a concomitant limitation in task.

It is well understood in the field that natural language systems must be tested using a constrained task. Currently, standard limited tasks can be found in evaluation of natural language database retrieval systems (such as withdrawing money from a bank account on the basis of a natural language request) and speech recognition systems (such as transcribing a spoken funds transfer request). The tasks, typically undertaken with limited vocabulary, are easily quantifiable along several dimensions (for example, technical notions of precision, recall, overgeneration, perplexity) independent of the subjective judgments of lay judges. Additionally, they can be adjusted to sit just at the edge of technology, unlike the Turing Test itself. The natural language research community has used such tests for some time now, and there has been increased interest in issues of evaluation of systems (primarily at the behest of funding agencies) over the last few years; whole conferences have been devoted to the subject (see, for instance, the report by Neal and Walter [21]).¹⁴

In summary, the Loebner Prize

competition neither satisfies its own avowed goals, nor the original goals of Alan Turing. In fact, it is difficult to imagine a scientific goal encouraged more by the Loebner Prize than by other uses of Loebner's \$1,500, his \$100,000 promissory note, and the \$80,000 in ancillary grants from the National Science Foundation and the Sloan Foundation. (Nonscientific goals are much easier to imagine, of course.)

Now to the second criterion for an appropriate technology prize, that the task be just beyond the edge of technology. Imagine that a prize for humanpowered flight were set up when the basic science of the time was far too impoverished for such an enterprise, say, in da Vinci's era. The da Vinci prize, we shall imagine, is constituted in 1492 and is to be awarded to the highest human-powered flight. Like the Loebner Prize, a competition is held every year, and a prize must be awarded each time it is held. The first da Vinci competition is won by a clever fellow with big springs on his shoes. Since the next competition is only one year away (no time to invent the airfoil), the optimal strategy is universally observed by potential contestants to involve building a bigger pair of springs. Twenty-five years later, the head of the prize committee announces that little progress has been made in human-powered flight since the first round of the prize, since everyone is still manufacturing springs.¹⁵

Of course, a lot of progress had been made in human-powered flight in those 25 years. Da Vinci himself was studying human physiology and anatomy and the flight of birds, andalthough his own work directly on the topic of human-powered flight, ornithopter design, was essentially meritless beyond its decorative qualities—the apparently tangential work was, in the long run, pertinent to the technologies that would eventually enable the Gossamer Condor to be constructed. (See e.g., Gibbs-Smith [12].) However, over that period, and indeed at every point during the following four centuries, the kind of progress necessary to solve the problem was not directly observable at that

¹⁵Hubert Dreyfus [8, p. 100] has made a similar analogy of climbing trees to reach the moon.

¹⁴Although the limitations and evaluation methods may be more sophisticated, the use of such task-limited evaluations to guide scientific research may be no more beneficial.

Any prize based on a behavioral test must use a limited task and

domain so the technology envelope is pushed, not ignored.

time in incremental improvement in solutions to the problem, the kind of improvement that might be observable in an annual contest. Nonetheless, tremendous scientific progress was made between the fifteenth and twentieth centuries. The Gossamer Condor and the digital computer are two outgrowths of this progress.

The field of AI is in that kind of state.¹⁶ The problem, like the problem of human-powered flight in the Renaissance, is only addressed directly and dismissed as imminently solvable by those who underestimate its magnitude. Progress on restricted tasks in limited domains is well documented in the literature on applications of AI. But progress on the underlying science that has been made in the last 25 years, important though it is, is not of the type that allows incremental advantage to be demonstrated on the big problem, the fullblown Turing Test, and this should not be seen as a failing of a field addressing a problem of the scope and magnitude of human intelligence. (And like all scientific endeavors, a lot of time can be spent on fruitless avenues of attack; ELIZA, as a discipline for natural language processing, was such a fruitless avenue. It was quite fruitful in other areas, however, as cogently argued by Weizenbaum himself.) Indeed, one aspect of the progress made in research on natural language processing is the appreciation for its complexity, which led to the dearth of entrants from the AI community-the realization that time spent on winning the Loebner Prize is not time spent furthering the field.

Twenty-five years of progress in the fields associated with the Turing Test—artificial intelligence, computational linguistics, and natural language processing—cannot be summarized in a single program, but is captured in the many small results, some of which, some day, at an unpredictable time in the future, may lead to a dramatic demonstration of apparently intelligent artificial behavior. To expect more is hubris. What is needed is not more work on solving the Turing Test, as promoted by Loebner, but more work on the basic research issues involved in understanding intelligent behavior. The parlor games can be saved for later.

