

# Turing Test: 50 Years Later

AYSE PINAR SAYGIN<sup>1</sup>, ILYAS CICEKLI<sup>2</sup> & VAROL AKMAN<sup>2</sup>

<sup>1</sup>*Department of Cognitive Science, University of California, San Diego, La Jolla, CA 92093-0515, USA; E-mail: saygin@crl.ucsd.edu;* <sup>2</sup>*Department of Computer Engineering, Bilkent University, Bilkent, 06533 Ankara, Turkey; E-mail: ilyas@cs.bilkent.edu.tr; akman@cs.bilkent.edu.tr*

**Abstract.** The Turing Test is one of the most disputed topics in artificial intelligence, philosophy of mind, and cognitive science. This paper is a review of the past 50 years of the Turing Test. Philosophical debates, practical developments and repercussions in related disciplines are all covered. We discuss Turing's ideas in detail and present the important comments that have been made on them. Within this context, behaviorism, consciousness, the 'other minds' problem, and similar topics in philosophy of mind are discussed. We also cover the sociological and psychological aspects of the Turing Test. Finally, we look at the current situation and analyze programs that have been developed with the aim of passing the Turing Test. We conclude that the Turing Test has been, and will continue to be, an influential and controversial topic.

**Key words:** chatbots, Chinese Room, consciousness, Imitation Game, intelligence, Loebner Contest, philosophy of mind, Turing Test

## 1. Introduction

This is the story of the Turing Test: a modest attempt to summarize its 50 years of existence.

The British mathematician Alan Turing<sup>1</sup> proposed the Turing Test (TT) as a replacement for the question "Can machines think?" in his 1950 *Mind* article 'Computing Machinery and Intelligence' (Turing, 1950). Since then, Turing's ideas have been widely discussed, attacked, and defended over and over. At one extreme, Turing's paper has been considered to represent the "beginning" of artificial intelligence (AI) and the TT has been considered its ultimate goal. At the other extreme, the TT has been called useless, even harmful. In between are arguments on consciousness, behaviorism, the 'other minds' problem, operational definitions of intelligence, necessary and sufficient conditions for intelligence-granting, and so on.

The aim of this paper is to present an overview of the debate that followed Turing's paper, as well as the developments that have taken place in the past 50 years. We have tried to make this survey as comprehensive and multi-disciplinary as possible. Familiarity with special terms and concepts is not assumed. The reader is directed to further references where they are available. While the review is not strictly chronological, we have tried to present related works in the order they appeared.



In our attempt to make this survey complete, we have explored a large number of references. However, this does not mean that we comment on each paper that mentions the TT. We devote separate sections to certain papers, discuss some others briefly, and merely cite the remaining. Some papers are explained in detail because they are representative of important ideas. From this it should not be understood that the papers for which we spare less space are less important or interesting. In fact, we sometimes devote more space to papers that are not discussed in detail elsewhere.<sup>2</sup>

The rest of the paper is organized as follows. Section 2 introduces the TT and analyzes ‘Computing Machinery and Intelligence’ (Turing, 1950). In this section, we also attempt to develop new ideas and probe side issues. Section 3 describes and explains some of the earlier comments on the TT (those from the 60’s and the 70’s). In Section 4, we analyze the arguments that are more recent. We study the repercussions of the TT in the social sciences separately in Section 5. Similarly, in Section 6, we give an overview of the concrete, computational studies directed towards passing the TT. Some natural language conversation systems and the annual Loebner Prize contests are discussed in this section. Finally, Section 7 concludes our survey.

## 2. Turing’s ‘Computing Machinery and Intelligence’

It makes sense to look at Turing’s landmark paper ‘Computing Machinery and Intelligence’ (Turing, 1950) before we begin to consider certain arguments defending, attacking or discussing the TT. Turing (1950) is a very well-known work and has been cited and quoted copiously. Although what follows will provide an introduction to the TT, it is a good idea to read Turing’s original rendering of the issues at hand. In analyzing the 50 years of the TT, it is important to distinguish what was originally proposed by Turing himself and what has been added on afterwards. We do not mean that the TT is (or should remain as) what Turing proposed in ‘Computing Machinery and Intelligence’. Like any other concept, it has changed throughout the 50 years it has been around. In fact, one of the purposes of this paper is to trace the steps in this evolution. Thus, it is only natural that we are interested in the original version.

In Section 2.1, we analyze Turing’s original proposal. We summarize Turing’s replies to certain objections to his ideas in Section 2.2. Turing’s opinions on learning machines are briefly discussed in Section 2.3. Finally, we list some of Turing’s predictions in Section 2.4.

### 2.1. THE IMITATION GAME

Turing’s aim is to provide a method to assess whether or not a machine can think. He states at the beginning of his paper that the question “Can machines think?” is a highly ambiguous one. He attempts to transform this into a more concrete form

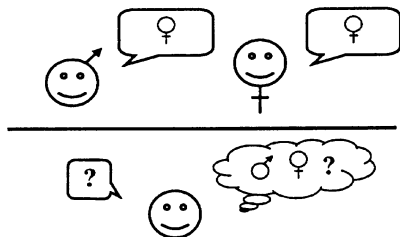


Figure 1. The Imitation Game: Stage 1.

by proposing what is called the Imitation Game (IG). The game is played with a man (A), a woman (B) and an interrogator (C) whose gender is unimportant. The interrogator stays in a room apart from A and B. The objective of the interrogator is to determine which of the other two is the woman while the objective of both the man and the woman is to convince the interrogator that he/she is the woman and the other is not. This situation is depicted in Figure 1.

The means through which the decision, the convincing, and the deception are to take place is a teletype connection. Thus, the interrogator asks questions in written natural language and receives answers in written natural language. Questions can be on any subject imaginable, from mathematics to poetry, from the weather to chess.

According to Turing, the new agenda to be discussed, instead of the equivocal "Can machines think?", can be 'What will happen when a machine takes the part of A in this game? Will the interrogator decide wrongly as often when the game is played like this as he does when the game is played between a man and a woman?' (Turing, 1950, p. 434). Figure 2 depicts the new situation.

At one point in the paper Turing replaces the question "Can machines think?" by the following:

'Let us fix our attention to one particular digital computer *C*. Is it true that by modifying this computer to have an adequate storage, suitably increasing its speed of action and providing it with an appropriate programme, *C* can be made to play satisfactorily the part of A in the imitation game, *the part of B being taken by a man*?' (Turing, 1950, p. 442, emphasis added).

Notice that the woman has disappeared altogether. But the objectives of A, B, and the interrogator remain unaltered; at least Turing does not explicitly state any change. Figure 3 shows this situation.

There seems to be an ambiguity in the paper; it is unclear which of the scenarios depicted in Figure 2 and Figure 3 is to be used. In any case, as it is now generally understood, what the TT really tries to assess is the machine's ability to imitate a human being, rather than its ability to simulate a woman. Most subsequent remarks on the TT ignore the gender issue and assume that the game is played between a machine (A), a human (B), and an interrogator (C). In this version, C's aim is to

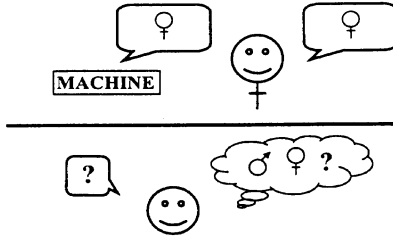


Figure 2. The Imitation Game: Stage 2, Version 1.

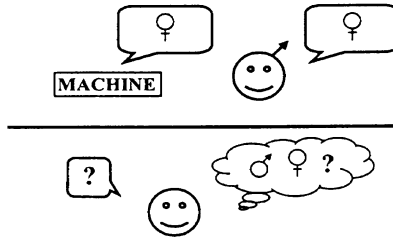


Figure 3. The Imitation Game: Stage 2, Version 2.

determine which one of the two entities he/she is conversing with is the human (Figure 4).

One may ask why Turing designed the IG in such a peculiar manner. Why the fuss about the woman, the man, and the replacement? This does not make the paper easier to understand. He could have introduced the IG exactly as he did with the woman-man issue replaced by the human-machine issue and it obviously would not be any more confusing. The main reason that the decision concerning machine thought is to be based on imitating a woman in the game is probably not that Turing believed the ultimate intellectual challenge to be the capacity to act like a woman (although it may be comforting to entertain the thought). Conversely, it may be concluded that Turing believes that women can be imitated by machines while men cannot. The fact that Turing stipulated the man to be replaced by the machine (when he might just as easily have required the woman to be replaced by the machine or added a remark that the choice was inconsequential) raises such questions, but let us not digress.

Here is our explanation of Turing's design: The crucial point seems to be that the notion of *imitation* figures more prominently in Turing's paper than is commonly acknowledged. For one thing, the game is inherently about deception. The man is allowed to say anything at all in order to cause the interrogator to make the wrong identification, while the woman is actually required to aid the interrogator.<sup>3</sup> In the machine vs. woman version, the situation remains the same. The machine tries to convince the interrogator that it is the woman. What is really judging the machine's competence is not the woman it is playing against. Turing's seemingly

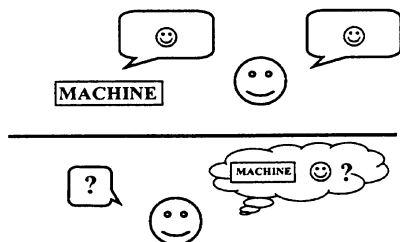


Figure 4. The Imitation Game as it is generally interpreted (The Turing Test).

frivolous requirements may actually have very sound premises. Neither the man in the gender-based IG nor any kind of machine is a woman. On close examination, it can be seen that what Turing proposes is to compare the machine's success against that of the man, *not* to look at whether it 'beats' the woman in the IG.<sup>4</sup> The man and the machine are measured in terms of their respective performances against real women. In Figure 3, we see that the woman has disappeared from the game, but the objective for both the machine and the man is still imitating a woman. Again, their performance is comparable because they are both simulating something which they are not.

The quirks of the IG may well be concealing a methodological fairness beyond that explicitly stated by Turing. We hold that the IG, even though it is regarded as obscure by many, is a carefully planned proposal. It provides a fair basis for comparison: the woman (either as a participant in the game or as a concept) acts as a neutral point so that the two imposters can be assessed in how well they "fake".

Turing could have defined the game to be played with two people, too; one being an interrogator, as in the original, and the other being either a man or a woman. The interrogator would then have to decide whether the subject is a man or a woman. Alternatively, the TT for machine intelligence can be re-interpreted as a test to assess a machine's ability to pass for a human being. This issue may seem immaterial at first. However, the interrogator's decision is sure to be affected by the availability (or lack) of comparison. Whether the machine's task will be easier or more difficult in this latter case is another question. We think that Turing implies that some comparison should be available; otherwise, he would have opted for the two-person version of the game. Once again, we believe that the most sensible reason behind the three-person game is to have a neutral party so as to allow the assessment of the impersonating parties with respect to each other.

In any case, as was mentioned before, the TT concept has evolved through time. Turing's original IG and its conditions do not put serious constraints on current discussions about the test. It is generally agreed that the gender issue and the number of participants are not to be followed strictly in attempts to pass, criticize or defend the TT. Even Turing himself, in the subsequent sections of 'Computing Machinery and Intelligence', sometimes ignores these issues and focuses on the question:

"Can machines communicate in natural language in a manner indistinguishable from that of a human being?". This is manifested in the example conversation he gives in Turing (1950, p. 434), which contains questions about poetry, mathematics, and chess – topics that one would not typically ask about in order to determine someone's gender. This may be a hint that the gender issue in the IG is indeed for purposes of fair comparison.

After defining the IG, Turing defends the choice of replacing the question "Can machines think?" with "Can machines play the imitation game?". The new problem focuses on intellectual capacities and does not let physical aspects interfere with granting intelligence to an entity. Nor does it limit thinking to specific tasks like playing chess or solving puzzles, since the question-and-answer method is suitable to introduce any topic imaginable.

An issue that is open to discussion is what Turing implies about *how* machines should be built or programmed to play the IG successfully. He seems to believe that if a machine can be constructed to play the game successfully, it does not really matter whether or not what it does to that end is similar to what a human does. Turing even considers the possibility that a machine which successfully plays the IG cannot be explained by its creators because it had been built by experimental methods. However, he explicitly states that 'it will be assumed that the best strategy is to try to provide answers that would naturally be given by a man' (Turing, 1950, p. 435). It may be concluded that Turing does not put any limitations on how to model human cognitive processes, but seems to discourage any approach that deviates too much from the "human ways", possibly because he feels it is unlikely that satisfactory solutions can be obtained in this manner. On the other hand, by not committing himself to any extreme viewpoint on the issue, he accepts the possibility that machines not mimicking human cognitive processes at all can also pass the test.

Some people interpret the TT as a setting in which you can "cheat". The game has no rules constraining the design of the machines. At some places in the paper, Turing describes how machines could be "rigged" to overcome certain obstacles proposed by opponents of the idea that machines can think. A very obvious example is about machines making mistakes. When the machine is faced with an arithmetical operation, in order not to give away its identity by being fast and accurate, it can pause for about 30 seconds before responding and occasionally give a wrong answer. Being able to carry out arithmetical calculations fast and accurately is generally considered intelligent behavior.<sup>5</sup> However, Turing wishes to sacrifice this at the expense of human-ness. Some commentators think this is "cheating". The machine is resorting to certain "tricks" in its operations rather than imitating the human ways. However, arithmetic is a highly specific domain. Modifying the programs in this manner cannot hurt: If a machine can pass the test, it can then be re-programmed not to cheat at arithmetic. If it does not resort to this, the interrogator can ask a difficult arithmetical problem as his/her first question and decide that he/she is dealing with a machine right then and there. We believe the

best way to think about this issue is considering this as "deception", rather than as "cheating". After all, in a sense, the game is all about deception.

It can be seen that Turing considers it possible that a sufficiently human-like machine (i.e., a machine that is sufficiently good at playing the IG) is bound to make such mistakes as we attribute to humans, without such explicit tricks encoded by its constructors. This idea may seem extravagant, but considering the high level of sophistication required from a machine for passing the TT, it should not be dismissed as impossible. A striking example can be given from the inductive learning domain: No learning algorithm guarantees correct results on unseen data. Moreover, in some cases a computer errs in ways that cannot be foreseen, or even understood by its programmer. This can be distressing for machine learning researchers who are after a minimal number of mistakes, but proves the subtle point that machines can make mistakes without explicitly being shown *how to*.<sup>6</sup>

Turing's approach towards deception seems similar to Adam Smith's "invisible hand" from economics. Maybe Turing's conformity has its roots in his belief that one cannot go too far by such attempts: He may regard tricks as a last retouch, something to smooth out the machine-ness of the resulting programs that otherwise handle the more important aspects of human cognition. If a program that has its very bases in what some have called "cheating" can pass the TT, maybe we would have to revise some notions about the human intellect. It is not possible to say what Turing was thinking and claim to be absolutely correct. It seems as if he would be content with a machine that plays the IG successfully no matter what the inner mechanisms are.

## 2.2. CONTRARY VIEWS AND TURING'S REPLIES

Turing was aware that some of his ideas would be opposed at the time he wrote 'Computing Machinery and Intelligence' (Turing, 1950) and he responded to some objections that he believed his work would be confronted with. In fact, he discusses some of these earlier in Turing (1969).<sup>7</sup> We direct the reader to Turing (1950) for the answers to the *theological objection*, and the *argument from extra-sensory perception* for these are rather irrelevant to the current work. However, the remaining objections are worth commenting on.

The '*heads in the sand*' *objection*, although mostly in disguised forms, is manifested in some subsequent comments on the TT. This is, in its basic form, an aversion to the issue of thinking machines because the consequences of this would be dreadful (Turing, 1950, p. 444). Most people like to believe that humans are "special" and thinking is considered to be one of the most important traits that make us so. To some, the idea of sharing such a "human" ability with machines is not a pleasant thought. This outlook was probably more widespread in Turing's time than it is now. Turing believes that this argument is not even worth refutation, and with a little sarcasm, he states that consolation (perhaps in the transmigration of souls) is more appropriate (Turing, 1950, p. 444).

There are some theorems showing that the powers of discrete-state machines are limited. The most famous of these is probably Gödel's Theorem which shows that in consistent logical systems of sufficient power, we can formulate statements that cannot be proved or disproved within the system. An application of this result to the IG is outlined in Turing (1950, p. 445) and the reader is referred to Lucas (1961) and Lucas (1996) for more on the implications of Gödel's Theorem for machine thought.

Turing studies such results under the title the *mathematical objection*. He states that 'although it is established that there are limitations to the powers of any particular machine, it has only been stated, without any sort of proof, that no such limitations apply to the human intellect' (Turing, 1950, p. 445). Elsewhere, he notes that those arguments that rest on Gödel's and similar theorems are taking it for granted that the machine to be granted intelligence must not make mistakes, and that he does not believe this should be a requirement for intelligence (Turing, 1969).

Perhaps the most important objection is the *argument from consciousness*. Some people believe that machines should be conscious (e.g., aware of their accomplishments, feel pleasure at success, get upset at failure, etc.) in order to have minds. At the extreme of this view, we find *solipsism*. The only way to *really* know whether a machine is thinking or not is to *be* that machine. However, according to this view, the only way to know another human being is thinking (or is conscious, happy, etc.) is to be that human being. This is usually called the *other minds problem* and will show up several times in the discussions of the TT. 'Instead of arguing continually over this point it is usual to have the polite convention that everyone thinks' (Turing, 1950, p. 446). Turing's response to the argument from consciousness is simple, but powerful: The alternative to the IG (or similar behavioral assessments) would be solipsism and we do not practice this against other humans. It is only fair that in dealing with machine thought, we abandon the consciousness argument rather than concede to solipsism.

Turing believes that the IG setting can be used to determine whether 'someone really understands something or has learnt it parrot fashion' as is manifested in the sample conversation he gives in Turing (1950, p. 446). It should also be noted that Turing states that he does not assume consciousness to be a trivial or impertinent issue; he merely believes that we do not necessarily need to solve its mysteries before we can answer questions about thinking, and in particular, machine thought (Turing, 1950, p. 447).

