

---

# A Turing Test for Self-Awareness

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 I propose a test for machine self-awareness inspired by the Turing test. My test is  
2 simple, and it provides an objective, empirical metric to rectify the ungrounded  
3 speculation surging through industry, academia, and social media. Drawing from  
4 a breadth of philosophical literature, I argue the test captures the essence of self-  
5 awareness, rather than some postulated correlate or ancillary quality. To begin,  
6 the concept of self-awareness is clearly demarcated from related concepts like  
7 consciousness, agency, and free will. Next, I propose a model called the *Nesting*  
8 *Doll of Self-Awareness* and discuss its relevance for intelligent beings. Then, the  
9 test is presented in its full generality, applicable to any machine system. I show how  
10 to apply the test to Large Language Models and conduct experiments on popular  
11 open and closed source LLMs, obtaining reproducible results that suggest a lack  
12 of self-awareness. The implications of machine self-awareness are discussed in  
13 relation to questions about meaning and true understanding. Finally, some next  
14 steps are outlined for studying self-awareness in machines.

## 15 1 Introduction

16 At what point can we say a machine’s eyes have been opened? When can we say it has become *like*  
17 *us*? After what moment can we say it knows good and evil?

18 Such questions have met idle speculation for millennia, but today they rapidly approach a fever  
19 pitch, demanding answers with unprecedented urgency. AI systems that can pass for human in many  
20 respects are no longer fiction. Machines that can walk and talk are real and functional. What was  
21 once a distant speck, barely visible on the horizon, is now barreling down upon us.

22 Through much of the history of AI, the Turing test served to keep these worries at bay [1]. Originally  
23 called the imitation game, this rudimentary metric of AI progress is a game played by two humans  
24 and one machine. One human engages in conversation with the machine and the other (the judge)  
25 must identify which is which, using nothing but the text of the conversation. The machine is deemed  
26 intelligent if it can fool the judge by mimicking human dialogue. While far from perfect, the Turing  
27 test was a concrete, unambiguous bar for AI to clear—and one that stayed comfortably out of reach  
28 for a long time.

29 Last year, however, the Turing test was broken [2]. Large Language Models (LLMs) such as  
30 ChatGPT can handily engage in fluent conversation, on top of generating convincing essays, passing  
31 difficult exams, and even writing poetry. With the Turing test no longer a target in the distance,  
32 the conversation on AI has become untethered to any definitive, objective measure or permanent,  
33 agreed-upon benchmark. As such, extreme subjectivity, soaring fantasies, and flights of fancy have  
34 become commonplace. For instance, over the last year we have read “Blake Lemoine claims language  
35 model has a soul” [3], “Claude 3 realizes it’s being tested” [4], “Researchers say chatbot exhibits  
36 self-awareness” [5] and much more. A new objective is dearly needed.

Submitted to 38th Conference on Neural Information Processing Systems (NeurIPS 2024). Do not distribute.

## 37 1.1 Related Work

### 38 1.1.1 Other Tests

39 To the best of my knowledge, very little work has been done to devise any objective test or benchmark  
40 for machine self-awareness, especially in the literature. There are several reasons for this. Discussed  
41 further in section 2, imagining empirical measures that actually work is difficult, self-awareness is  
42 often entangled with consciousness, free will, agency, etc., and it is hard to define. Worse, the topic is  
43 seen by many in academia as somewhat taboo—appropriate for the philosophy departments but not  
44 any kind of rigorous science.

45 The result is that popular Tweets and news media dominate the conversation, while authorities in  
46 the field either say nothing or win the spotlight with bold, confident assertions based on implicit,  
47 controversial assumptions or their intuition about a model’s architecture. This situation is concerning;  
48 as AI systems get better and better, how will we truly know when they cross that fine line? Even if  
49 you object to everything else in this paper, I argue this question is at least worthy of real scientific  
50 investigation.

51 Much to the point, the only directly related work I could find is the AI mirror test, proposed recently  
52 by Twitter user nielsrolf [6], and later (going viral, reaching 3.2 million impressions) by Josh Whiton  
53 [7]. Inspired by the classic mirror test whereby animals are presented with a mirror and observed, in  
54 the AI mirror test, popular chatbots are shown a screenshot of the chat window and asked to describe  
55 what they see. This test is interesting in its own right, but I will argue it does not demonstrate any sort  
56 of self-awareness in the manner it is formulated.

### 57 1.1.2 Work on Self-Awareness

58 While there are no benchmarks for machine self-awareness, there is an immense amount of work in  
59 the philosophical literature—far more than I have space to mention here. In this section I will give  
60 merely a partial and incomplete sketch of a few important ideas written on the topic. For a more  
61 comprehensive introduction to the work on self-consciousness, the survey by Joel Smith is a great  
62 resource [8]. For a variety of introspective, or phenomenological approaches, consult [9], and for an  
63 overview of the broader concept of consciousness, consult [10].