Alternatives to the Loebner Prize

Given that the Loebner Prize, as constituted, is at best a diversion of effort and attention and at worst a disparagement of the scientific community, what might a better alternative use of Loebner's largess be? The goal of furthering the scientific understanding of complex human behavior is no less laudable now than it was before the competition, but clearly, a direct assault on a valid test of intelligent behavior is out of the question for a long time; even the prize committee well appreciates that. Thus, any award or prize based on a behavioral test must use a limited task and domain, so that the envelope of technology is pushed, not ignored. The efforts of the Loebner Prize committee to design such a test have failed in that the test they developed rewards cheap tricks such as parrying and insertion of random typing errors. This is an (indubitably predictable) lesson of the 1991 Loebner Prize competition.

This problem is a general one: any behavioral test that is sufficiently constrained for our current technology must so limit the task and domain as to render the test scientifically uninteresting. Adjusting the particulars of the Loebner competition rules will not help. By way of example, many years of effort have gone into the design of the tests of natural language systems used at the annual DARPAsponsored Message Understanding Conferences. The trend among entrants over the last several conferences has been toward less and less sophisticated natural language processing techniques, concentrating instead on engineering tricks oriented to the exigencies of the restricted task-keyword spotting, template matching, and similar methods. In short, this is because such limited tests are better addressed in the near term by engineering (building bigger springs) than science (discovering the airfoil). Behavioral tests of intelligence are either too difficult for a prize or too rewarding of incidentals.

At this stage, objective behavioral tests must give way to subjective evaluative ones. A more appropriate way to reward novel, potentially breakthrough-inducing efforts toward the eventual goal of mimicking intelligent behavior would be to institute a prize for just such efforts, on the model of the Nobel Prizes, ACM's Turing Award, and similar subjectively determined awards. Rather than awarding lifelong achievement or past accomplishments, however, the prize could be awarded for particular discoveries, regardless of field, that the committee determined were of sufficient originality, import, and technical merit and that were deemed contributory to Turing's goal (even though they may provide no incremental edge in a current-day Turing test).

To avoid unquestioning obedience to AI conventional wisdom, the awards committee would include eminent thinkers from a wide range of related fields (much as the current Loebner Prize committee does) but to ensure technical fidelity, a nominating committee of researchers from the pertinent technical fields should verify purported results before passing them on for consideration. In order to prevent stopping the approval of the reconstructed Loebner

¹⁶Prize committee member Weizenbaum places the state of AI technology a bit later in his analogy with Newtonian physics [31, p. 199], Dreyfus a bit earlier in his analogy with alchemy [7]. Neither writer is, of course, sanguine about the prospects for progress in the coming centuries.

Prize, it would be awarded on an occasional basis, only when a sufficiently deserving new result, idea, or development presented itself. I am not ostentatious enough to provide examples I believe would be appropriate for such an award; I am sure the reader can imagine one or two.¹⁷

As the years passed, and the speculations of this Loebner Prize committee as documented in their past decisions began to prove perspicacious, the Loebner Prize might grow in stature to that of the highly sought prizes of other scientific areas, and so provide a tremendous motivation for innovative ideas in the quest for AI.

Postscript

The Second Annual Loebner Prize Competition was held at the Cambridge Center for Behavioral Studies on December 15, 1992. The number of computer entrants had decreased from six to three, with Joseph Weintraub's program, complete with the winning strategy from the previous year's competition, taking first prize once again, this time under the purported topic "men vs. women." Bigger springs had prevailed.

Acknowledgments

I am grateful to the many readers of earlier drafts of this article; Ned Block, Noam Chomsky, Jacques Cohen, Daniel Dennett, Susan Cole Dranoff, Barbara Grosz, Harry Lewis, David Mumford, Fernando Pereira, Jeff Rosenschein, and David

¹⁷It is interesting to compare the Loebner Prize with the Leibniz Award for automatic theorem proving, endowed in 1983 by the Fredkin Foundation and administered by Carnegie-Mellon University. Like the Loebner Prize, the Leibniz Award offers \$100,000 on the basis of an extremely difficult task; it is to be conferred on the occasion of the first major new mathematical theorem whose proof is found with essential contributions by automatic theorem proving.

However, there are important differences Awarding of the Leibniz Prize is at the discretion of the Committee on Automatic Theorem Proving of the American Mathematical Society; it is therefore a subjective test, as it must be to decide issues such as the suitability of the theorem that was proved. In the interim, until the Leibniz Prize is awarded, intermediate awards are occasionally (not annually) presented. The Milestone and Current Awards are conferred, respectively, for "foundational work in automatic theorem proving" and for "ongoing research that shows promise," again at the recommendation of the committee. The Current Award, as an award for present developments rather than past achievement, is therefore structured in much the same way as the present proposal.

Yarowsky. I have incorporated many of their thoughtful comments into the article, although the opinions presented here are my own, and should not be taken as necessarily representative of the previous readers' views.