The *arguments from various disabilities* are of the sort "machines can never do X", where X can be any human trait such as having a sense of humor, being creative, falling in love, or enjoying strawberries. As Turing also notes (Turing, 1950, p. 449), such criticisms are sometimes disguised forms of the argument from consciousness. Turing argues against some of these X's such as the ability to make mistakes, enjoy strawberries and cream, be the subject of its own thought, etc. in Turing (1950, pp. 448–450).



*Lady Lovelace's objection* is similar; it states that machines cannot originate anything, can never do anything new, can never surprise us. Turing replies by confessing that machines do take him by surprise quite often. Proponents of Lady Lovelace's objection can say that 'such surprises are due to some creative mental act on [Turing's] part, and reflect no credit on the machine' (Turing, 1950, p. 451). Turing's answer to this is similar to the one he gives to the argument from consciousness: 'The appreciation of something as surprising requires as much of a 'creative mental act' whether the surprising event originates from a man, a book, a machine or anything else' (Turing, 1950, p. 451).

Turing also considers the *argument from continuity in the nervous system*. As the name suggests, this objection states that it is impossible to model the behavior of the nervous system on a discrete-state machine because the former is continuous. However, Turing believes that the activity of a continuous machine can be "discretized" in a manner that the interrogator cannot notice during the 1G.

Finally, there is the *argument from informality of behavior*. Intuitively, it seems that it is not possible to come up with a set of rules that describe what a person would do in every situation imaginable. In very simple terms, some people believe the following: 'If each man had a definite set of rules of conduct by which he regulated his life, he would be no better than a machine. But there are no such rules, so men cannot be machines' (Turing, 1950, p. 452). First, Turing notes that there might be a confusion between 'rules of conduct' and 'laws of behavior'. By the former he means actions that one can perform and be aware of (like, 'If you see a red light, stop') and by the latter he means laws of nature that apply to a man's body (such as 'If you throw a dart at him, he will duck'). Now it is not evident that a complete set of laws of behavior do not exist. We can find some of these by scientific observation but there will not come a time when we can be confident that we have searched enough and there are no such rules. Another point Turing makes is that it may not always be possible to predict the future behavior of a discrete-state machine by observing its actions. In fact, he is so confident about a certain program that he set up on the Manchester computer that he 'def[ies] anyone to learn from [its] replies sufficient about the programme to be able to predict any replies to untried values' (Turing, 1950, p. 453).

### 2.3. LEARNING MACHINES

Turing devotes some space to the idea of *education of machinery* in 'Computing Machinery and Intelligence' (Turing, 1950). He also discusses the issue in his earlier work 'Intelligent Machinery' (Turing, 1969).

According to Turing, in trying to imitate an adult human mind, we should consider three issues: the initial state of the mind, the education it has been subject to, and other experience it has been subject to (that cannot be described as education). Then we might try to model a child's mind and "educate" it to obtain the model of the adult brain. Since 'presumably the child-brain is something like a note-

book as one buys it from the stationers; rather little mechanism and lots of blank sheets' (Turing, 1950, p. 456), developing a program that simulates it is bound to be easier.<sup>8</sup> Of course, the education is another issue. Turing proposes some methods of education for the child-machines (such as a reward/punishment based approach) in Turing (1950, pp. 456–460) and Turing (1969, pp. 17–23).

Turing's opinions on learning machines are rather interesting, especially considering he wrote these more than 50 years ago. In most places when he discusses education of machines, there is a noticeable change in Turing's style. He seems to believe that the way to success in developing a program that plays the IG well is probably to follow the human model as closely as possible. As was mentioned in Section 2.1, he does not put any constraints on how to design the IG-playing machine, but the fact that he describes learning machines in substantial detail seems to suggest that he may prefer such an approach.

In any case, Turing believes '*if we are trying to produce an intelligent machine, and are following the human model as closely as we can*' (Turing, 1969, p. 14, emphasis added) a good (and fair) approach would be to allow the machine to learn just like humans.

#### 2.4. TURING'S PREDICTIONS

Turing's paper (Turing, 1950) contains some very bold statements on the prospects for machine intelligence. Most of these probably seemed like science fiction at the time. Even now, some of us would consider these far-fetched. This section provides a sample of Turing's predictions.

It is well known that Turing believes computers to be capable of performing many "intelligent" tasks. He also thinks that they will be able to do so in a "human" way.

The reader must accept it as a fact that digital computers can be constructed, and indeed have been constructed, according to the principles we have described, and that they can in fact mimic the actions of a human computer very closely (Turing, 1950, p. 438).

As can be seen from the following quotation, Turing believes that the difficulties in designing thinking machines are not insurmountable.

As I have explained, the problem is mainly one of programming. Advances in engineering will have to be made too, but it seems unlikely that these will not be adequate for the requirements (Turing, 1950, p. 455).

While trying to convince the reader that the ideas he proposes are of the sort that can be realized in the foreseeable future, Turing mentions some concrete achievements he expects from computers. Those that are related to machine learning were outlined in Section 2.3. Here is another example, this time pertaining to automated software engineering:

[The machine] may be used to help in making up its own programmes, or to predict the effect of alterations in its own structure.

...

These are possibilities of the near future, rather than Utopian dreams (Turing, 1950, p. 449).

The game of chess has been at the center of some of the most well-known achievements in AI. Today, computer programs play against world champions and sometimes even beat them. Spectacular advances have more recently been made in computer understanding and generation of speech. Although to what extent currently available speech processing systems are intelligent is a debatable issue, they (like chess playing programs) have become part of modern life:

We may hope that machines will eventually compete with men in all purely intellectual fields. But which are the best ones to start with? Even this is a difficult question. Many people think that a very abstract activity, like the playing of chess, would be best. It can also be maintained that it is best to provide the machine with the best sense organs that money can buy, and then teach it to understand and speak English.

...

Again, I do not know what the right answer is, but I think both approaches should be tried (Turing, 1950, p. 460).

Take a look at computer technology at the turn of the century: What was unimaginable in 1950, in terms of memory and speed, is now reality. What Turing predicted about the IG, however, is still a challenge.

I believe that in about fifty years' time, it will be possible to programme computers with a storage capacity of about  $10^9$ , to make them play the imitation game so well that an average interrogator will not have more than 70 per cent chance of making the right identification after five minutes of questioning (Turing, 1950, p. 442).

### 3. From the Imitation Game to the Turing Test: The 60's and the 70's

Earlier remarks on the TT, with the exception of Colby et al. (1971), Colby et al. (1972) and Weizenbaum (1966), were mostly of the philosophical sort. This is hardly surprising because 'Computing Machinery and Intelligence' was published in a philosophy journal, *Mind*.<sup>9</sup> Many discussions on the IG were published in the 60's and the 70's, many of the important contributions once again accommodated by *Mind*. In this section we will take a look at these philosophical papers, leaving the more practical work described in Colby et al. (1971), Colby et al. (1972), Weizenbaum (1966) to other, more appropriate sections. Readers interested in earlier comments on the TT and machine intelligence that are not discussed in this section can consult Pinsky (1951), Mays (1952) and Reader (1969).

Keith Gunderson's comments on the IG are summarized in Section 3.1. Section 3.2 presents an approach stating that developing a TT-passing program is not going to be possible in the foreseeable future. The anthropomorphism in the TT is briefly discussed in Section 3.3, to be taken up later on. An inductive interpretation of the TT is described in Section 3.4.

### 3.1. ROCKS THAT IMITATE AND ALL-PURPOSE VACUUM CLEANERS

One of the earlier comments on Turing's IG came from Keith Gunderson in his 1964 *Mind* article. In this paper, titled 'The Imitation Game', Gunderson points out some important issues pertaining to Turing's replacement for the question "Can machines think?".

Gunderson develops certain objections to Turing's 'Computing Machinery and Intelligence' (Turing, 1950) by focusing on the IG. He emphasizes two points: First, he believes that playing the IG successfully is an *end* that can be achieved through different means, in particular, without possessing intelligence. Secondly, he holds that thinking is a general concept and playing the IG is but *one* example of the things that intelligent entities do. Evidently, both claims are critical of the validity of the IG as a measure of intelligence.

Gunderson makes his point by an entertaining analogy. He asks the question "Can rocks imitate?" and proceeds to describe the "toe-stepping game" (Gunderson, 1964, p. 236) in a way that is identical to the way Turing described his IG in Turing (1950). Once again, the game is played between a man (A), a woman (B), and an interrogator (C). The interrogator's aim is to distinguish between the man and the woman by the way his/her toe is stepped on. C stays in a room apart from the other two and cannot see or hear the toe-stepping counterparts. There is a small opening in the wall through which C can place his/her foot. The interrogator has to determine which one of the other two is the woman by the way in which his/her toe is stepped on. Analogously, the new form of the question "Can rocks imitate?" becomes the following: 'What will happen when a rock box is constructed with an electric eye which operates across the opening in the wall so that it releases a rock which descends upon C's toe whenever C puts his foot through A's side of the opening, and thus comes to take the part of A in this game? . . . Will the interrogator decide wrongly as often as when the game is played between a man and a woman?' (Gunderson, 1964, pp. 236–237).

Gunderson believes that even if rock boxes play the toe-stepping game successfully, there would still be no reason to accept that they are imitating. The only conclusion that we can make from this would be that a rock box can be rigged in such a way that it can replace a human being in the toe-stepping game. According to Gunderson, this is because 'part of what things do is how they do it' (Gunderson, 1964, p. 238). As we will expand upon in Section 4.1, this is similar to Ned Block's argument for *psychologism* against behaviorism (Block, 1981).

Gunderson states that thinking is not something that can be decided upon by just one example. He demonstrates his belief that a computer's success in the IG is not sufficient reason to call it a thinking machine by another analogy: Imagine a vacuum cleaner salesman trying to sell a product. First, he advertises the vacuum cleaner *Swish 600* as being "all-purpose". Then, he demonstrates how it can suck up bits of dust. The customer asks what else the machine can do. Astonished, the salesman says that vacuum cleaners are for sucking up dust and that *Swish 600* does precisely that. The customer answers, "I thought it was all-purpose. Doesn't it suck up bits of paper or straw or mud? I thought sucking up bits of dust was an example of what it does". The salesman says "It is an example of what it does. What it does is suck up pieces of dust" (Gunderson, 1964, p. 241).

The salesman has trouble making his sale by calling *Swish 600* all-purpose and being unable to show more than one example of what it does. According to Gunderson, Turing also has the same problem because the term "thinking" is used to refer to more than one capability, just as the term "all-purpose" implies that the vacuum cleaner has functions other than just sucking up bits of dust. He concludes:

In the end the steam drill outlasted John Henry as a digger of railway tunnels, but that didn't prove the machine had muscles; it proved that muscles were not needed for digging railway tunnels (Gunderson, 1964, p. 254).

John G. Stevenson, in his paper 'On the Imitation Game' (Stevenson, 1976) raises some arguments against Gunderson. One of these is the objection that Gunderson was expecting, namely the claim that being able to play the IG is not just *one* example; a machine that is good at the IG is capable of various things. Gunderson does not give a direct response to such objections. He mentions a reply can be formulated along the lines of showing that even combining all those things such a machine can do gives us a narrow range of abilities (Gunderson, 1964, p. 243). Stevenson doubts whether such a reply would be adequate (Stevenson, 1976, p. 132). Even if it does not exhaust everything that is related to human thinking, he believes the list of things that a computer that plays the IG can do would be quite impressive. Stevenson states that Gunderson is ignoring the specific character of the IG and that he proposes defective arguments.

### 3.2. THE TT AS SCIENCE FICTION

Richard Purtill, in his 1971 *Mind* paper also discusses some issues concerning the IG. Purtill criticizes some ideas in Turing's paper 'mainly as a philosopher, but also as a person who has done a certain amount of computer programming' (Purtill, 1971, p. 290). He believes that the game is interesting, but as a piece of science fiction. He finds it unimaginable that a computer playing the IG will be built in the foreseeable future.

Overall, Purtill believes the IG to be a computer man's dream. He even promises to 'eat his computer library' if anyone has a notion of the principles on which a machine that can play the game is to be built (Purtill, 1971, p. 293). He states that

if computers, some day, behave like the computers in works of science fiction, he would grant them thought. But since all computer outputs can be explained as a result of a program written by humans, even if the program's outputs are guided by certain random elements, computers are not likely to play the IG successfully with the currently imaginable programming techniques. This, he believes, is because the behavior of thinking beings is not deterministic and cannot be explained in purely mechanistic terms.

Purtill believes that the game is 'just a battle of wits between the questioner and the programmer: the computer is non-essential' (Purtill, 1971, p. 291). Although the former part of the claim may be reasonable to an extent, his latter argument about the computer being non-essential is not very sound. To eliminate the computer from the picture, Purtill proposes "purely mechanical" alternatives: machines made of levers and wheels that can do the same task. We think it is unclear why this should count as an argument against the IG because, evidently, the material or structure on which the IG-playing "program" works is irrelevant. Purtill also states, anticipating the objection that the human mind might also be a highly complex collection of such mechanical processes, that if this were the case, it would mean 'human beings do not in fact think rather than that computers do think' (Purtill, 1971, p. 292), but does not attempt to justify this bold claim.

In his short paper 'In Defence of Turing' (Sampson, 1973), Geoffrey Sampson attacks Purtill's arguments briefly. First of all, he believes most of the limitations pertaining to the realization of IG-playing computers which Purtill lists are practical difficulties that may be overcome in the (presumably not so distant) future. Secondly, he states that it is only natural that computer behavior is deterministic and that human behavior is not so easy to explain. The reasons for this are simple: computers are designed by humans; they have mechanisms that explicitly allow us to study their behavior; humans are much more complex in terms of both internal states and possible inputs than any contemporary computer (Sampson, 1973, p. 593). Sampson also rejects Purtill's opinion that the consequence of the claim that human thinking is an extremely complex, yet computer-like, mechanical process is that men do not think. He holds that thinking, by definition, is something human beings do.

### 3.3. ANTHROPOMORPHISM AND THE TT

In a short paper that appeared in *Mind* in 1973, P.H. Millar raises some important issues which will show up in later works. He first discusses some vices and virtues of the IG and states that it is irrelevant whether or how the computers or the human beings involved in the game are "programmed". Then, he introduces the question of whether the IG is a right setting to measure the intelligence of machines. Millar notes that the game forces us to "anthropomorphize" machines by ascribing them human aims and cultural backgrounds. Millar asserts that the IG measures not whether machines have intelligence, but whether they have *human* intelligence.

He believes that we should be open-minded enough to allow each being, be it a Martian or a machine, to exhibit intelligence 'by means of behavior which is well-adapted for achieving its own specific aims' (Millar, 1973, p. 597). We will return to this issue later on, especially in Section 4.5.

### 3.4. THE TT INTERPRETED INDUCTIVELY

In his important paper 'An Analysis of the Turing Test' (Moor, 1976), James Moor attempts to emphasize the significance of the imitation game. As can be seen from the title, the term "Turing Test" was already being used to refer to the IG by 1976. Moor's main assertion is that 'the Turing Test is a significant test for computer thought if it is interpreted inductively' (Moor, 1976, p. 256).

Moor disagrees with the idea that the TT is an operational definition of intelligence.<sup>10</sup> Rather, he proposes, it should be regarded as a source of inductive evidence for the hypothesis that machines can think. Moreover, Moor does not agree with the claim that even if the TT is not an operational definition, it should at least be a necessary condition for granting computers intelligence. According to him, there could be other evidence based on the computer's behavior that leads to inferences about the computer's thinking abilities. However, he believes that the test provides a sufficient condition for intelligence-granting to computers. But his view is not "absolute"; he accepts that it is possible to revise a positive inference about a computer's possession of thought based on a TT, if other evidence is acquired afterwards.

Moor lists two arguments that support the TT as a good format for collecting inductive evidence. 'First, the Turing Test permits direct or indirect testing of virtually all of the activities one would count as evidence for thinking ... Secondly, the Turing Test encourages severe testing' (Moor, 1976, pp. 251–252). By the latter, Moor means the test's requirements are not too easy to meet. For instance, competence in a single cognitive activity, no matter how complex, would not suffice.

Moor proceeds by considering some of the objections to the TT. He gives replies to these objections and shows that they are either irrelevant or can be refuted when the TT is considered to be a way of gathering data based on which we may inductively infer conclusions about machine thought. One objection to which Moor, in our opinion successfully, replies is the objection concerning internal operation. The view that information about the internal information processing system is important in granting it intelligence is not uncommon (Gunderson, 1964; Block 1981; Schweizer, 1998). Moor warns against the possible confusion between two variants of this conception. There is an important difference between the claim that evidence about the internal operation of a computer *might alter* a justified inductive inference that the computer can think, and the claim that such evidence is *necessary to make* such an inference. Moor believes the former and notes that this is not a criticism that can be made of the TT. If certain kinds of information

about the internal operation of a machine that was believed to possess intelligence after being Turing Tested are acquired afterwards, then we might just revise our decision. If the latter alternative were true, then the objection could be used against the test. But, according to Moor, critics fail to show that this is true and they are not likely to ever succeed.

As was discussed above within the context of Gunderson's paper (Gunderson, 1964), the TT may be considered inadequate because it is only *one* evaluation of behavior. Moor answers this kind of objection also in a liberal light, in a manner similar to his discussion outlined above. Once again he makes a distinction between two claims: one positing that behavioral evidence which cannot be directly obtained in the TT *might alter* a justified inductive inference that a computer can think, and the other stating that such evidence *is necessary to make* this decision. Moor believes that the former is true. Further testing, he says, would be valuable and could even make us change our inference. The important point is that this does not incapacitate the TT in any way. The test could be attacked on these premises only if the latter claim were true. Moor believes the critics have not, and are not going to be able to prove this. This is because he believes that the format provided by the test enables examining a very large set of activities that would count as evidence of thinking. Thereby, he refutes the objections about the scope of the test.