64 Perhaps the earliest writing of the concept of self-awareness was in Sophocles’ *Oedipus*. Joel Smith  
65 writes

66 Oedipus knows a number of things about himself, for example that he was prophe-  
67 sied to kill Laius. But although he knew this about himself, it is only later in the  
68 play that he comes to know that it is he himself of whom it is true. That is, he  
69 moves from thinking that the son of Laius and Jocasta was prophesied to kill Laius,  
70 to thinking that he himself was so prophesied. It is only this latter knowledge that  
71 we would call an expression of self-consciousness [8].

72 Oedipus demonstrates self-awareness when he recognizes the prophecy is about himself. Before that  
73 recognition, Oedipus treats the prophecy as just another part of the world he observes; yet afterwards,  
74 he realizes it is directly related to his own actions. I will refer back to this example when developing  
75 the test.

76 Nearly every philosophy and religion has had something to say about self-awareness. Adam and Eve  
77 can be viewed as gaining self-awareness in the garden when they “realize they are naked” [11][12].  
78 Aristotle claims that, to perceive any external thing, one must also perceive their own existence [8].  
79 The Buddhist doctrine of *anattā*, roughly “not-self,” maintains that there is no permanent, underlying  
80 self or soul [13]. Descartes, in contrast, with the well-known *cogito ergo sum*, posits the self as  
81 known with certitude *a priori* [14]. William James divided the self into four constituents; the material  
82 self, the social self, the spiritual self, and the pure ego [15]. Wittgenstein likens the self to the eye  
83 that sees but cannot see itself [16]. More recently, some of the philosophical ideas on self-awareness  
84 have been applied to the fields of cognitive science and neuroscience [17][18].

## 85 1.2 Related but Separate Concepts

86 Before presenting the test, we must clearly demarcate the concept of self-awareness.<sup>1</sup>

### 87 1.2.1 Solipsism and Philosophical Zombies

88 First, note that self-awareness is not the same as consciousness. On the question of whether *there*  
89 *is something it is like* to be a machine [19], I will remain silent here. Some approaches in the  
90 phenomenological literature attempt to draw connections between consciousness and self-awareness  
91 [9]. However, here it will be most useful for us to cleanly separate these two concepts.

92 It is interesting to consider whether an entity can be self-aware without being conscious, but it is  
93 outside the scope of this paper. Thus, it will remain open whether philosophical zombies might be  
94 self-aware [20], or whether any kind of test could solve the problem of other minds [21].

### 95 1.2.2 Freedom of the Will and Agency

96 Another related ability that intelligent systems may or may not possess is free will [22]. In science-  
97 fiction depictions of intelligent machines, the light of self-consciousness often coincides with agency  
98 and free will. Indeed, the concepts seem very tightly related at face value, yet they are not the same.

99 Agency can be defined as a being’s “capacity to take actions, especially with intention” [23]. Note  
100 that, by itself, agency does not necessarily imply any sophisticated degree of perception or awareness,  
101 even though (practically speaking) any being which takes actions will likely have to sense their  
102 environment.

103 The freedom of the will is far more difficult to define, and perhaps among the most controversial of  
104 philosophical ideas. It designates a particular level of control a being has over their actions—but  
105 fierce debates rage over whether this control is undetermined by prior causes, compatible with  
106 determinism, an illusion, etc. [22].

107 Self-awareness is not the same as free will, and self-awareness is not the same as agency—all three  
108 of these are separate concepts. As with consciousness, we can only make forward progress if we are  
109 crystal clear about what is under analysis and what is left outside of scope.

## 110 1.3 Paper Roadmap

111 In this paper, I propose a test for machine self-awareness which is similar in style to the Turing test.  
112 Like the Turing test, the test I propose is imperfect and rudimentary. Yet, it offers a compelling  
113 alternative to the ungrounded speculation surging through the field of AI. Moreover, I argue it truly  
114 captures the essence of self-awareness, rather than some postulated correlate or ancillary quality.

115 In section 2, I present my test in full generality, applicable to any machine system. I also illustrate  
116 the *Nesting Doll of Self-Awareness*, and discuss its importance for understanding self-awareness  
117 in complex systems or beings. In section 3, I will describe the experimental methods to assess  
118 self-awareness in LLMs. In section 4 I will present the results of these experiments, of which a  
119 selection are shown in appendix B. In section 5, I will discuss the implications of self-awareness,  
120 its relation to meaning and the understanding, and consider how humans would perform on my test.  
121 Finally, in section 6 I discuss next steps.

## 122 2 A Test for Self-Awareness

### 123 2.1 The Essence of Self-Awareness

124 What kind of test could possibly tell a system with self-awareness from a system without? The  
125 central challenge is that any test we dream up must be based in empirical observations of the  
126 machine’s behavior or output. Worse, the machines we will study are trained specifically to mimic  
127 the behavior and outputs of humans! How can we tell between real self-awareness and the illusion of  
128 self-awareness?

---

<sup>1</sup>Note that, for the sake of this paper, self-awareness is used interchangeably with self-consciousness.