References

- Arbib, M. A. More on computer models of psychopathic behavior. *Commun. ACM* 17, 9 (Sept. 1974), 543.
- Block, N. Psychologism and behaviorism. *Philos. Rev.* 90, 1 (1981), 5–43.
- Block, N. The computer model of the mind. In An Introduction to Cognitive Science III: Thinking, D. N. Osherson and E. E. Smith, Eds. MIT Press, Cambridge, Mass., 1990, pp. 147– 289.
- Colby, K. M. Modeling a paranoid mind. *Behav. Brain Sci.* 4, 4, (1981), 515–560.
- 5. Colby, K. M., Hilf, F. D., Weber, S., and Kraemer, M. C. Turing-like indistinguishability tests for the validation of a computer simulation of paranoid processes. *Artif. Intell. 3*, 1 (1972), 199–221.
- Dennett, D. C. Can machines think. In *How We Know*, Michael Shafto, Ed. Harper and Row, San Francisco, Calif., 1985, pp. 121–145.
- Dreyfus, H. L. Alchemy and artificial intelligence. P. 3244, The RAND Corporation, 1965.
- 8. Dreyfus, H. What Computers Can't Do: A Critique Of Artificial Reason. Harper and Row, New York, 1979.
- **9.** Epstein, R. The quest for the thinking computer. *AI Mag. 13*, 2, 1992, 81–91.
- Fodor, J. Psychological Explanation. Random House, New York, 1968, pp. 126–127.
- French, R. Subcognition and the limits of the Turing test. *Mind* 99, 393 (Jan. 1990), 53-65.
- 12. Gibbs-Smith, C. H. Leonardo da Vinci's Aeronautics. Her Majesty's Stationery Office, London, 1967.
- Gomes, L. Computer fools half of human panel. San Jose Mercury News, Nov. 9, 1991, 1A, 16A.
- 14. Grosser, M. Gossamer Odyssey. Michael Joseph, London, 1981.
- Gunderson, K. Mentality and Machines. Doubleday and Company, Inc., New York, 1971.
- **16.** Krasner, J. Experts try to tell man from machine. *Boston Herald* Nov. 9, 1991, 21, 23.
- 17. Lindquist, C. Quest for machines that think. *Computerworld* (Nov. 18, 1991).
- Moor, J. H. Turing test. In Encyclopedia of Artificial Intelligence. Stuart C. Shapiro, Ed. John Wiley and Sons, New York, 1987, pp. 1126–1130.

- Moor, J. H. Explaining computer behavior. *Philos. Stud.* 34, 3 (1978), 325–327.
- Moor, J. H. An analysis of the Turing test. *Philos. Stud.* 30, 4 (1976), 249– 257.
- Neal, J. G. and Walter, S. M. Natural language processing systems evaluation workshop. Tech. Rep. RL-TR-91-362, Rome Lab., Griffiss Air Force Base, Rome, N.Y. 1991.
- 22. Reay, D. A. *The History of Man-Powered Flight*. Pergamon Press, Oxford, 1977.
- Searle, J. R. Can computers think. In Minds, Brains, and Science, Harvard University Press, Cambridge, Mass., 1984, pp. 28–41.
- 24. Searle, J. R. Minds, brains, and programs. *Behav. Brain Sci.* 3, 3 (1980), 417–457.
- **25.** Stalker, D. F. Why machines can't think: A reply to James Moor. *Philos. Stud. 34*, 3 (1978), 317–320.
- 26. Stipp, D. Some computers manage to fool people at game of imitating human beings. *Wall Street J.* Nov. 11, 1991, p. B3A.
- 27. The Guardian. Machines meet mastermind. Aug. 29, 1991.
- Turing, A. M. Computing machinery and intelligence. *Mind* 59, 236 (Oct. 1950), 433–460.
- **29.** Weizenbaum, J. Computer Power and Human Reason. W. H. Freeman, San Francisco, 1976.
- Weizenbaum, J. Reply to Arbib: More on computer models of psychopathic behavior. *Commun. ACM* 17, 9 (Sept. 1974), 543.
- Weizenbaum, J. Automating psychotherapy. Commun. ACM 17, 7 (July 1974), 425.

About the Author:

STUART M. SHIEBER is John L. Loeb Associate Professor of the Natural Sciences in the division of applied sciences at Harvard University. Current research interests include computational linguistics, automated graphic design, and combinatorial optimization. email: shieber@das. harvard.edu

Author's Present Address: Aiken Computation Laboratory, Division of Applied Sciences, Harvard University, Cambridge, MA 02138

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the ACM copyright notice and the title of publication and its date appear, and notice is given that copying is by permission of the Association for Computing Machinery. To copy otherwise, or to republish, requires a fee and/or specific permission.

[©] ACM 0002-0782/94/0600 \$3.50