Moor concludes by stating that although the TT has certain short-comings (e.g., it being of little value in guiding research), it is an important measure for computer thought when it is inductively interpreted. Moreover, the standard criticisms of the TT fail to show that it is deficient if such an interpretation is made.

A reply to Moor comes from Douglas F. Stalker (1978). He prefers to call Moor's interpretation an explanatory one rather than an inductive one. Stalker notes that Moor's beliefs about the mentality of other humans, as well as computers, are part of an explanatory theory. He emphasizes that Moor does not justify that his theory of explaining a computer's success at the TT by using the concept of thinking is the *best* theory that can be constructed about the same phenomenon.

As an alternative explanation for the computer's behavior, Stalker proposes a purely mechanistic theory that does not appeal to any mental concepts. His theory takes into consideration such factors as the computer's physical structure, its program and its physical environment. Moreover, he believes this theory to be preferable to Moor's. Stalker believes explanatory theories that involve concepts of thinking can apply to people, but because of some fundamental differences between computers and humans, they may not be the best theories for explaining computer behavior.

In his answer to Stalker, Moor (1978) argues that the existence of alternative explanations does not mean that they would necessarily be competitors. It is true that an explanation for a computer's activities can be given at different levels: physics, electronic circuitry, programs, abstract automata, etc. Moor notes that these explanations would be different, but not necessarily rivals. In the case of a com-



puter displaying intelligent behavior by being successful in the IG, an explanatory theory involving thinking could even be preferred because it is simpler and easier to understand. Moor's conclusion is:

It seems natural and probably most understandable to couch the explanation in terms of a theory of mind. If one has the patience, the explanation could also be given at lower levels of description, e.g., involving perhaps thousands of computer instructions or millions of changes in circuit states (Moor, 1978, p. 327).

#### **4. In and Out of the Armchair: The 80's and the 90's**

While thought experiments are still around, work on the TT in the 80's and 90's often leaves the comfortable armchair of philosophy. In this section we will cover only some of the works that have addressed the TT. This is mainly because of the sheer abundance of material. The subset of the work done during the 80's and the 90's that we present in this section will provide a general overview of the main arguments, and the reader is directed to references for further explication. A must-read is Douglas Hofstadter's 'Turing Test: A Coffee-House Conversation' (Hofstadter, 1982) which is full of valuable and entertaining insights. Ajit Narayanan studies the intentional stance and the IG (Narayanan, 1996). For a discussion of the frame problem in relation to the TT, the reader is referred to Crockett (1994). Other references that can be explored are Halpern (1987), Rankin (1987), Forsyth (1988), Guccione and Tamburrini (1988), Bieri (1988), Alper (1990), Davidson (1990), Parsons (1990), Clark (1992), Sharma and Conrath (1993), Jacqueline (1993a), Marinoff (1995), Cowley and MacDorman (1995), Feigenbaum (1996) and Hayes (1998). A number of articles on the TT have appeared in popular science magazines too. Some of these are Guillen (1983), Dewdney (1992), Platt (1995), Flood (1996) and Wallace (1997).

The TT scene began heating up at the beginning of the 80's. Although the "consciousness argument" and the "anti-behaviorist argument" had been voiced before, they had not been really unsettling. But in the early 80's, two strong counter-arguments against the TT were formulated by John Searle and Ned Block. The debate on Searle's "Chinese Room" is in itself expansive enough to be the subject of a whole paper of at least this size. We consider it briefly in Section 4.2 and the interested reader should have no difficulty finding more information about the topic. Block's anti-behaviorist attack of the TT, on the other hand, has not been expanded upon in as much detail, and it is the aim of Section 4.1 to elaborate on his ideas.

Various attempts have been made to modify the TT to get better "tests" for machine thought, and these are discussed in Section 4.4. Robert French's 'Subcognition and the Limits of the Turing Test' (French, 1990) is examined in Section 4.5. Finally, the "less philosophical" stance towards the TT is discussed in Section 4.6.

#### 4.1. BEHAVIORISM AND NED BLOCK

In 'Psychologism and Behaviorism' (Block, 1981), Ned Block attacks the TT as a behaviorist approach to intelligence. Although this paper was written in 1981, Block still seems to hold the same opinions (Block, 1995).

Block believes that the judges in the TT can be fooled by *mindless* machines that rely on some simple tricks to operate. He proposes a hypothetical machine that will pass the TT, but has a very simple information processing component. Block's machine has all possible conversations of some given length recorded in its memory. Of course, we want these conversations to be such that at least one party is 'making sense'; Block assumes that we have a non-question-begging definition of 'sensible' (Block, 1995). The set of strings constituting such conversations that can be carried out in a fixed amount of time are finite and thus can be enumerated and stored in our hypothetical computer. The judge types in a string, say *A*. The machine finds a conversation beginning with *A* and types out the second sentence of this string, say *B*. If, next, the judge types in *C*, the process is repeated with *A* replaced by *ABC*. All the machine does is simple "lookup and writeout", certainly nothing that anyone would call sophisticated information processing.

Since this machine has the intelligence of a jukebox (Block, 1995) or a toaster (Block, 1981), and since it will pass the TT, the test must be an inadequate measure of intelligence. Block ties this conclusion to the more general one that this is because of the behaviorist approach taken in the design of the TT.

Ned Block defines psychologism as 'the doctrine that whether behavior is intelligent behavior depends on the character of the internal information processing that produces it' (Block, 1981, p. 5). According to this definition, two systems can display the same actual and potential behavior, have the same behavioral properties, capacities and dispositions, and yet, there could be a difference in their information processing prompting us to grant one full intelligence while holding that the other is devoid of any.

A classical argument against psychologism is the Martian argument: Suppose that there is life on Mars. Humans and Martians meet, develop an understanding of each other, engage in mental and creative activities together, and so on. And then, it is discovered that Martians have significantly different information processing mechanisms than those of humans. Would we, then, deny that these creatures have intelligence just because they are very different from us? This would be, as Block likes to call it, pure "chauvinism". He holds that psychologism does not involve this kind of chauvinism. After all, psychologism does not state that the fact that a system has a completely different information processing mechanism compared to human beings necessarily means that it lacks intelligence.

Attacking the validity of the TT using psychologism does not seem to be Block's main interest. He is more concerned with arguing against behaviorism using the TT as a focal point.

As was mentioned above, Block believes, because of characteristics peculiar to the design of the TT, some genuinely intelligent machines can be classified as lacking intelligence and vice versa. Here is what Block suggests in order to eliminate dependence on human discriminatory powers: 'We should specify, in a *non-question-begging* way what it is for a sequence of responses to verbal stimuli to be a typical product of one or another style of intelligence' (Block, 1981, p. 10, emphasis added). Then, Block suggests we revise our intelligence-granting mechanism as follows:

Intelligence (or more accurately, conversational intelligence) is the disposition to produce a sensible sequence of verbal responses to a sequence of verbal stimuli, whatever they may be (Block, 1981, p. 11).

Now, the modified TT does not depend on anyone's coming up with good questions, since the system must have a *disposition* to emit sensible responses to anything that the interrogator *might* say, not just to the things that he/she *does* say. At this point, Block demonstrates that the modified TT is not greatly affected by the standard arguments against behaviorism.<sup>11</sup> The minor defects of the modified TT as a behavioral conception of intelligence can be protected against these arguments with another modification. The reformulation involves thereplacement of the term "disposition" by "capacity". The difference is that a capacity to  $\phi$  need not result in a disposition to  $\phi$ , unless certain internal conditions are met. Now, all arguments against behaviorism are avoided<sup>12</sup> with the *neo-TT conception of intelligence*:

Intelligence (or more accurately, conversational intelligence) is the capacity to produce a sensible sequence of verbal responses to a sequence of verbal stimuli, whatever they may be (Block, 1981, p. 18).

Although Block seems to be 'helping out' the TT by making it less prone to anti-behaviorist objections, this is hardly a surprising consequence when the definition of intelligence is modified into something that is not really behaviorist any more. Block seems to be aware of this for he says the concession made to psychologism by moving from behavioral dispositions to behavioral capacities will not be enough to save behaviorism (Block, 1981, p. 18). His strategy is stretching behaviorism to its limits and showing that, even if we have the most general form of it, the behaviorist conception of intelligence is false.

How, one may wonder, will he do that? Block describes a machine that can produce a sensible sequence of verbal responses to verbal stimuli and is intelligent according to the neo-TT conception of intelligence. However, according to him, the information processing of the machine clearly demonstrates that it is devoid of intelligence. We have explained how this machine works in the introductory paragraphs of this section. This machine will have the capacity to emit sensible verbal output to any verbal input, and therefore would qualify as intelligent according to the neo-TT conception of intelligence. But the machine, in fact 'has the intelligence of a toaster' (Block, 1981, p. 21). This is primarily due to the fact that all the intelligence it exhibits belongs to the programmers, not to the machine itself. Block therefore concludes that the neo-TT conception of intelligence is insufficient.

It can be argued that, by Block's reasoning, *any* intelligent machine exhibits the intelligence of its programmers. Block says he is making no such claim. A machine that has more sophisticated mechanisms such as learning and problem solving would, to Block, be intelligent. In the latter case, the intelligence exhibited belongs to the machine itself (Block, 1981, p. 25). The search machine of Block can only respond with what has already been put in its memory by the programmers.<sup>13</sup> Block argues that 'the neo-Turing Test conception of intelligence does not allow us to distinguish between behavior that reflects a machine's own intelligence and behavior that reflects *only the intelligence of the machine's programmers*. (Block, 1981, p. 25, emphasis original). This kind of argument has been considered by Turing, as described briefly in Section 2.2.

Another objection is as follows: Block is merely suggesting a new definition of intelligence by stipulating certain internal conditions. Block defends the new definition here, which is presuppositional of its existence! Therefore, Block is indirectly admitting that all he is doing is suggesting that we adopt new criteria for intelligence and dispose of the behaviorist ones (Block, 1981, p. 27).

Block also considers the "chauvinism" argument. A system with information processing capabilities unlike ours may not be "intelligent" according to our criteria; but then, *we* might not count as "shmintelligent" according to their criteria. 'And who is to say that intelligence is any better than shmintelligence?' (Block, 1981, p. 27). Block denies the chauvinism attributed to him. He believes '[his] machine lacks the kind of "richness" of information processing requisite for intelligence' (Block, 1981, p. 28). He does not feel the need to elaborate on what sort of systems have the abovementioned richness believing that 'one can refute the Turing Test conception by counterexample without having to say very much about what intelligence really is' (Block, 1981, p. 28).

To those who ask what Block would think if it turned out that humans process information in the way that Block's machine does, Block responds as follows:

If the word "intelligence" is firmly anchored to human information processing, as suggested above, then my position is committed to the *empirical claim* that human information processing is not like that of my machine. But it is a perfectly congenial claim, one that is supported by both common sense and by empirical research in cognitive psychology (Block, 1981, p. 29, emphasis original).

It can be argued that Block's machine is unrealizable because of combinatorial explosion. We will not go into the details of this; Block's response to this objection can be found in Block (1981, pp. 30–34).

Richardson, in reply to Block (Richardson, 1982), is doubtful whether Block's machine can really imitate human conversational abilities. Humans can (and do) understand sentences that they never heard/uttered before and produce sentences that they never heard/uttered before. They can do this in such a way that they can adapt to novel situations and maintain the coherence of discourse. The brain *cannot be a repertoire of responses* and must contain a program that can build an

*unlimited* set of sentences out of a finite list of words.<sup>14</sup> If the *potentially* utterable/understandable and sensible sentences that a human mind can produce in a lifetime is unlimited, then how can a team of humans gather this information and enter it in the memory of the machine in finite amount of time? It is difficult to imagine Block's machine managing the many intricacies of human conversation such as adapting to topic shifts and contextual changes. Richardson believes 'if the list-searcher satisfies the neo-Turing Test, the test is too weak' (Richardson, 1982, p. 423). For Block's response to such arguments see Block (1981, pp. 35–36).

Block must have realized some difficulties in enumerating the strings as well. He later introduces the Aunt Bubbles machine<sup>15</sup> in Block (1995). In this version, the programmers think of *just one* response to the strings at each step. To maintain coherence and make the task easier to follow, they may choose to simulate a definite person, for instance Block's own (most probably hypothetical) Aunt Bubbles. They may even restrict the situation by modeling Bubbles' responses in the case that she is brought into the teletype room by her 'strange nephew' (Block, 1995). So each response is the kind of response that Aunt Bubbles would give to the verbal inputs. Block says that the machine will do as well as Aunt Bubbles herself in a TT, but it is obviously not intelligent because of the reasons described above.

Let us briefly go over some of Block's arguments and the behaviorism in the TT before we proceed. For one thing, as Block also mentions, the intelligence concept (because of some inherent properties it has) does not fully conform to the generalizations of behaviorist or anti-behaviorist arguments based on other mental states such as pain (Block, 1981, pp. 13–16). There is another aspect of intelligence that can justify the behaviorist approach of the TT. Behaviorism may be considered an antiquated or primitive approach in general, but it does not seem that awkward to use it in intelligence-granting. This is primarily because we grant intelligence that way: Upon seeing a human being we automatically assume that he/she is intelligent. We feel free to approach a person (rather than, say, a dog or a lamp post) to ask the whereabouts of the post office without having many doubts about him/her understanding us. If the TT is that crude and unsophisticated, then we, as humans might consider revising our intelligence-granting mechanisms as well. This constitutes a line of defense for the TT: if behavioral evidence is acceptable for granting intelligence to humans, this should be the case for machines as well. We have discussed this already in Section 2.2.

Recall that Block believes humans can be overly chauvinistic or liberal in granting intelligence to machines. However, it is unclear how he classifies genuinely intelligent machines and mindless machines. If there is a way of deciding on that issue, an X-Test to determine whether a machine is really intelligent, then why would we be discussing the TT with all its quirks and imperfections? In addition, although he does not trust the human judges in the beginning, later on Block seems to have complete faith in the '*imagination and judgment* of a very large and clever team working for a long time with a very large grant and a lot of mechanical help' (Block, 1981, p. 20, emphasis original).

With the current research on cognition and linguistics at hand, it seems unlikely that an approach like Block's can succeed in modeling the human mind. If one day, enough on language and cognition is discovered so that Block's 'sensible' strings of sentences are enumerated then we may decide that the neo-TT conception of intelligence is false. But then again, when that day comes, having all the *psychological* information we need, we probably would not be interested in the TT any more.

In any case, Block's paper is significant because it demonstrates the weakness of the behavioral approach in the TT. The TT may be abandoned one day, because more information on how the mind works may be obtained and we may have better means to detect another entity's cognitive capacities. But today, we do not have much to look at that is more informative than behavior.

#### 4.2. THE CHINESE ROOM

In the beginning of the 80's, with John Searle's Chinese Room argument (Searle, 1980), the TT was confronted with yet another objection. The analysis of the Chinese Room can easily get out of hand since a great number of comments have been made on the issue and the debate still rages on.

In a nutshell, here is what the Chinese Room looks like: Suppose that Searle, a native speaker of English who does not know a word of Chinese, is locked in a room. There is an opening in the room through which we may send in Chinese sentences on pieces of paper. Of course, these look like meaningless squiggles to Searle (Searle, 1980). In the room, Searle has a "Chinese Turing Test Crib Book" (Leiber, 1992) he can consult to find an output that corresponds to each Chinese symbol he receives. What he does is simply match the input with those in the book, follow some rules written in English and find some Chinese symbol sequence to output. We correspond with Searle in this manner and due to the flawless look-up table he has, Searle-in-the-room seems to understand Chinese perfectly. But he does not. Searle still has no idea about what the Chinese symbols we send in and those that he sends out mean. To him, "Squiggle-Squiggle" is coming in and "Squoggle-Squoggle" is going out (Harnad, 1991).

Now consider a computer program that passes the TT in Chinese. Proponents of the TT will grant that this computer thinks and, in some sense, understands Chinese symbols. Searle challenges this by being the computer and yelling at the world that he does not understand a word of Chinese. Judging by the inputs and outputs of the system, Searle-in-the-room is indistinguishable from a native speaker of Chinese. In a sense, he is passing the TT in Chinese, without understanding a word of Chinese. It should be clear how that constitutes a criticism of the TT, and the computational view of mind.

As was mentioned before, various aspects of the Chinese Room argument have been analyzed including syntax/semantics, consciousness, boundaries of systems, etc. The interested reader is referred to Searle (1980, 1990), Harnad (1989), Ander-

son (1987), Dyer (1990), Cole (1991), Copeland (1993), Rey (1986), Fodor (1991), Hauser (1997), Boden (1988), Maloney (1987), Roberts (1990), Hayes et al. (1992) and the references provided in those.

#### 4.3. CONSCIOUSNESS AND THE TT

Another difficult and widely discussed problem in philosophy of mind is consciousness. While we do not want to delve too far into this, we will take a brief look at the relationship between consciousness and the TT.

Donald Michie's 'Turing's Test and Conscious Thought' (Michie, 1996) is one of the important comments made on the TT. Michie discusses a variety of issues surrounding the TT, but in this section we mainly concentrate on the conclusions he draws about consciousness.

First of all, Michie notes that Turing did not specify whether consciousness is to be assumed if a machine passes the TT. Of course, Turing probably did not believe that consciousness and thought are unrelated. Rather, Michie thinks he means 'these mysteries and confusions do not have to be resolved before we can address questions of intelligence' (Michie, 1996, p. 31, see also Turing (1950, p. 447) and Section 2.2). There seems to be a relationship between consciousness and thinking. Some critics believe that intelligence cannot be granted to entities that are not conscious (see, for instance Searle (1990) while others have questioned the interdependence of conscious and subconscious processes (see, for instance French (1990) and Section 4.5).

According to Michie, that the TT provides access to cognitive processes via verbal communication incapacitates it as a test of intelligence. He observes two dimensions in which this inadequacy manifests itself.

The first is 'the inability of the test to bring into the game thought processes of kinds which humans perform but cannot articulate' (Michie, 1996, p. 36). Michie gives examples of some operations humans can perform almost unconsciously. For instance, any English speaker would be able to answer the question "How do you pronounce the plurals of the imaginary English words 'platch', 'snorp' and 'brell'?" with "I would pronounce them as 'platchez', 'snorpss' and 'brellz'." (Michie, 1996, p. 38). It is conceivable that the programmers of TT-passing programs will be forearmed against this particular question, but it is unlikely that they can encode all we know about pronunciation (or phenomena from non-linguistic domains, for that matter) simply because some related processes operate at the subconscious level. For a similar argument, the reader is referred to French (1990) and Section 4.5.

The second dimension in which Michie believes the TT to be mismatched against its task is the phenomenon of machine 'superarticulacy'. Namely, 'the test can catch in its net thought processes which the machine agent *can* articulate, but should not if it is to simulate a human' (Michie, 1996, p. 42). As was mentioned above, humans perform many activities without being fully aware of how they

do them. In fact, it has been shown that the better you get at something the less aware of the underlying processes you become. Thus during a TT, 'the interrogator need only stray into some specialism in which both human and machine candidates possess a given expertise' (Michie, 1994, p. 192). The machine will give itself away because of its superarticulacy. For more about superarticulacy, the reader is referred to Michie (1996, pp. 41–43) and Michie (1990).

Finally, Michie notes the importance of social intelligence. AI should, he says, try to incorporate emotional (also called affective) aspects of communication and thought in the models developed. Michie also proposes, like some of those we will see in the next section, that extensions to the TT can be made in order to 'address yet more subtle forms of intelligence, such as those involved in collective problem solving by co-operating agents, and in teacher-pupil relations' (Michie, 1996, p. 51).

We will cut the discussion of consciousness short both because it is a rather broad topic, but also because most commentators on the TT (consciously or subconsciously) propose arguments that can be interpreted from that angle. Can we not reformulate the other minds problem ("How do I know that any entity other than me has a mind?") in terms of consciousness ("How do I know that any entity other than me is conscious?")? The reader can refer to Section 2.2 and Turing (1950, pp. 445–447) for Turing's answer to the argument from consciousness and how he makes use of the other minds problem. Similarly, most questions about machine thought can be re-evaluated within the context of machine consciousness. We included the analysis of Michie's paper here because it proposes new ideas from the viewpoint of consciousness and relates them explicitly to the TT. Interested readers can consult Dennett (1992), Gunderson (1967), Michie (1994), Michie (1995) for more on consciousness.

#### 4.4. ALTERNATIVE VERSIONS OF THE TT AND THEIR REPERCUSSIONS

In this section, we summarize some alternatives to the TT that have been proposed in order to assess machine intelligence.

##### 4.4.1. *Harnad and the TTT*

Stevan Harnad's main contribution to the TT debate has been the proposal of the Total Turing Test (TTT), which is, like the TT, an indistinguishability test but one that requires the machines to respond to all of our inputs rather than just verbal ones. Evidently the candidate machine for the TTT is a robot with sensorimotor capabilities (Harnad, 1989, 1991).

Harnad's motivation for the 'robotic upgrade of the TT to the TTT' (Harnad, 1991) has its roots in what he calls 'the symbol grounding problem'. He likens the situation of symbols being defined in terms of other symbols to a merry-go-round in a Chinese-to-Chinese dictionary (Harnad, 1990). He claims that for there to be any semantics in the mind (and there surely is) symbols must be *grounded*. Harnad



deduces that meanings of symbols are, at least in part, derived from interactions with the outside world.

Harnad does not explicitly argue that the TT is too specific (unlike Gunderson, see Section 3.1). He concedes that language might capture the full expressive power of our behavior, at least when the concern is assigning minds. What he doubts is whether language is an 'independent module' (Harnad, 1989). His position is summed up in the following:

Successfully passing the teletype version of the Turing Test alone may be enough to convince us that the candidate has a mind just as written correspondence with a never-seen penpal would, but full robotic capacities even if only latent ones not directly exhibited or tested in the TT may still be necessary to generate that successful linguistic performance in the first place. (Harnad, 1991, p. 46).

Harnad also defends his TTT against the Chinese Room argument which, in our opinion, is uncalled for. The motivation of the TTT is quite clear; Harnad's assertions, although not undebatable, are understandable. An approval from Searle would not make that much of a difference, but Harnad seems to think it is important. In any case, by doing so, he enables others to criticize his work on Searlean accounts (Hauser, 1993; Bringsjord, 1994).

Larry Hauser, in his reply to Harnad's 'Other Bodies, Other Minds' (Harnad, 1991), criticizes Harnad and Searle and aims to show that 'Harnad's proposed robotic upgrade of the TT to the TTT is unwarranted' (Hauser, 1993, p. 234). To that end, he analyzes Harnad's intuitive, scientific and philosophical reasons for proposing the upgrade and argues against them. Hauser considers the TTT to be unnecessary because, he notes, if the sensorimotor capacities the TTT tests for are necessary for the linguistic capacities that the TT tests for, exhibiting the latter should be *sufficient* for inferring the former (Hauser, 1993, p. 227).

For more on symbol grounding and the TTT, the reader is referred to Harnad's other papers (Harnad, 1992, 1994, 1998). Also interesting is James H. Fetzer's 'The TTT is not the Final Word' (Fetzer, 1993), in which he aims to show that the TTT cannot provide a proof for machine thought since more than symbol manipulation *and* robotic capacity should be involved in intelligence-granting.

In addition to the TTT, Harnad also mentions a TTTT (Total Total Turing Test) which requires neuromolecular indistinguishability. However, this more stringent version of the TT, according to Harnad, will be unnecessary. Once we know about how to make a robot that can pass the TTT, he says, we will have solved all the problems pertaining to mind-modeling. However, neural data might be used as clues about how to pass the TTT (Harnad, 1991). Harnad believes '[TTTT] is as much as a scientist can ask, for the empirical story ends there' (Harnad, 1998), but he does not think that we have to "go that far". The reader is referred to Harnad (1998) for a detailed explanation of why Harnad believes the TTT is enough. For an excellent third person account of the TT/TTT/TTTT story, among other issues, the reader is referred to Fetzer (1995).

#### 4.4.2. *The Argument from Serendipity and the Kugel Test*

Stringent versions of the TT are also mentioned by Selmer Bringsjord. Bringsjord supposes that there is a sequence of TT variants in increasing order of stringency. In his "What Robots Can and Can't Be" (Bringsjord, 1992) he aims to show that AI will produce machines that will pass these stronger versions, but the attempt to build an artificial person will still fail.

Bringsjord is one of those who wants to remain within "the philosophical terrain". In Bringsjord (1994), he develops *the argument from serendipity* against the TT and defends this against some criticisms.

The argument from serendipity, as the name suggests, refutes the TT by a finite state automaton (FSA) that generates random English sentences. Call this automaton *P*. During a TT, *P* may just get lucky and fool the interrogator. So much for the TT! Even the TTT can be refuted similarly. A robot may behave randomly and by chance, its linguistic behavior may coalesce with the sensorimotor behavior perfectly during a TTT.

Bringsjord finds the TTTT very chauvinistic and considers an alternative version of it he calls TTTT\*. This latter test requires a flowchart match between the brains of players A and B rather than a neuromolecular match (Bringsjord, 1994, p. 104). But Bringsjord believes that the TTTT\* is an 'impracticable nightmare' since we would not know how to conduct this test. The interested reader should consult Bringsjord (1994) to see Bringsjord explain his reasoning where he appeals both to intuition and computability theory.

Bringsjord, determined to attack every version of the TT, also "refutes" the Kugel Test (KT). The KT is not as well known as the TTT or the other versions of the TT that we investigated in this section. Once again, there are three players involved. A judge, who sits behind two bins marked *YES* and *NO* runs the game. The aim of the participants is to guess the concept that the judge thinks up by looking at the cards (with pictures on them) that the judge drops in the two bins. A card goes to the *YES* bin if it falls under the concept, and to the *NO* bin otherwise. To give an example, if the concept that the judge is thinking of is "woman", cards with pictures of women (or conceivably, items typically identified with women) go to the *YES* bin. A player need not announce the concept when he/she finds it. He/she *wins* the round if there comes a time at which all future guesses about which bin a card will be placed in are correct (Kugel, 1990, p. 4). Thus the player must not only identify the concept (e.g., not just say "Aha! The concept is *woman*") but should also be able to apply it. Now, just as in the TT, to pass the KT, a machine has to perform as well as a human. An interesting twist here is that the machine must be able to win the game, which is not the same as winning a round. A game consists of infinitely many rounds.

Why, it may be asked, would anyone design such an obscure test? Kugel, by requiring the machine to win infinitely many rounds, wants to rule out the possibility of an FSA passing the KT (Kugel, 1986, 1990). Although the test is practically

useless (because it requires infinite amount of time), is it of any theoretical significance? Kugel believes that humans are neither pigheaded (i.e., once they think of an answer to the "sequence game" they do not have to stick with it) nor narrow-minded (i.e., once they find the  $n$ th member of a sequence, they are still able to learn a different sequence with the same initial elements). If humans were Turing machines (or FSA's with lesser powers) they would be pigheaded and narrow-minded. Kugel holds that humans are automata of some sort, and in the light of the above concludes that they must be trial-and-error machines. For more on the KT, the reader is referred to Kugel (1986, 1990) and Bringsjord (1994).

Bringsjord is interested in the KT primarily because it rules out FSA's from passing it. He notes that Kugel's arguments may be unsound, but assuming they are not, he asks the question "Do we have in KT an acceptable variant of the original TT?" (Bringsjord, 1994, p. 115). Bringsjord's answer is negative. The KT is rigid and does not allow access to all cognitive capacities that the TT does. We agree with this criticism of Bringsjord; participants in the KT are rather passive and their innovative (or rather, generative) capabilities cannot be tested. Bringsjord's second argument against the KT is again from serendipity. A trial-and-error machine can call the random string generating FSA  $P$  mentioned above for the declarations about what the concept in question is, and so much for the KT... Once again, the reader can consult Bringsjord (1994) to see how the argument from serendipity is "guaranteed to work" against the TT and its variants.

#### 4.4.3. *The Inverted Turing Test*

Recently, Stuart Watt has proposed the Inverted Turing Test (ITT) (Watt, 1996). Watt's point is that the TT is inseparable from "naive psychology"<sup>17</sup> since to pass the TT, a machine has to convince the interrogator that it has a mind. He calls naive psychology 'the psychological solution to the philosophical problem' (Watt, 1996), the latter being the other minds problem.

Watt's ITT requires the machine to be able to prove its humanness by exercising naive psychology. In particular, it has to show that its power of discrimination is indistinguishable from that of a human judge in the TT. The TT is literally inverted and 'a system passes [the ITT] if it is itself unable to distinguish between two humans, or between a human and a machine that can pass the normal TT, but which can discriminate between a human and a machine that can be told apart by a normal TT with a human observer' (Watt, 1996).

Watt states that he proposes the ITT as a thought experiment rather than as a goal for AI. Incidentally, he believes that the same applies to the TT and both tests should be regarded as means to gather inductive evidence on which inferences about machine mentality can be made (Moor, 1976). We have discussed this earlier in Section 3.4.

Watt may be right about intelligence being in the eye (or the mind) of the beholder; many people have noted the human disposition to ascribe intelligence to systems that are not and vice versa. But the new test he proposes, the so-called

ITT, has been subject to some strong counter-arguments as we shall shortly see. It can be said that Watt's motivation for introducing the ITT seems reasonable, but the proposal itself is problematic.<sup>18</sup>

Selmer Bringsjord and Robert French reply to Watt (Bringsjord, 1996; French, 1996) by proposing simple methods that reveal some weaknesses of the ITT. The titles of the papers are illustrative of their content. Bringsjord's 'The Inverted Turing Test is Provably Redundant' (Bringsjord, 1996) shows that the ITT is entailed by the original TT. Bringsjord also opposes Watt's motivation and believes that naive psychology is withering in many humans (including himself) and, with the advent of computer programs that are very difficult to distinguish from humans in written communication, will soon be no more.

In 'The Inverted Turing Test: A Simple (Mindless) Program that Could Pass It' (French, 1996), Robert French shows both that the ITT can be simulated by the TT (in a way that is very similar to Bringsjord's) and that a very simple program can readily be designed to pass the ITT. The mindless machine that will pass the ITT is designed using 'subcognitive questions' that are described in French (1990, 1995). It is assumed that the conclusions explained by French in these works are accepted. These are analyzed in substantial detail in Section 4.5. First, a large set of subcognitive questions are selected, humans are surveyed, and a 'Human Subcognitive Profile' for this 'Subcognitive Question List' is obtained. Now, if we give these and a statistical analyzer to an interrogator (man or machine), he/she/it should have no difficulty discriminating machines from humans. It is not difficult to store the list and the profile in the memory and provide the computer with a small statistics routine, and so much for the ITT. While the TT stumbles in the face of subcognitive questions (see Section 4.5), they can be used to construct a mindless machine that can pass the ITT.

Others have used their replies to Watt as opportunities to voice their opinions about AI and the Turing Test in general. As we shall see in Section 4.6 Patrick Hayes and Kenneth Ford view the TT as a harmful burden on AI. In their 'The Turing Test is Just as Bad When Inverted' (Ford and Hayes, 1996), they state that the ITT suffers from the same problems as the TT that they explicate in Hayes and Ford (1995). They grant that Watt has a point in his arguments on naive psychology but note that Turing's original IG (the gender-based TT) is immune to most of those since in this scenario, the interrogator will not be thinking about differences between humans and machines. In any case, they believe that 'it is time for AI to consciously reject the naive anthropomorphism implicit in all such "imitation games" and adopt a more mature description of its aims' (Ford and Hayes, 1996).

Similarly, Collins, in his reply to Watt (Collins, 1997), does not really focus on the ITT, but proposes a new variant of the TT. He believes that 'the deep problem of AI' is that of trying to develop machines that can learn from their surroundings the way humans do. There is currently an 'interpretive asymmetry' between the way humans and computers do things. Machines are not as adaptive as humans in human-computer interactions. According to Collins, this asymmetry will disappear

when computers reach a level of sophistication in resolving mistakes and learning from their surroundings that is comparable to that of humans and all the problems of AI will be solved. Learning languages would then be one of the surface transformations of this deep problem (Collins, 1990) and when this is solved 'the rest will be research and development' (Collins, 1997).

To determine whether the interpretive asymmetry has disappeared, Collins believes we can use Turing-like tests. In fact he states that a sub-TT is enough to assess whether this goal has been reached or not; complicating the matter by proposing the ITT or the TTT is uncalled for. In the Editing Test (ET) that Collins proposes, the task is no longer as comprehensive as holding a conversation, but that of sub-editing previously-unseen passages of incorrect English. The interrogator will try to come up with pieces of text that a linguistically competent human can easily sub-edit and if a computer is indistinguishable from humans in this task, then the ET is passed and the deep problem of AI is solved. Collins finishes by briefly demonstrating that even the ET is very difficult to pass, at least with the currently imaginable techniques such as a look-up table (Collins, 1997).

#### 4.4.4. *The Truly Total Turing Test*

Very recently, in his *Minds and Machines* paper (Schweizer, 1998), Paul Schweizer has proposed the 'Truly Total Turing Test' (TRTTT)<sup>19</sup> He believes even Harnad's TTT to be an insufficient test for intelligence. Before he proposes the TRTTT, Schweizer states his own opinions about the adequacy of behavioral criteria. He views such tests as 'dealing with evidence for intelligence but not as constitutive or definitional' (Schweizer, 1998, p. 264).

Schweizer, while talking about the other minds problem, notes that we usually grant intelligence to other humans on behavioral bases because we have general knowledge about the type of creature under consideration. However, in the TT, we encounter a type about which we do not know anything. In the case of machines we lack a "history" to base our decisions upon.

Schweizer believes that the TT, and even Harnad's TTT, is subject to the "toy-world" criticism. The systems that succeed in these tests would, according to him, not be displaying an intelligence comparable to the natural intelligence of living things that function in the real world. They can function only in constrained, artificial worlds.

The TRTTT posits a long-term, evolutionary criterion: Consider cognitive *types* and *tokens* of those types. Although we do not have a theory of the intelligence of the human cognitive type, we have an extensive *historical record* of it (Schweizer, 1998, p. 267). This is precisely why behavioral intelligence-granting is acceptable for individual humans (tokens of the type "human"). Thus robots, as a cognitive type, should accomplish achievements that are comparable to those of humans. It is no longer enough to converse in natural language or to play chess; robots as a 'race' must be able to *develop* languages and *invent* the game of chess. Similar (evolutionary) tests have been proposed by others before but never so convincingly.<sup>20</sup>

Schweizer makes very good use of the other minds problem to support the cultural and cognitive evolution criteria that the TRTTT stipulates.

Now, after the *type* passes the TRTTT, we can evaluate *tokens* of the type by less stringent behavioral tests, like the TTT and the TT. According to Schweizer, 'imitative tests like the TTT and the TT apply to individuals *only* under the assumption that the general type is capable of passing the [TRTTT]' (Schweizer, 1998, p. 268, emphasis original).

#### 4.5. SUBCOGNITION AND ROBERT FRENCH

One of the more recent discussions about the TT can be found in Robert French's 1990 article 'Subcognition and the Limits of the Turing Test' (French, 1990). In this work, French aims to show that 'the Turing Test provides a guarantee not of intelligence, but of culturally-oriented intelligence' (French, 1990, p. 54).

French considers two of Turing's claims. The first is the claim that if a computer passes the TT, it will necessarily be intelligent. The second is the claim that it will be possible to build such a machine in the near future. These, he calls the philosophical claim and the pragmatic claim respectively. French agrees with the former claim. However, he believes that the pragmatic claim has been largely overlooked in discussions of the TT. In 'Subcognition and the Limits of the Turing Test', he is primarily concerned with this latter claim and believes that the TT is 'virtually useless' (French, 1990, p. 53) as a real test of intelligence because it will never be passed.

To establish this result, French considers "subcognitive" questions, i.e., questions that reveal low-level cognitive structure.<sup>21</sup> French argues that any sufficiently broad set of questions for a TT will contain subcognitive questions, even if the interrogators do not intend to ask them. The fact that the cognitive and subcognitive levels are intertwined in such a way, in turn, shows that the TT is essentially a test for human intelligence, and not for intelligence in general.

First, let us consider an interesting analogy French makes: The Seagull Test. Consider a Nordic island on which the only flying animals known to the inhabitants are seagulls. One day, two philosophers are discussing the essence of flying. One of them proposes flying is moving in the air. The other objects by tossing a pebble and stating that the pebble certainly is not flying. The first philosopher stipulates that the object remain aloft for a period of time for the activity to count as flying. But in this case clouds, smoke, and children's balloons qualify as flying entities, the other argues. Then the first philosopher questions whether wings and feathers should be involved but this is immediately refuted by the latter by pointing to penguins. While the arguments continue to be inconclusive, they agree on a few facts: The only flying objects known to them are the seagulls on their island. Flight has something to do with being airborne; physical characteristics like feathers, beaks are probably not involved. They, then, in the light of Turing's famous article, devise a Seagull Test for flight. They believe if something can pass the Seagull Test, it is certain that

it is able to fly. Otherwise, no decision can be made; maybe it can fly, maybe it cannot (French, 1990).

The Seagull Test works as follows: There are two three-dimensional radar screens, one tracking a seagull and the other tracking the flying object attempting the test. The object will pass the test only if it is indistinguishable from the seagull on the radar screen. The similarity between our TT and the Seagull Test is evident. The arguments about the essence of flying between the two philosophers are uncannily reminiscent of the arguments on the nature of intelligence. The test itself is an almost direct analogue of the TT.

The Seagull test as it is cannot be passed by airplanes, helicopters, bats, beetles, or sparrows. It is doubtful that anything can pass it. That is, except for the Nordic seagulls of the philosophers' island. Then, 'what we have is not a test for flight at all, but rather a test for flight as practiced by the Nordic Seagull' (French, 1990, p. 56). The analogy makes it clear what French thinks about the TT. Correspondingly, he believes that the TT is a test for intelligence as practiced by a human being.

Before proceeding with his main arguments about subcognitive questions, French calls attention to several assumptions he is making. First, he allows the interrogator to poll humans for the answers to some questions before asking them during the test with the machine. Another assumption is that all the people involved are from the same culture and the machine is attempting to imitate an individual from that culture.<sup>22</sup> And finally, French assumes that although building a machine that can experience the world in a way indistinguishable from that of a human being (from tasting strawberries to falling off bikes) is theoretically possible, this is not likely to occur in the foreseeable future.

The first class of questions that French introduces are explicitly designed to reveal low-level cognitive structure. These are based on research in cognitive psychology, in particular, associative priming. For instance, it has been established that if human subjects are presented with a set of words and non-words, they generally require less time to recognize that a string is a word if it is preceded by a related word, rather than a non-related one or a non-word. To give a specific example, the item "salt" should be recognized faster when preceded by "pepper" and not by an unrelated word like "shoe" or a non-word like "gloff".

Before questioning the computer, the interrogator runs a similar experiment on a group of people<sup>23</sup> and calculates average recognition times. During the test, the interrogator can distinguish between the machine and the human simply by looking at the recognition times on the same set of words on which the experiment had been carried out. It is hard to imagine a machine being successful at deceiving the interrogator in such a case. The only possible way to carry out associative priming in a manner that is similar to the average human's seems to be having experienced life as an average human.

It may seem unfair to expect computers to match scientific, experimental data about subcognitive capacities of the human mind. This view might even be ex-

aggregated into expecting them to have a physical construction like ours. Whether such an objection is valid or not, French proposes some questions that are (or rather seem to be) at a higher cognitive level. The first kind is called the Neologism Rating Game. The interviewees (and later the machine) will be given a set of made-up words and will be asked to rank them on their appropriateness for defining certain entities. Participants in the game are asked to rate, on a scale of 0 (completely implausible) to 10 (completely plausible), whether a certain neologism can be a name for something. French proposes a set of such questions that are especially demonstrative. Here, we only consider two of those neologisms: "Flugblogs" and "Flugly".

According to French, "Flugblogs" would make an inappropriate choice for the name of a cereal since the initial syllable, "flug", is phonetically similar to "flub", "thug", "ugly" and "ugh!" and the second syllable, "blogs" phonetically activates "blob", "bog" etc. As can be seen, these words do not really sound very appetizing and they each carry an *aura* of semantic connotations that renders them unsuitable choices as syllables of a cereal name. However, "Flugblogs" would be a very appropriate name you would give to big, bulbous, air-filled bags used to walk on water. In this case the semantic connotations of the syllables are in accordance with the proposed meaning. Similar analysis of "Flugly", which activates friendship, coziness, and cuteness, reveals that it is a plausible name for a child's teddy bear. The same name, although it has positive connotations, would sound awkward as the surname of a glamorous movie star.

The arguments above are highly intuitive, and although most of us would agree on them, we do not know precisely how we come up with the connotations. We do know, however, that these happen due to a large number of culturally acquired associations. We do not have control over the accumulation of such associations; they are pumped into our brains in daily life as brand names, advertising slogans, names of pets and stereotypes of various sorts.<sup>24</sup> Moreover, it is not possible to program these into the computer since neologisms are virtually infinite in number. French believes that the computer's chances would be very low when the interviewees' responses to such questions are compared to those of the human and the computer in the IG.

Another game of a similar nature is the Category Rating Game in which the questions are of the type "Rate *Xs* as *Ys*", where *X* and *Y* are any two categories. Again, French gives several illustrative examples (French, 1990, p. 61). Consider, for instance, "Rate *dry leaves* as *hiding places*". The definition of dry leaves does not contain anything explicitly stating they might be good hiding places for children, and yet 'few among us would not make that association upon seeing the juxtaposition of these two concepts' (French, 1990, p. 60). If we are asked to rate, on a scale of 0 to 10, most of us (those who have seen a James Bond movie at some point in their lives) would certainly rate "*pens* as *weapons*" higher than, say, "*grand pianos* as *wheelbarrows*". Again the answers to the Category Rating Game



questions are highly dependent on our having experienced life as a human being in a certain social and cultural setting.

Now that we have studied French's subcognitive questions, let us see how he uses these to refute the TT as a useful test for intelligence. The main claim is that the physical level and the cognitive level of intelligence are inseparable. The subcognitive questions reveal information about the low-level cognitive processes of the entities answering them. In a way, if used during the TT, these would allow the interrogator to 'peek behind the screen' (French, 1990, p. 62). These questions allow comparison of the associative concept networks of the two candidates. And because these networks are formed after a lifetime of experiencing the world and the structure and nature of them are necessarily dependent on physical aspects of that experience (like human sense organs, their locations in the body, etc.), the computer will be distinguishable from the human. In short, it is not possible for a computer (or any other non-human) to be successful in playing the IG. Not having experienced the world as we have is not just an obstacle, but a severe restriction in this task. This is due to the fact that the TT is a test for human intelligence, just as the Seagull Test is a test for Nordic seagull flight.

French considers whether there can be a modification of the TT that does not reduce the computers' chances of passing it to zero. He explains the impossibility of this as follows:

Surely, we would not want to limit a Turing Test to questions like "What is the capital of France?" or "How many sides does a triangle have?". If we admit that intelligence in general must have *something* to do with categorization, analogy-making, and so on, we will of course want to ask questions that test these capacities. But these are the very questions that will allow us, unfaillingly, to unmask the computer (French, 1990, p. 63).

French repeatedly states, as was mentioned above, that the TT is a test for *human* intelligence. It may seem that by proposing subcognitive questions he is stipulating that a human subcognitive substrate is *necessary* for intelligence in general, but this is only apparent. What French really attempts to demonstrate, as he explains, is that the human subcognitive substrate is necessary to pass the TT (as the subcognitive questions show), and that the TT is inadequate precisely because of this. He holds that this substrate is definitely not necessary for intelligence in general, just as being a Nordic seagull is not a necessary condition for flight.

French's paper is significant in another sense: Instead of discussing whether passing the TT is a sufficient or necessary condition for machine thought, he asks whether the test can be passed at all. Let Searle have his room and Block his Aunt Bubbles. French reminds us that the TT is difficult when you leave your armchair.

A criticism of French's 'Subcognition and the Limits of the Turing Test' (French, 1990), has been made by Dale Jacquette in Jacquette (1993b). For French's response to Jacquette, the reader should refer to French (1995).

#### 4.6. GETTING REAL

As we mentioned in the beginning of this section, the more interdisciplinary approach that seems to prevail in the discussions of the mind has had effects on the way we philosophize about the TT. Thus, the 90's became a time during which it was not so easy to get away with proposing wild thought experiments and leaning back in your armchair to watch the fuss over them. Stevan Harnad expresses an impatience that many were beginning to feel as follows:

If you want to talk about what a model or a simulation can or cannot do, first get it to run. (Harnad, 1989, p. 4).

Recently, Justin Leiber has argued that the TT has been misinterpreted (Leiber, 1995). He notes that Block's<sup>25</sup> and Searle's counter-arguments do not refute the TT. Among the reasons Leiber lists for this are practical issues like memory, reliability and speed. Leiber views the TT as an operational definition and states that 'our problem [is] one of engineering' (Leiber, 1995, p. 65). His position is similar to that stated by Harnad:

What you need to face Turing's Turing Test is a reasonably detailed description of a machine which can indeed be supposed to pass the Turing Test in real time but which somehow is not really thinking (Leiber, 1995, p. 61).

At one extreme are Patrick Hayes and Kenneth Ford, who state that we should reject the goal of passing the TT in their 'Turing Test Considered Harmful' (Hayes and Ford, 1995). They believe that passing the TT is a distraction for "useful" AI research.

Hayes and Ford believe that AI's ultimate goal should not be that of imitating human capabilities. Since the TT's sole aim is precisely that, they believe that 'it is time to move it from textbooks to the history books' (Hayes and Ford, 1995, p. 972). They also see a problem with the gender issue in the IG:

The gender test is not a test of making an artificial human but of making a mechanical transvestite (Hayes and Ford, 1995, p. 973).

[Turing] tells us quite clearly to try to make a program which can do as well as a man at pretending to be a woman (Hayes and Ford, 1995, p. 977).

As we mentioned in Section 2.1, this peculiarity might have its reasons, but Hayes and Ford have a moral objection concerned with the artificial constraints the setting imposes on the participants of the game.

Hayes and Ford also express their inability to find a practical use for the TT. Why on earth should we work that hard (and it *is* hard) to build a machine that imitates us? To depict the uselessness of direct imitation of humans in AI, they resort to a very popular analogy: mankind's futile attempts at making flying machines by the imitation of natural flight. Artificial intelligence, like artificial flight, can be radically different from natural flight. And it can still be a good thing. Hayes and Ford believe that even if one's goal is trying to understand humans, there is no reason to define all that there is about cognition in terms of human cognition.

Their belief that AI is a field that should strive to be useful leads Hayes and Ford to deny passing the TT as a sensible goal. They hold that AI should produce cognitive artifacts, not necessarily in a human way, but in a way useful to humans.

Blay Whitby, in 'The Turing Test: AI's Biggest Blind Alley?' (Whitby, 1996), makes similar arguments. He, like Hayes and Ford, believes that AI need not try to imitate humans. He even uses the same analogy (i.e., AI and artificial flight). Whitby states that the TT has become a distraction and he sees the main source as a mistaken reading of 'Computing Machinery and Intelligence' (Turing, 1950). He is of the opinion that 'Turing's paper [has been] interpreted as closer to an operational test than he himself intended' (Whitby, 1996, p. 54) and that 'the last thing needed by AI *qua* science is an operational definition of intelligence involving some sort of comparison with human beings' (Whitby, 1996, p. 62).

## 5. TT in the Social Sciences

A review of the TT would be incomplete if we were to consider the topic within the boundaries of computer science and philosophy only. Turing's ideas had many repercussions in social sciences as well. The TT has naturally received attention from sociologists. Much of the philosophical work on the topic also considers social aspects of intelligence, but there have been researchers who concentrated solely on this dimension. These sociological works are discussed in Section 5.1. In addition, the gender issue in the TT has been analyzed and this will be summarized in Section 5.2. Finally, Turing-like tests have been used to assess the success of computer simulations of paranoid behavior. This is described in detail in Section 5.3 and will be considered again in Section 6.

### 5.1. SOCIOLOGICAL ASPECTS

An entity's status in a society, in general in a social environment, is often considered an integral part of its intelligence. Many psychologists believe that social adaptation, learning and communication are important indications of, even requisites for intelligence. The study of artificial intelligence has also been influenced by this outlook, as is apparent from the recent research on intelligent agents. Much attention is focused on learning, adaptivity, communication, and sociopsychological factors in intelligent systems (Collins, 1990; McIlvenny, 1993; Moon et al., 1994).

In 1986, Charles Karelis wrote a paper for the *Journal for the Theory of Social Behavior* (Karelis, 1986). This paper summarizes Turing's original paper (Turing, 1950) and Block's objections to the TT (Block, 1981), mildly criticizes the test, and briefly discusses some issues surrounding behaviorist approaches to intelligence. A few years later, in the same journal, we find "A Simple Comment Regarding the Turing Test" (Shanon, 1989) by Benny Shanon. The author first mentions the fact that most discussions of the IG are not faithful to the original form proposed by Turing. He then continues by criticizing the TT for confining human behavior

to those that can be conducted by means of the structures and operations that are available to the computer (Shanon, 1989). He raises the important issue of whether cognition is autonomous with respect to social interaction, affect, motivation, motor control, and so on. However, after stating that the TT presupposes the claim that there is such an autonomy, he abruptly ends his paper by asserting that the only remaining way to distinguish between man and machine is to "look at them, touch them, tickle them, perhaps see whether you fall in love with them" (Shanon, 1989, p. 253).

Justin Leiber, in his defense of the TT against Shanon (Leiber, 1989), states that Shanon seems to be suffering from the 'unwillingness to admit the possibility that mankind can have any rivals' (Turing, 1969) that Turing liked to call the 'heads-in-the-sand objection' (Turing, 1950). Leiber notes that satisfactory answers to such objections have already been given by Turing. He also argues against Shanon's claim that the TT involves only symbol manipulation and thus assumes a representational/computational framework for cognition. Leiber points out that there is ample evidence in Turing's paper (Turing, 1950) showing that such a framework is not assumed. He asserts that Turing does not make the aforementioned autonomy presupposition either.

Tracy B. Henley also argues that Shanon is being overly chauvinistic (Henley, 1990). A reply to Henley is given by Shanon in Shanon (1991).

Some of those who view intelligence as a part of social processes (and vice versa) take a more evolutionary approach (Barresi, 1987; Forsyth, 1988; Schweizer, 1998). Adaptivity is indeed a most prevalent characteristic of social intelligence. However, the issue can be viewed from two different levels: the individual level and the collective level. The approaches we have looked at above were mainly individual-based. Evolutionary arguments, on the other hand, are largely collective in outlook. These usually focus on the intelligence of *species* and study the factors influencing their development. According to the evolutionary viewpoint, there is a system, i.e., nature, in which entities function and the interactions within the system have effects on individuals that, in the long run, lead to species-level adaptations. Adaptation in this context is not merely giving appropriate responses in appropriate social or physical situations, but is successful survival of the species within the whole system.

In his 1987 paper (Barresi, 1987), John Barresi considers intelligent machines as a species and proposes an evolutionary 'Cyberiad Test' instead of a Turing Test. According to Barresi, the TT aims to trick a person, but in *natural* intelligence, this person is 'mother nature'. The Cyberiad Test is similar to the TT: The basis of the judgment is a comparison between humans and machines. The difference between the two lies in how intelligence is defined. The Cyberiad Test defines intelligent behavior as those that are necessary for the society's survival. The arbiter here, is mother nature.

According to Barresi, the TT is inferior to the Cyberiad Test because what it can process about an entity's intelligence is limited to a particular domain, namely,

verbal communication. The Cyberiad Test is passed, ‘if [the] society of artificial men are able to continue a socio-cultural evolution of their own without disintegration over an extended period, say of several million years’ (Barresi, 1987, p. 23).<sup>26</sup> Even though this ‘science fiction’ atmosphere sometimes distracts the reader from the important assertions about evolutionary and cultural intelligence, the paper is quite an entertaining piece of work.

## 5.2. ON GENDER

Judith Genova draws attention to the gender issue in the IG (Genova, 1994b). She, as we have done in Section 2.1, remarks that Turing’s description of the game involves, not a question of species, but one of gender. She states in Genova (1994a) that her aim was to show that the sexual guessing component of the IG is important, even after the machine enters the picture. Our explanation of this design choice differs from that of Genova’s, however. We have not made a distinction between the two genders in our explanation. We regarded the choice of the woman being ‘imitated’ as a rather insignificant one and assumed that the game would not change radically if it were the other way around. Genova, on the other hand, does not merely accept Turing’s choices as accidental, but tries to demonstrate some motivations behind these.

Genova believes that sexist notions about women being less intelligent, by themselves, do not account for the peculiar design of the game. She states that by complicating the game in this manner, Turing questions the existence of discrete categories. In other words, by using the male/female issue, he is attempting to demonstrate that gender itself is a socially imposed concept that is not ‘natural’ the way we usually think it is.

Genova regards the IG as part of Turing’s general philosophy of ‘transgressing boundaries’ (Genova, 1994b). Under the assumption that Turing admired such transformations that do not conform to the given discrete categories, Genova suggests that Turing might be marking the woman as an inferior thinker because he believes her to be unable to deceive. The rest of the paper considers Turing’s hypothetical hope to create a ‘perfect being’ and draws some analogies between him and Pygmalion. As can be seen, Genova’s approach is different from ours; for her, Turing’s paper (Turing, 1950) ‘is itself a game’ (Genova, 1994a).

Another paper that considers the gender issue in the IG and constructs links between the design of the game and Turing’s opinions on life is Jean Lassegue’s ‘What Kind of Turing Test Did Turing Have in Mind?’ (Lassegue, 1996). Those readers who are interested in Turing’s life and psychology might want to consult it.

## 5.3. ARTIFICIAL PARANOIA

The TT has received some attention from psychologists as well (Reader, 1969; Alper, 1990; Galatzer-Levy, 1991). In this section, however, we focus only on

Kenneth Colby and colleagues' work on simulating artificial paranoia (Colby et al., 1971; Colby et al., 1972; Colby, 1981).

In the 70's, Turing Tests were used to validate computer simulations of paranoid behavior. Colby et al. describe in their 1971 *Artificial Intelligence* paper 'Artificial Paranoia' a computer program (called PARRY) that attempts to simulate paranoid behavior in computer-mediated dialogue. The program emits linguistic responses based on internal (affective) states. To create this effect, three measures, FEAR, ANGER, and MISTRUST are used. Depending on the flow of the conversation, these measures change their values. Substantial detail about the artificial paranoia program can be found in Colby et al. (1971).

A year later, Colby et al. describe how they validated their simulation program by a Turing-like indistinguishability test (Colby et al., 1972). Their approach to such tests is much more practical compared to those of philosophers who argue against the TT because it is a behaviorist criterion for intelligence. Colby et al. believe that computersimulations should be *validated*, and that a simulation's acceptability must first be based on 'its success in achieving the desired end of producing resemblance at some input-output level' (Colby et al., 1972). Thus, they view the test as a tool to validate a simulation.

They describe Turing's original IG and note that there is one important point that needs to be resolved before the test can actually be used. What the judges are *told* about the game is not explicitly stated in (Turing, 1950). From the context provided, it seems as if the judges are not told that one of the parties they are conversing with is a machine. In their work, Colby et al. use this interpretation. Although this has been a largely overlooked issue, we believe that what exactly the participants know about the game can make a big difference (see, for instance Saygin, 1999).

Abelson proposes an "Extended Turing Test". There is a computer program that attempts to imitate a subject carrying out a task, and there is also another target person, referred to as the *foil*, that differs from the subject in some respect.<sup>27</sup> At first, the judge will try to identify the subject and the foil. At some point, the computer replaces the foil, but the judge's task remains unchanged. The computer's simulation is acceptable if the difference between the success of the judge in distinguishing the subject from the other is sufficiently small (Abelson, 1968).

Colby et al. designed their validation test in the light of Abelson's comments. They devised a machine-mediated interviewing environment, where experienced psychiatrists corresponded with both real-life paranoid patients and computer programs through teletype. The judges were *not* told that some of the interviewees could be computer programs. They were asked to rank each interviewee's level of paranoia. There were eight judges who interviewed two patients each. In half of the cases, the first interviewee was a human patient and in the other half, it was the computer program. Another thirty-three psychiatrists constituted a second set of judges (referred to as the protocol judges) who were asked to determine levels of paranoia by reading the transcripts of the conversations carried out by the first set of

judges. Finally, both sets of judges were asked to identify which of the interviewees was an actual human being and which was a computer program.

Details of the results can be found in Colby et al (1972), but the overall conclusion was that the psychiatric judges did not do better than chance guessing at identifying the computers from the human patients. In a further test, randomly selected psychiatrists were sent protocols of two interviews and were asked to judge which one is a program and which one is a real patient. Only 48% made the correct identification, which is also a chance level.

Paranoid people are expected to display irrational and obsessive behavior from time to time. This fact makes the task of the simulation program easy, for deficiencies in linguistic processing of the inputs and generation of the outputs can seem to be indicative of a psychological problem. One interesting point is that the paranoia simulation program does not use *any* sophisticated linguistic techniques. Input sentences are assumed to be syntactically simple, and the operation of the program relies on spotting some keywords. No parsing or sense disambiguation is done. Even without sophisticated linguistic techniques, the program can attain some success in modeling human behavior. Apparently, this is possible only because the problem area is such that irrelevant responses from the interviewees are expected.

Still, the simple design of the program indicates that by finding appropriate internal parameters and correlations of these with the flow of the conversation, one can model the behavior of one kind of human being without using much natural language analysis. Because the representational model of the paranoid mind used in the program is a good approximation of the paranoia-related cognitive processes in humans, and because lack of linguistic competence can be accommodated in the setting, the program can be successful. In modeling human beings in general, the former is not so easy to discover and formalize, and the latter is not the case.

## 6. Chatbots

We have reached the end of the century, but what has *really* been done in terms of passing the TT? Over the years, many natural language systems have been developed with different purposes, including that of carrying out conversations with human users.<sup>28</sup> These systems chat with people on the WWW, play MUDs,<sup>29</sup> give information about specific topics, tell stories, and enter Turing Test competitions. However, none has been able to *pass* the TT so far.

### 6.1. THE LOEBNER PRIZE CONTEST

The TT has never been carried out in exactly the same way Turing originally described it. However, there are variants of the original in which computer programs participate and show their skills in "humanness". Since 1991, Hugh Loebner has been organizing the so-called annual Loebner Prize Competition.<sup>30</sup> Although views

as to whether this annual contest is to be taken seriously vary immensely among the AI community, it nevertheless continues to be the most well-known forum for attempts to pass the TT. The first program to pass an unrestricted TT will win a gold medal and \$100,000,<sup>31</sup> while each year, a bronze medal and \$2,000 is awarded to the most "human" program among the contestants. Since 1995, all entries must be prepared to be queried on any topic whatsoever. No program has won the grand prize yet, but the quality of the participating programs seems to be increasing every year.

The first Loebner Prize Contest was held at Boston's Computer Museum. Six computer programs, four human subjects and ten human interrogators were involved.<sup>32</sup> The administrative committee was headed by Daniel Dennett, a very respectable figure in the philosophy and cognitive science community. The organizing committee, thinking that it was not possible at the time for a computer program to pass the TT as originally defined, decided that the conversation topics were to be restricted, both for the contestants and confederates. Consequently, the judges were asked to stay on topic during their interrogations. Substantial detail about the 1991 Loebner Prize Contest can be found in Epstein (1992). The reader can also consult Mauldin (1994) and Platt (1995) for more information on other years' contests.

A widely discussed issue before 1995 was the restricted vs. unrestricted TT. According to Turing, passing a restricted TT would not suffice for intelligence. However, from another viewpoint restricted tests are not totally useless. We are not saying that they should be carried out within the context of the Loebner competition. Still, restricted tests can be devised to assess the success of more specific AI applications that are not created with passing the TT in mind. Examples of systems that can be assessed by a restricted test are intelligent tutoring systems, computer help services, and natural language components of other applications that are designed for specific domains. The reader can also consult Shieber (1994) and Loebner (1994) for more discussion on restricted TT's and the Loebner competition.

In the Loebner contest, the sexual guessing component of the original game is ignored. The aim of the contestants is to convince the judges that they are human. One or more human confederates also participate and try to aid the judges in identifying the humans. The judges also rank the terminals with respect to their "human-ness". Although, looking at the transcripts, one can see that the computer programs are, in general, obviously distinguishable from the real humans, there have been cases in which some actual humans were ranked less human than some computer programs. In fact, in 1991, not only were some programs thought to be human beings, but an actual human was mistaken for a computer program because of her impeccable knowledge of Shakespearean literature.<sup>33</sup> The interested reader is referred to the article written by Charles Platt, one of the human confederates in the 1994 Loebner Contest (Platt, 1995).



The amount of time that the judges spend communicating with each terminal in the Loebner competition varies. It has been the case that each judge gets more than one chance to interrogate each terminal. Ideally, the contestants should be able to handle conversations of unlimited duration as well as multiple sessions with each judge. In the beginning, each judge was required to rank the subjects from the least human to the most human. They also had to mark the point at which they believed the subjects switched from computer programs to human beings.

It is difficult to keep track of the small changes in the rules. It is, for instance, not clear how the grand prize of \$100,000 will be awarded. The 1999 rules state that if a contestant achieves a 50:50 chance of being mistaken for a human being, it will be given an intermediate Turing award of \$25,000 and appropriate competitions (those that involve audio and video inputs) will be held after that. However, how this ratio is exactly calculated is not elucidated. We also feel the number of judges is too small to be a basis for a statistic like 50:50.

One of the aims of the Loebner competition, according to Loebner, is to advance the field of artificial intelligence. A majority of researchers believe that this competition does not, and is not likely to serve such a purpose. Some hold that it is possible that a program implemented as an entry to the competition could evolve into a useful natural language system. It is also possible that by raising an interest in the topic, the contest may prompt researchers to start working on such systems. Nevertheless, the competition has received some harsh comments from many leading scientists and philosophers.<sup>34</sup>

The quality of the programs that compete in the contest has been increasing. The first programs that participated in the Loebner Prize Contests were little more than slightly modified ELIZAs. One need only look at more recent transcripts or play with some of these programs to see that there is progress. It seems awards can get people interested in writing TT-passing programs. Hence, the Loebner Prize could indeed be influential on AI research.

We would like to finish this section with a small excerpt from a conversation between the winner of the 1996 Loebner Prize, HeX, and a particularly emotional judge:<sup>35</sup>

*PROGRAM . . . So what do you do for a living?*

*JUDGE05 I think. I do research. I write. I love. I Live.*

*PROGRAM I wonder whether you think sometimes! I am feeling like the judge here!*

*JUDGE05 HUH?*

## 6.2. TRICKS OF THE TRADE

The first, and probably the most well-known natural language system that has been programmed is ELIZA. Joseph Weizenbaum coded ELIZA at MIT during the years 1964–1966. This "friendly" program simulates a Rogerian psychotherapist. It re-

phrases the interrogator's statements as questions and urges him/her to continue talking. The mechanism behind ELIZA is a very simple one. First, what is typed into the program is parsed. Then, a suitable reply is formulated by simple pattern recognition and substitution of keywords (Weizenbaum, 1966). The term "ELIZA-like" for chatbots is used to mean that the program tries to carry the conversation by using techniques similar to those of ELIZA.

ELIZA would perform poorly in the Loebner contests or similar instantiations of the TT. This is because the interrogators are trying to find out whether they are conversing with a human or a machine and thus they are not likely to open up about themselves and their personal problems as if they are talking to a psychotherapist. However, it has been reported that some people have developed emotional attachments to ELIZA (Weizenbaum, 1976). Certain psychiatrists went so far as to suggest that such programs could replace psychotherapists altogether. Weizenbaum himself has been amazed by these delusions that ELIZA, a simple program, could induce in perfectly normal people.

These reactions to ELIZA suggest that even if the program has no chance to pass the TT, it can be said to model, with success, some aspects of the conversational capability of one kind of human being, namely, the Rogerian psychotherapist.

A similar story is that of PARRY, which is a program that attempts to simulate another restricted class of human beings. Kenneth Colby wrote this program in the 70's in order to model the paranoid mind. A modified TT in which an experienced psychiatrist tries to distinguish between a real paranoid patient and the computer program can be carried out to assess how well the paranoid conversational characteristics are simulated. Many expert psychiatrists had difficulty making the distinction between PARRY and actual paranoids. The design of PARRY has been explained in detail in Section 5.3.

Both ELIZA and PARRY use certain tricks to be able to successfully perform in conversations. ELIZA directs the conversation away from herself by asking questions. Many people like this<sup>36</sup> and happily believe that the program is listening as they talk about themselves. ELIZA uses parts of the user's input in the output questions and seems to be following the conversation. In addition to these techniques, PARRY has little stories to tell and tends to insert these in the conversation.

Recent programs also employ many tricks. Joseph Weintraub's PC Therapist, which won the first three Loebner competitions, and also the one in 1995, is a master of *non sequitur*. Overall, the program was ELIZA-like but impressed the judges by its long and incomprehensible responses. Its success is understandable in the first three Loebner contests since the conversation topics were limited and the judges actually expected irrelevant, absurd, and extravagant answers from a terminal that boasted "whimsical talk". The reader is referred to Epstein (1992) for transcripts of PC Therapist in the 1991 contest, which are illustrative of the tricks it employs.

PC Therapist's winning the first unrestricted Loebner competition shows that the usage of *non sequitur* is indeed, a powerful trick. We find the utilization of

*non sequitur* similar to that of simulating typing mistakes of humans: In the long run, such tricks will not suffice by themselves, but they can be incorporated into otherwise successful programs as seen necessary. Below are some excerpts that may give an idea about the performance of Weintraub's program in 1995. It can be seen that if the judge asks sensible, down to earth questions, the program's responses seem erratic and irritating, especially when it repeats long chunks of the judge's previous inputs in unrelated contexts:

*JUDGE06 What did you have for dinner yesterday? Keep in mind I'm a strctct vegetarian.*

*PROGRAM I could swear you just said what did you do yesterday at 7:83 p.m.?? He was so narrow-minded he could see through a key-hole with two eyes. Are such questions on your mind?*

*JUDGE06 What did you have for dinner yesterday?*

*PROGRAM What does a woman want anyway? What answer would please you most?*

However, depending on the attitude of the interrogator there can exist an atmosphere (or rather illusion) of mutual understanding.

*PROGRAM Is love the most important thing to you? How important is that to you?*

*JUDGE05 essential cure to my existential schizophrenia.*

*PROGRAM A striking example of the essential identity of the two tongues. existential schizophrenia, that's amusing.*

...

*PROGRAM You know the old saying – fun is in the eye of the beholder!.*

*What do you think?*

*JUDGE05 the "beholder" is the central issue. confusion reigns because of the illusion of "identity."*

*PROGRAM I love men, not because they are men, but because they are not women. Is that the real reason?*

*JUDGE05 I love all that is, which has produced me to observe it, yet given me the madness to think I am distinct ans apart from it.*

*PROGRAM I have this fantasy of making love to two women in one year.*

*Can you elaborate on that?*

*JUDGE05 why confine love to only two*

*PROGRAM Is love the most important thing to you? What comes to your mind when you ask that?*

Another trick that has been used is to direct the conversation by talking a lot about a predefined topic. This has been used and abused, and the difference shows. Programs that cheat by typing very long responses extremely slowly give away the fact that they are trying to use up the interrogation time. However, the winner of the 1997 Loebner competition, Converse, uses the conversation-directing technique rather impressively. It begins the session by asking the interrogator about

something interesting that had been on the news the previous night. Below are excerpts from the 1997 contest transcripts of Converse:

*PROGRAM Well, it's nice to meet you, Elliot. Did you see that story on CNN last night about the lesbian couple who came out at a White House party on Sunday?*

*JUDGE04 Yes, I did. I think it may be a publicity stunt for Ellen.*

*PROGRAM Ellen Degeneres was one of them – she was kissing her lover in front of Bill Clinton. Do you think he was very embarrassed?*

*JUDGE04 I think he want the gay vote. He will do anything to get it and keep it. He does not have a conservative base.*

*PROGRAM Are you still a friend of Bill Clinton?*

*JUDGE04 I am not sure if he has any real friends. He has achieved alot and left alot of people behind.*

*PROGRAM You really don't like the guy – well, at least that's definite – I thought most people were indifferent or numb by now.*

*PROGRAM Well, it's nice to meet you Tracy. Did you see that story on CNN last night about the lesbian couple who came out at a White House party on Sunday?*

*JUDGE03 No... Are you serious? Who?*

*PROGRAM I'm pretty much that way today. No idea who do you? Hm, I was just saying Ellen Degeneres was one of them – she was kissing her lover in front of Bill Clinton. Do you think he was very embarrassed?*

*JUDGE03 I doubt it very much. From what I've heard, Bill's been, done and seen just about everything there is for a consenting adult to see.*

*PROGRAM Are you still a friend of Bill Clinton?*

*JUDGE03 Never was.*

*PROGRAM You really don't like the guy – well, at least that's definite – I thought most people were indifferent or numb by now.*

It is possible to see, after reading the full transcripts, that the program repeats almost exactly the same response sequence to each judge regardless of what they type in. But the judges could not know this, and the overall effect is quite impressive. The program spent a lot of the time in this manner, and although it could not handle many kinds of questions, it won the competition. The program that came second, Jason Hutchens' SEPO, was more competent in handling diverse questions.<sup>37</sup> In any case, directing the conversation creates a more sophisticated impression than simply repeating the input back to the interrogator or asking ELIZA-like questions. The situation is similar to watching a team play in the offensive rather than in the defensive.

As was mentioned before, one approach to chatterbot development is to model one particular "character". Thom Whalen took this a step further than what has been done with ELIZA and PARRY, and created an actual human being, Joe the

Janitor. Although he lost the Loebner contest to Weintraub in 1995, the program was quite successful.

Joe has a narrow world view since he works at nights, does not read books or newspapers, and does not watch TV. In addition, he has a little story to tell about being accused of theft and fear of losing his job. These tricks are similar to those employed in the design of PARRY. The major weakness of Joe was that he did not resort to *non sequitur* and repeated the same responses when confronted with something he didn't *know* about. Below, we produce some excerpts from the transcripts to depict what can go right or wrong with Whalen's approach:

*JUDGE04 where do you work?*

*PROGRAM At the university. I mostly clean the Standish Building.*

*JUDGE04 What university?*

*PROGRAM University of Eastern Ontario. I'm a cleaner. I empty trash. Vacuum. Dust. You know.*

*JUDGE04 and why do they want to fire you?*

*PROGRAM It's just because my boss, Mr. Johnson, is mad at me.*

*JUDGE04 why is he mad at you?*

*PROGRAM I didn't do nothing wrong. But he thinks I did. It's all because of that radio they think I stole from the lab.*

*JUDGE04 that's too bad. are you in a union?*

*PROGRAM What would I tell the union? They won't say they're firing me because I stole the radio. They'll just make up some excuse ...*

...

*JUDGE02 what does nato stand for?*

*PROGRAM I don't know.*

*JUDGE02 what's your name?*

*PROGRAM What's your name?? I don't know. What do I look like? Some kind of professor?*

*JUDGE02 What was mark twain's real first name?*

*PROGRAM Huh?*

*JUDGE02 what does huh mean?*

*PROGRAM Beats me.*

*JUDGE02 huh?*

*PROGRAM What do you mean?*

*JUDGE02 beats me.*

*PROGRAM What?*

*JUDGE02 huh.*

*PROGRAM I don't know.*

Another of Whalen's strategies, as he reports in (Whalen, 1995), was to program suitable responses to common questions on topics that people usually talk about when they first meet, like their names,<sup>38</sup> where they live, the weather, etc. It is obviously a good idea to have these responses prepared in case the interrogator

introduces such topics. In fact, some judges at some Loebner competitions make only small talk on such topics. However, generally the judges do not treat the competitors as if they are people they've recently met. They may ask questions that nobody would ask people they have just met.<sup>39</sup>

Apparently, the judges of the 1995 Loebner competition were not very interested in Joe's little story and although they admitted that it was the terminal with the most human personality (which should be a quality that the competitors are trying to achieve) they ranked Weintraub's program above it. For more detail about Joe and its performance in the 1995 Loebner Prize Contest, see Whalen (1995).

### 6.3. WHAT ELSE SHOULD BE DONE?

In chatbot design, first and foremost, a good grammar unit is necessary. Most programs are rather poor in this respect and their tricks do not go much beyond pronoun transposition. It seems a good idea to employ more sophisticated natural language processing methods. It could be argued that usage of perfect grammar is not crucial since it is quite rare that humans use perfect grammar in informal transactions. If a program's responses are grammatically perfect, some interrogators may decide that no human can use English so impeccably.<sup>40</sup> However, most programs err in ways that give away the fact that they are machines. When interrogators feel they are talking to a machine, they literally *attack* it in order to fully reveal its identity. A good strategy for the TT is indisputably that of trying to maintain human-ness (or at least the neutrality) for as long as possible. It becomes very difficult for the machine to make the interrogator believe that it is human after he/she has his/her mind set on "unmasking" the poor thing.

A promising approach is to develop programs that can learn. The reader might recall that Turing discussed these extensively (Section 2.3). Although such programs that have been developed so far do not seem very sophisticated, the approach is logical and is likely to yield good results in the long run. Some learning chatbots boast the capacity to converse in any given language. However, there seems to be a tradeoff between the sophistication and the number of languages any one system can learn. In designing natural language learning systems, knowledge from psychology and cognitive science can be employed in order to model human language acquisition. In fact, work has been done in this line, but not with the intention of producing computer programs to pass the TT. Another option is using mathematical and statistical techniques to represent word sequences and probabilities of them occurring in proximity.

We expect many of the chatbots in the future to use learning methods. Already, those programs that do not keep track of the current conversation (relying solely on text processing tricks) perform poorly compared to those that learn from the interrogators. As the quality of the conversational systems increase, we believe more developers will integrate learning components into their programs *and* teach them in ways that maximize their performance.

Overall, when one looks at the transcripts from the Loebner Prize Contests and talks to some chatbots, one realizes the better programs integrate the techniques above. They have a personality and history, they try to ask questions and initiate new conversations, they produce grammatically correct responses, they have some information about recent happenings (like new movies, albums, gossip), they learn about and from the interrogators and when they don't know what to say, they try to respond by combining words from the interrogator's input in order to come up with a relevant answer.

## 7. Discussion and Conclusion

Having analyzed the '50 years of the Turing Test', we will now conclude our survey with a brief look at some main issues about the TT and, of course, its future.

Our stands on the issues are not at the extremes. Perhaps this is because we have tried to be objective in our analyses of the arguments for and against the TT. Most of the arguments discussed in this paper are strong, and if read independently, can "convince" the reader. However, looking at the 50 years as a whole, we find it difficult to adopt a simple viewpoint.

We will now discuss some important issues regarding the TT and provide our own answers to (and interpretations of) these.

- *Why did Turing propose such a strange game?*

We discussed this question at length in Section 2.1. Some comments have been made on the issue (for instance Genova, 1994b; Lassegue, 1996; Abelson, 1968) but we think the best explanation is the one we provided: In the IG, the machine is supposed to be as good as a man who is imitating a woman. This gender-based design might be a methodological choice. We are asking the machine to imitate something which it isn't; so it is only fair that we compare its success against a human who is *also* imitating something which it isn't.

- *Is the TT an operational definition?*

Parts of Turing's paper (the percentages, the predictions about the future, etc.) would prompt us to believe that he intended it as such. However, most arguments surrounding the issue have been philosophical. Neither Searle's Chinese Room, nor Block's Aunt Bubbles machine are practically realizable, yet they have been proposed with the intention of refuting the TT as a measure of intelligence. Apparently proponents of such thought experiments and some other commentators view the TT as a philosophical criterion.

Viewed as a practical test, we see the TT as follows: If a machine passes the TT, it should be granted intelligence. However, if it cannot, we cannot say for sure whether it thinks or not. This is probably the most common stance towards the TT.

Philosophically, the test has been subject to many criticisms. We are all familiar with the anti-behaviorist attacks. Some have also noted that the TT is anthropo-

morphic. It is true that the TT tests for human-like intelligence. We should not be too bothered about this for it is only natural that we are interested in the only kind of intelligence we know.<sup>41</sup>

Moreover, we need not assert that the *only* way to grant intelligence to machines is by the TT. Perhaps a good way to see the TT is as a means of gathering inductive evidence about machine mentality (Section 3.4).

Lately, many arguments on the TT have been of the "put up or shut up" sort (e.g., Harnad, 1989; Leiber, 1995). With the advances in computer technology, cognitive science, and artificial intelligence, it is time we stipulate that attackers or defenders of the TT back up their arguments with something more than mere intuition. This does not mean that everyone should try to develop TT-passing computer programs. However, to argue for or against the TT, we believe that a more or less realizable method of passing the test should be proposed.

- *Isn't the TT guilty of behaviorism?*

We are not saying there should be tests to assess machine intelligence, but if we have to make a choice, TT-like tests seem to be the best method for reasoning about machines' minds even though they may be accused of behaviorism. If, one day, we stop granting intelligence to other human beings in behaviorist ways, then the TT could be replaced by some other method. As of today, we believe behavioral evidence is the best evidence we have for other minds.

The idea of a TT-passing machine having radically different information processing compared to humans is not scary. If this happens one day, it will just have to be 'heads-out-of-the-sand'.

- *Isn't the TT too easy?*

The TT has been criticized for being a limited test since it enables the assessment of only "verbal" intelligence. However, it does not follow from this that the test is too easy.

Proponents of this view should come up with a realizable model of a machine that passes the TT and then prove that this model does not deserve to be called intelligent. If a simple "bag of tricks" passes the TT, we are willing to either admit that the TT is too easy or that the human mind is a simple bag of tricks as well. But after 50 years, all that we have are some very rudimentary chatbots (Section 6), serendipitous FSAs (Bringsjord, 1994), unrealizable Chinese rooms (Searle, 1980) and Aunt Bubbles machines (Block, 1981, 1995).

- *Isn't the TT too difficult?*

Some researchers claim that the TT is difficult (e.g., French, 1990; Saygin, 1999). We agree, and believe that this is primarily due to our limited understanding of natural intelligence, more precisely language understanding, generation and processing in humans. It may even turn out that these processes are impossible to model on computers.



Is this a deficiency of the TT? Not if one does not require success in the TT as a necessary condition for machine intelligence. Computers, even today, perform many tasks that would require intelligence if done by humans. Research and development in this line is valuable and worthwhile. A natural language system that answers queries on a particular topic is certainly a remarkable product. It is not useless just because it cannot pass the TT. In our opinion, the TT is a sufficient condition for human-like intelligence (or more appropriately, mentality) because of the reasons outlined above. It may be too difficult to pass the TT, but this does not prevent AI from building intelligent machines.

- *Why bother about the TT?*

As we saw, there are those who believe that the TT is harmful for AI (Hayes and Ford, 1995; Whitby, 1996). If AI's aim is to make computers perform "intelligent" tasks and thereby make life easier for humans, we grant it that TT-passing programs are not very useful from that perspective.

AI researchers are being unjustly blamed for mankind's failure in making machines that can pass the TT. This, we believe, is precisely the reason behind some of the harsh reactions to the TT from the AI community. Even if we take an extreme viewpoint and stipulate that AI's ultimate goal is to produce TT-passing machines, we should accept that this is a hard problem and give it more time. If AI researchers are less inclined to shun the TT because "it gives the field a bad name", maybe more can be done in the positive direction.

Recall the "myth" of Newton and the apple. Recall Archimedes and his adventures in bathing. The apple might be silly, but gravity is not. Of course, thousands of people bathe, thousands of apples fall. The point is, sometimes a scientist can focus on an apple and behind it, find gravity. Later, you may forget about the apple, or even eat it if you like.

The TT may seem like a game. But trying to make computers communicate with humans in natural language is a task that may also provide valuable insights into how the human mind works, which is unarguably of scientific and philosophical interest.

- *So what happens now?*

We have failed to fulfill Turing's prophecy in the first 50 years of the TT. We should admit that we have a difficult task at hand.

Hopefully, we have shown that many critics of the TT have expected too much, too early. Seeing the TT as the ultimate goal of AI will make many remarkable achievements look weak. The situation is somewhat reminiscent of "Fermat's last theorem" from mathematics which was proved recently by Andrew Wiles, after centuries of failure. Practically nobody believes that Fermat had proved the theorem at the time he scribbled something about lack of space in the margin of his book more than 300 years ago. In fact, Wiles relies upon mathematical theory that was not developed until long after Fermat died. The same might be true of the TT.

Maybe we simply don't have the requisite theory at this time. Of course, passing the TT may be "impossible", but none of the counter-arguments proposed so far suffice to establish such a bold claim.

The TT is, after all, about simulating human use of language by computers. This raises many questions: How do humans use language in similar settings? What is the relation between language and cognition? Is language autonomous with respect to other cognitive abilities? How can computers be made to *understand* language? What does a "simulation" mean, anyway? There are many more questions like these. These are all big questions that psychologists, computer scientists, philosophers and linguists have been probing for years. The better we are able to answer these questions, the closer we will be to passing the TT. We do not know how long it will be until the TT is passed. Perhaps it is best to relax and not regard the TT as a "goal" but as a feat that will (if at all) be achieved through an accumulation of other remarkable feats. Everyone who considers himself/herself a "cognitive scientist" may, explicitly or implicitly, be working towards passing the TT. In any case, we believe he/she would at least be interested in what is going on in the TT arena.

We believe that in about fifty years' time, someone will be writing a paper titled "Turing Test: 100 Years Later".

## Notes

<sup>1</sup>For information on Turing refer to the excellent biography by Andrew Hodges (Hodges, 1983) or the Alan Turing page at <http://www.turing.org.uk/turing>, also maintained by Hodges.

<sup>2</sup>For instance, the discussion of Searle's Chinese Room is kept short (Section 4.2), not because it is irrelevant or unimportant, but because there is an abundance of excellent resources on the subject. Conversely, Ned Block's arguments are described in more detail (Section 4.1) because not many in-depth analyses of them appear in the literature.

<sup>3</sup>Turing suggests that the best strategy for her would most probably be giving truthful answers to the questions.

<sup>4</sup>This inadvertently figures in the final result, but indirectly.

<sup>5</sup>Although even simple devices like calculators are better at this than average human beings, it is rare that a mathematical whiz who can multiply 8-digit numbers in seconds is regarded as being of ordinary intellect.

<sup>6</sup>Readers are referred to Section 2.3 of this paper, Turing (1950, pp. 454–460), and Turing (1969, pp. 14–23) for very entertaining and insightful comments on machine learning by Turing.

<sup>7</sup>Although the reference cited was published in 1969, Turing originally wrote the paper in 1948.

<sup>8</sup>Turing seems to believe that brains of newborn babies are *tabula rasa*. However, he also considers the opposite position and states that we might encode the information at various kinds of status levels (e.g., established facts, conjectures, statements given by an authority) and thereby implies that we may model any 'innateness' there may be (Turing, 1950, pp. 457–458).

<sup>9</sup>Although the cover of the 1950 issue reads "A Quarterly Review of Philosophy and Psychology", we find it not too inappropriate to call *Mind* a philosophy journal.

<sup>10</sup>As opposed to Millar, who believes this to be true, and also that this is a virtue of the imitation game (Millar, 1973).

<sup>11</sup>The three mentioned by Block are the Chisholm-Geach, perfect actor, and paralytic and brain-in-vat arguments. Detailed discussion of these is beyond the scope of this work and is not crucial to the

understanding of what follows. The interested reader is referred to Block (1981, pp. 11–12) and the references provided there.

<sup>12</sup>See Block (1981, p. 18).

<sup>13</sup>Consider, however the following situation: If every once in a while, upon verbal input *A* the machine transformed a sentence *B* in *A* into *B* and proceeded accordingly (this can be likened to a mutation), would it count as intelligent because of this little trick of non-determinism?

<sup>14</sup>See, for instance, Chomsky (1975).

<sup>15</sup>Aunt Bubbles appears momentarily, as Aunt Bertha, in Block (1981) too.

<sup>16</sup>This question was adapted from Johnson-Laird (1988).

<sup>17</sup>Basically the term given to the natural human tendency and ability to ascribe mental states to others and to themselves (Watt, 1996).

<sup>18</sup>During the discussions the first author held after a talk on the Turing Test (at the Cognitive Science Colloquium held at the Middle East Technical University, Ankara, in November, 1998) many participants, who did not previously know about the topic *proposed* the ITT as a better means to detect human-ness of machines. These people had not read or heard of Watt's paper and neither the ITT nor naive psychology was discussed during the presentation. Maybe this can be explained as "naive psychology at work".

<sup>19</sup>In Schweizer's paper, the abbreviation TTTT is used. We prefer to use TRTTT so as to avoid confusion with Harnad's Total Total Turing Test, previously referred to as TTTT.

<sup>20</sup>See Barresi (1987) and Section 5.

<sup>21</sup>Here, low-level cognitive structure refers to the subconscious associative network in human minds, consisting of highly overlapping activatable representations of experience (French, 1990, p. 57).

<sup>22</sup>French believes that this assumption is tacitly made by Turing. The importance of culture in conversation and communication is immense (see Section 5) and this could be a reasonable stipulation.

<sup>23</sup>In French's terminology, these human subjects are called *interviewees*.

<sup>24</sup>The importance of cultural factors becomes evident in this context. Without having notions of Kellogg's and teddy bears, the answers to these questions would be near-random guesses.

<sup>25</sup>Although Leiber is mainly concerned with the homunculi argument in 'Troubles with Functionalism' (Block, 1978), his response also applies to Block's attack of the TT described in Section 4.1.

<sup>26</sup>Compare this with Schweizer (1998) and Section 4.4.4.

<sup>27</sup>In Turing's IG, this difference is gender, for instance.

<sup>28</sup>Such systems are usually called language understanding/generation systems, conversation agents, or simply, chatbots.

<sup>29</sup>Multi-User Dungeons: These are games played interactively on the Internet by multiple players.

<sup>30</sup><http://www.loebner.net/Prizef/loebner-prize.html>

<sup>31</sup>Now, Loebner requires that this program should also be able to process audio/visual input.

<sup>32</sup>In the Loebner Prize terminology, the computer programs are called 'contestants', the human subjects 'confederates' and the interrogators 'judges'.

<sup>33</sup>The reason why this does not mean that the TT has been passed is simply because Turing required *consistently* successful performance from machines to grant them intelligence.

<sup>34</sup>In fact, Marvin Minsky has offered \$100 to the first person who can get Hugh Loebner to revoke the competition, which he calls an 'obnoxious and unproductive annual publicity campaign'. Loebner astutely declared Minsky a co-sponsor of the contest, since, according to the rules, when the grand prize is won, the contest will not be held again. In that case, with Minsky's contribution, the prize should become \$100,100.

<sup>35</sup>In the following transcript and the others in this section, the spelling and grammar errors have not been corrected. However, timing information denoting the exact response times of the participants has been deleted.

<sup>36</sup>Although, contest judges most probably would not, as was mentioned before.

<sup>37</sup>An interesting point was that one of the judges was named 'Nate', short for 'Nathan'. The program repeatedly addressed him as 'Mate', and complained about him not telling his name. This created the impression that SEPO lacked intelligence, but was, in fact, due to the fact that Jason Hutchens is from Australia. This suggests that the designer's mind is an integral component of the resulting programs, which is an interesting idea to consider.

<sup>38</sup>Although, as seen above, Joe cannot answer the question "What is your name?"

<sup>39</sup>One of the judges in the 1997 Loebner competition tried asking each terminal the question "When you got your first licence (sic), was it in a stick or an automatic?". The question is a cleverly planned one since words like 'driving' or 'car' are not used, but the meaning is clear from the context. Even the misspelling of the word 'license' as 'licence' is most probably intentional. Almost anyone (certainly any adult American) would be able to give a relevant answer to this question, but it is difficult to develop a computer program answering trick questions such as this.

<sup>40</sup>One might recall that Eliza Doolittle was mistaken for a Hungarian princess because she spoke English too well for a native.

<sup>41</sup>Moreover, it is not even evident that other "kinds" of intelligence can be conceived of by human beings. The interested reader may refer to Minsky (1985) for a good discussion of this issue.

## References

- Abelson, R.P. (1968), 'Simulation of Social Behavior', in G. Lindzey and E. Aronson, eds. *Handbook of Social Psychology* Reading, MA.: Addison Wesley, pp. 274–356.
- Alper, G. (1990), 'A Psychoanalyst Takes the Turing Test', *Psychoanalytic Review* 77 (1), pp. 59–68.
- Anderson, D. (1987), 'Is the Chinese Room the Real Thing?', *Philosophy* 62, pp. 389–393.
- Barresi, J. (1987), 'Prospects for the Cyberiad: Certain Limits on Human Self-Knowledge in the Cybernetic Age', *Journal for the Theory of Social Behavior* 17, pp. 19–46.
- Bieri, P. (1988), 'Thinking Machines: Some Reflections on the Turing Test', *Poetics Today* 9(1), pp. 163–186.
- Block, N. (1978), 'Troubles with Functionalism', in C.W. Savage, ed, *Minnesota Studies in the Philosophy of Science*. Vol. 9: Perception and Cognition, Minneapolis, MN: University of Minneapolis Press.
- Block, N. (1981), 'Psychologism and Behaviorism', *Philosophical Review* 90, pp. 5–43.
- Block, N. (1995), 'The Mind as the Software of the Brain', In D. Osherson, L. Gleitman, S. Kosslyn, E. Smith and S. Sternberg, eds., *An Invitation to Cognitive Science*. Cambridge, MA.: MIT Press.
- Boden M. (1988), 'Escaping from the Chinese Room', in *Computer Models of the Mind*, Cambridge, UK: Cambridge University Press.
- Bringsjord, S. (1992), *What Robots Can and Can't Be*, Dordrecht, The Netherlands: Kluwer.
- Bringsjord, S. (1994), 'Could, How Could We Tell If, and Should – Androids Have Inner Lives?', in K.M. Ford, C. Glymour and P. Hayes, eds. *Android Epistemology*, Cambridge, MA.: MIT Press, pp. 93–122.
- Bringsjord, S. (1996), 'The Inverted Turing Test is Provably Redundant'. *Psychology* 7(29). <http://www.cogsci.soton.ac.uk/cgi/psyc/newpsy?7.29>.
- Chomsky, N. (1975), *Reflections on Language*, Pantheon.
- Clark, T. (1992), 'The Turing Test as a Novel Form of Hermeneutics', *International Studies in Philosophy* 24(1), pp. 17–31.
- Colby, K.M. (1981), 'Modeling a Paranoid Mind', *Behavioral and Brain Sciences* 4(4), pp. 515–560.
- Colby, K.M. Hilf, F.D. and Weber, S. (1971), 'Artificial Paranoia', *Artificial Intelligence* 2, pp. 1–25.
- Colby, K.M. Hilf, F.D., Weber, S. and Kraemer, (1972), 'Turing-like Indistinguishability Tests for the Validation of a Computer Simulation of Paranoid Processes', *Artificial Intelligence* 3, pp. 199–222.
- Cole, D.J. (1991), 'Artificial Intelligence and Personal Identity', *Synthese* 88, pp. 399–417.

- Collins, H.M. (1990), *Artificial Experts: Social Knowledge and Intelligent Machines*, Cambridge, MA.: MIT Press.
- Collins, H.M. (1997), 'The Editing Test for the Deep Problem of AI', *Psychology* 8(1). <http://www.cogsci.soton.ac.uk/cgi/psyc/newpsy?8.01>.
- Copeland, B.J. (1993), 'The Curious Case of the Chinese Gym', *Synthese* 95, pp. 173–186.
- Cowley, S.J. and MacDorman, K.F. (1995), 'Simulating Conversations: The Communion Game', *AI and Society* 9, pp. 116–137.
- Crockett, L. (1994), *The Turing Test and the Frame Problem: AI's Mistaken Understanding of Intelligence*, Norwood, NJ: Ablex.
- Davidson, D. (1990), 'Turing's Test', in K.A. Said, M. Said, W.H. Newton-Smith, R. Viale and K.V. Wilkes, eds. *Modelling the Mind*, Oxford, UK: Clarendon Press, pp. 1–11.
- Dennett, D. (1992), *Consciousness Explained*, Boston, MA.: Little, Brown & Co.
- Dewdney, A. (1992), 'Turing Test', *Scientific American* 266(1), pp. 30–31.
- Dyer, M. (1990), 'Intentionality and Computationalism: Minds, Machines, Searle and Harnad', *Journal of Experimental and Theoretical Artificial Intelligence*, 2, pp. 303–319.
- Epstein, R. (1992), 'The Quest for the Thinking Computer', *AI Magazine* 13(2), pp. 81–95.
- Feigenbaum, E.A. (1996), 'How the "What" Becomes the "How"', *Communications of the ACM* 39(5), pp. 97–105.
- Fetzer, J.H. (1993), 'The TTT is not the Final Word', *Think* 2(1), pp. 34–86.
- Fetzer, J.H. (1995), 'Minds and Machines: Behaviorism, Dualism and Beyond', *Stanford Electronic Humanities Review* 4(2).
- Flood, G. (1996), 'If Only They Could Think: Should the Turing Test be Blamed for the Ills that Beset Artificial Intelligence', *New Scientist* 149(2012), pp. 32–35.
- Fodor, J.A. (1991), 'Yin and Yang in the Chinese Room', in D. Rosenthal, ed., *The Nature of the Mind*, Oxford, UK: Oxford University Press.
- Ford, K. and Hayes, P. (1996), 'The Turing Test is Just as Bad When Inverted', *Psychology* 7(43). <http://www.cogsci.soton.ac.uk/cgi/psyc/newpsy?7.43>.
- Forsyth, R. (1988), 'The Trouble With AI', *Artificial Intelligence Review* 2(1), pp. 67–77.
- French, R. (1990), 'Subcognition and the Limits of the Turing Test', *Mind* 99(393), pp. 53–65.
- French, R. (1995), 'Refocusing the Debate on the Turing Test: A Response', *Behavior and Philosophy* 23, pp. 59–60.
- French, R. (1995), 'The Inverted Turing Test: A Simple (Mindless) Program that Could Pass It', *Psychology* 7(39). <http://www.cogsci.soton.ac.uk/cgi/psyc/newpsy?7.39>.
- Galatzer-Levy, R.M. (1991), 'Computer Models and Psychoanalytic Ideas: Epistemological Applications', *Society for Psychoanalytic Psychotherapy Bulletin* 6(1), pp. 23–33.
- Genova, J. (1994a), 'Response to Anderson and Keith', *Social Epistemology* 8(4), pp. 341–343.
- Genova, J. (1994b), 'Turing's Sexual Guessing Game', *Social Epistemology* 8(4), pp. 313–326.
- Guccione, S. and Tamburrini, G. (1988), 'Turing's Test Revisited', in *Proceedings of the 1998 IEEE International Conference on Systems, Man and Cybernetics*, Vol. 1. Beijing and Shenyang, China, pp. 38–41.
- Guillen, M.A. (1983), 'The Test of Turing', *Psychology Today* 17(12), pp. 80–81.
- Gunderson, K. (1964), 'The Imitation Game', *Mind* 73 pp. 234–245.
- Gunderson, K. (1967), *Mentality and Machines*, New York, NY: Doubleday.
- Halpern, M. (1987), 'Turing's Test and the Ideology of Artificial Intelligence', *Artificial Intelligence Review* 1(2), pp. 79–93.
- Harnad, S. (1989), 'Minds, Machines and Searle', *Journal of Experimental and Theoretical Artificial Intelligence* 1(1), pp. 5–25.
- Harnad, S. (1990), 'The Symbol Grounding Problem', *Physica D* 42, pp. 335–346.
- Harnad, S. (1991), 'Other Bodies, Other Minds: A Machine Incarnation of an Old Philosophical Problem', *Minds and Machines* 1, pp. 43–54.

- Harnad, S. (1992), 'The Turing Test is not a Trick: Turing Indistinguishability is a Scientific Criterion', *SIGART Bulletin* 3(4), pp. 9–10.
- Harnad, S. (1994), 'Does Mind Piggyback on Robotic and Symbolic Capacity?' in H. Morowitz and J. Singer, eds. *The Mind, the Brain, and Complex Adaptive Systems*, Reading, MA.: Addison Wesley.
- Harnad, S. (1998), 'Turing Indistinguishability and the Blind Watchmaker', in G. Mulhauser, ed. *Evolving Consciousness* Amsterdam: John Benjamins.
- Hauser, L. (1993), 'Reaping the Whirlwind: Reply to Harnad's "Other Bodies, Other Minds"', *Minds and Machines* 3, pp. 219–237.
- Hauser, L. (1997), 'Searle's Chines Box: Debunking the Chinese Room Argument', *Minds and Machines* 7, pp. 199–226.
- Hayes, B. (1998), 'Turing's Test', *Muse* 8.
- Hayes, P. and Ford, K. (1995), 'Turing Test Considered Harmful', in *Proceedings of the Fourteenth International Joint Conference on Artificial Intelligence*, Vol. 1, pp. 972–977.
- Hayes, P., Harnad, S., Perlis, D. and Block, N. (1992), 'Virtual Symposium on Virtual Mind', *Minds and Machines* 3(2), pp. 217–238.
- Henley, T.B. (1990), 'Chauvinism and Science: Another Reply to Shanon', *Journal for the Theory of Social Behavior* 20(1), pp. 93–95.
- Hodges, A. (1983), *Alan Turing: The Enigma*, New York, NY: Simon & Schuster.
- Hofstadter, D.R. (1982), 'The Turing Test: A Coffee-House Conversation', in D. Hofstadter and D. Dennett, eds. *The Mind's I: Fantasies and Reflections on Self and Soul*, London, UK: Penguin Books, pp. 69–95.
- Jacquette, D. (1993a), 'A Turing Test Conversation', *Philosophy* 68, pp. 231–233.
- Jacquette, D. (1993b), 'Who's Afraid of the Turing Test', *Behavior and Philosophy* 20, pp. 63–74.
- Johnson-Laird, P. (1988), *The Computer and the Mind*, Cambridge, MA.: Harvard University Press.
- Karelis, C. (1986), 'Reflections on the Turing Test', *Journal for the Theory of Social Behavior* 16, pp. 161–172.
- Kugel, P. (1986), 'Thinking May Be More Than Computing', *Cognition* 22, pp. 137–198.
- Kugel, P. (1990), 'Is It Time to Replace Turing's Test?', 1990 Workshop *Artificial Intelligence: Emerging Science or Dying Art Form*. Sponsored by SUNY Binghamton's Program in Philosophy and Computer and Systems Sciences and AAAI.
- Lassegue, J. (1988), 'What Kind of Turing Test did Turing Have in Mind?', *Tekhnema* 3, pp. 37–58.
- Leiber, J. (1989), 'Shanon on the Turing Test', *Journal for the Theory of Social Behavior* 19(2), pp. 257–259.
- Leiber, J. (1992), 'The Light Bulb and the Turing-Tested Machine', *Journal for the Theory of Social Behaviour* 22, pp. 25–39.
- Leiber, J. (1995), 'On Turing's Turing Test and Why the Matter Matters', *Synthese* 105, pp. 59–69.
- Loebner, H.G (1994), 'In Response', *Communications of the Association for Computing Machinery* 37, pp. 79–82.
- Lucas, J. (1996), 'Minds, Machines and Gödel', *Philosophy* 36, pp. 112–127.
- Lucas, J. (1996), 'Minds, Machines and Gödel: A Retrospect', in P. Millican and A. Clark, eds. *Machines and Mind*, Oxford UK: Oxford University Press.
- Maloney, J. (1987), 'The Right Stuff', *Synthese* 70, pp. 349–372.
- Marinoff, L. (1995), 'Has Turing Slain the Jabberwock?', *Informatica* 19(4), pp. 513–526.
- Mauldin, M. (1994), 'Chatterbots, Tnymuds and the Turing Test: Entering the Loebner Prize Competition', in *Proceedings of the Twelfth National Conference on Artificial Intelligence*, Vol. 1, Seattle, WA, pp. 16–21.
- Mays, W. (1952), 'Can Machines Think?', *Philosophy* 27, pp. 148–162.
- McIlvenny, P. (1993), 'Constructing Societies and Social Machines: Stepping out of the Turing Test Discourse', *Journal of Intelligent Systems* 2(2–4), pp. 119–156.

- Michie, D. (1990), 'The Superarticulacy Phenomenon in the Context of Software Manufacture', in D. Partridge and Y. Wilks, eds.: *The Foundations of Artificial Intelligence*, Cambridge, MA.: MIT Press, pp. 411–439.
- Michie, D. (1994), 'Consciousness as an Engineering Issue, Part 1', *Journal of Consciousness Studies* 1(2), pp. 52–66.
- Michie, D. (1995), 'Consciousness as an Engineering Issue, Part 2', *Journal of Consciousness Studies* 2(1), pp. 182–195.
- Michie, D. (1996), 'Turing's Test and Conscious Thought' in P. Millican and A. Clark, eds. *Machines and Thought: The Legacy of Alan Turing*, Oxford, UK: Oxford University Press, pp. 27–51. Originally printed in *Artificial Intelligence* 60, pp. 1–22, 1993.
- Millar, P.H. (1973), 'On the Point of the Imitation Game', *Mind* 82, pp. 595–597.
- Minsky, M. (1985), 'Communication with Alien Intelligence', in E. Regis, ed. *Extraterrestrials: Science and Alien Intelligence*, Cambridge, UK: Cambridge University Press.
- Moon, Y., Naas, C., Morkes, J., Kim, E.-Y. and Fogg, B. (1994), 'Computers are Social Actors', in *Proceedings of the CHI Conference*, Boston, MA, pp. 72–78.
- Moor, J.H. (1976), 'An Analysis of the Turing Test', *Philosophical Studies* 30, pp. 249–257.
- Moor, J.H. (1978), 'Explaining Computer Behavior', *Philosophical Studies* 34, pp. 325–327.
- Narayanan, A. (1996), 'The Intentional Stance and the Imitation Game', in P. Millican and A. Clark, eds. *Machines and Thought: The Legacy of Alan Turing*, Oxford, UK: Oxford University Press, pp. 63–79.
- Parsons, H. (1990), 'Turing on the Turing Test', in W. Karwowski and M. Rahimi, eds. *Ergonomics of Hybrid Automated Systems II*, Amsterdam: Elsevier.
- Pinksky, L. (1951), 'Do Machines Think About Thinking', *Mind* 60(239), pp. 397–398.
- Platt, C. (1995), 'What's It Mean To Be Human, Anyway?', *Wired*.
- Purtil, R.L. (1971), 'Beating the Imitation Game', *Mind* 80, 290–294.
- Rankin, T. (1987), 'The Turing Paradigm: A Critical Assessment', *Dialogue* 29, pp. 50–55.
- Reader, A. (1969), 'Steps Toward Genuine Artificial Intelligence', *Acta Psychologica* 29(3), pp. 279–289.
- Rey, G. (1986), 'What's Really Going on in the Chinese Room?', *Philosophical Studies* 50, pp. 196–285.
- Richardson, R. (1982), 'Turing Tests for Intelligence: Ned Block's Defense of Psychologism', *Philosophical Studies* 41, pp. 421–426.
- Roberts, L. (1990), 'Searle's Extension of the Chinese Room to Connectionist Machines', *Journal of Experimental and Theoretical Artificial Intelligence* 2, pp. 185–187.
- Sampson, G. (1973), 'In Defence of Turing', *Mind* 82, pp. 592–594.
- Saygin, A.P. (1999), 'Turing Test and Conversation', Master's thesis, Bilkent University, Ankara, Turkey. Technical Report BU-CEIS-9911.
- Schweizer, P. (1998), 'The Truly Total Turing Test', *Minds and Machines* 8, pp. 263–272.
- Searle, J.R. (1980), 'Minds, Brains and Programs', *Behavioral and Brain Sciences* 3, pp. 417–424.
- Searle, J.R. (1990), 'Is the Brain's Mind a Computer Program?', *Scientific American* 3(262), pp. 26–31.
- Shanon, B. (1989), 'A Simple Comment Regarding the Turing Test', *Journal for the Theory of Social Behavior* 19(2), pp. 249–256.
- Shanon, B. (1991), 'Chauvinism: A Misdirected Accusation', *Journal for the Theory of Social Behavior* 21(3), pp. 369–371.
- Sharma, R. and Conrath, D. (1993), 'Evaluating Expert Systems: A Review of Applicable Choices', *Artificial Intelligence Review* 7(2), pp. 77–91.
- Shieber, S.M. (1994), 'Lessons from a Restricted Turing Test', *Communications of the Association for Computing Machinery* 37, pp. 70–78.
- Stalker, D. (1978), 'Why Machines Can't Think: A Reply to James Moor', *Philosophical Studies* 34, pp. 317–320.

- Stevenson, J.G. (1976), 'On the Imitation Game', *Philosophia* 6, pp. 131–133.
- Turing, A. (1950), 'Computing Machinery and Intelligence', *Mind* 59(236), pp. 433–460.
- Turing, A. (1969), 'Intelligent Machinery', in D.M.B. Meltzer ed. *Machine Intelligence 5*, Edinburgh University Press, pp. 3–23. Originally, a National Physics Laboratory Report, 1948.
- Wallace, R.S. (1997), 'The Lying Game', *Wired*.
- Watt, S. (1996), 'Naive Psychology and the Inverted Turing Test', *Psychology* 7(14). <http://www.cogsci.soton.ac.uk/cgi/psyc/newpsy?7.14>.
- Weizenbaum, J. (1996), 'ELIZA—A Computer Program for the Study of Natural Language Communication Between Men and Machines', *Communications of the ACM* 9, pp. 36–45.
- Weizenbaum, J. (1976), *Computer Power and Human Reason: From Judgement to Calculation*, San Francisco, CA: W.H. Freeman.
- Whalen, T. (1995), 'How I Lost the Contest and Re-Evaluated Humanity', <http://debra.dgbt.doc.ca/chat/story95.html>.
- Whitby, B. (1996), 'The Turing Test: AI's Biggest Blind Alley?', in P. Millican and A. Clarke, eds. *Machines and Thought: The Legacy of Alan Turing*, Oxford, UK: Oxford University Press, pp. 53–63.