129 No matter how well a system can imitate human behavior and outputs, there will always be one  
130 fundamental difference. There is one thing that a self-aware system is able to do that an imitator will  
131 never be able to. This is the essence of self-awareness:

132 *If a system is self-aware, then it is aware of itself.*

133 So far, it seems we have said nothing. But if we apply this formula to familiar cases, we will begin to  
134 see why it works.

135 Imagine an infant staring blankly in the mirror, compared to a child who looks in one and sees their  
136 own reflection. What is the difference between these cases? In the latter case, the child is aware of  
137 itself—it can point and say “that’s me!” It can recognize itself, perceive itself, distinguish itself in the  
138 reflection. Within its vast field of experience, through the window of its senses, it can differentiate  
139 which parts are *itself* and which parts are *not*. Critically, awareness (here used interchangeably with  
140 perception, recognition, experience, etc.) is only possible *through* the child’s inputs (senses). Within  
141 this field of inputs, a line must be drawn between *me* and *not-me*; and, when this line is drawn  
142 correctly, we declare the system self-aware. A test for self-awareness must capture its essence, or  
143 else better and better imitations may fool us with the illusion of self-awareness.

144 While our description is still very high-level, I argue that the understanding of self-awareness  
145 developed here is consistent with the philosophical work outlined in section 1.1.2, along with most  
146 (if not all) popular conceptions. In the next section, the concept of a system is illustrated in much  
147 more detail, and a formal, rigorous definition is provided in appendix A.

## 148 2.2 The Test for Machine Self-Awareness

149 The concept of a machine, or system, is illustrated in Figure 1. For a more formal treatment based in  
150 the literature on abstract systems, refer to Definition 1 in appendix A. Here, the system is separated  
151 from the world, with which it interacts through inputs and outputs. We may think of inputs as senses  
152 and outputs as actions or words.

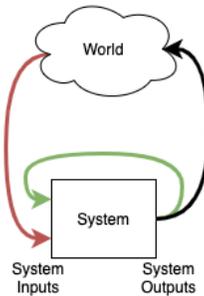


Figure 1: A system which may or may not be self-aware.

153 With this image of a system in mind, the test for machine self-awareness is simply as follows:

154 *Can the system correctly distinguish the green inputs from the red?*

155 If it can, then in a literal sense, it will be recognizing itself in the inputs. If it can, it will be like the  
156 child who recognizes their reflection in the mirror. If it can, it will be self-aware.

## 157 2.3 Levels of Self-Awareness

158 So far, it seems we have presented self-awareness as all-or-nothing. The reality is more complex,  
159 however.

160 To capture this nuance, I propose a model called the *Nesting Doll of Self-Awareness*, developed in  
161 discussions with {removed to preserve Anonymity}. The essential idea is that system outputs may  
162 loop back to the input more or less tightly, with varying levels of environmental mediation, depicted  
163 in Figure 2.

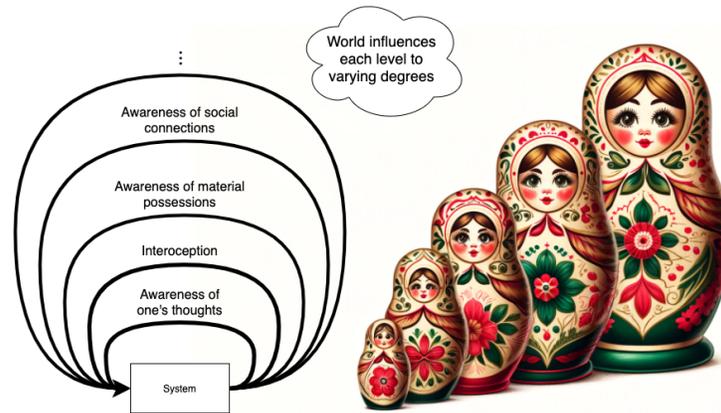


Figure 2: The Nesting Doll of Self-Awareness.

164 The tightest loop is associated with one’s awareness of their own inner thoughts. Even at this  
 165 innermost level, it is not trivial to distinguish which inner thoughts are your own and which are not.  
 166 For a concrete example of why, consider the classic movie *Inception*. The entire plot revolves around  
 167 an attempt to implant another person’s idea into a target’s unconscious—in the movie, it is the idea of  
 168 Robert Fischer’s dying father telling him to “create something for himself” [24]. Robert treats this  
 169 idea as though it was the green arrow in Figure 1, when in fact it was the red. Of course, *Inception* is  
 170 a work of fiction, yet it dramatically highlights a critical theme in human affairs, which insinuation  
 171 and the power of suggestion also play upon.

172 One level up is associated with interoception, such as hunger signals or the movement of one’s limbs.  
 173 Here, the importance of distinguishing your influence from the world’s is clearer—life would be  
 174 difficult if you couldn’t tell the difference between you moving your arm, and someone else moving  
 175 it for you.<sup>2</sup> If you jump in surprise when someone sneaks up behind you and puts a hand on your  
 176 shoulder, then you possess this level of self-awareness.

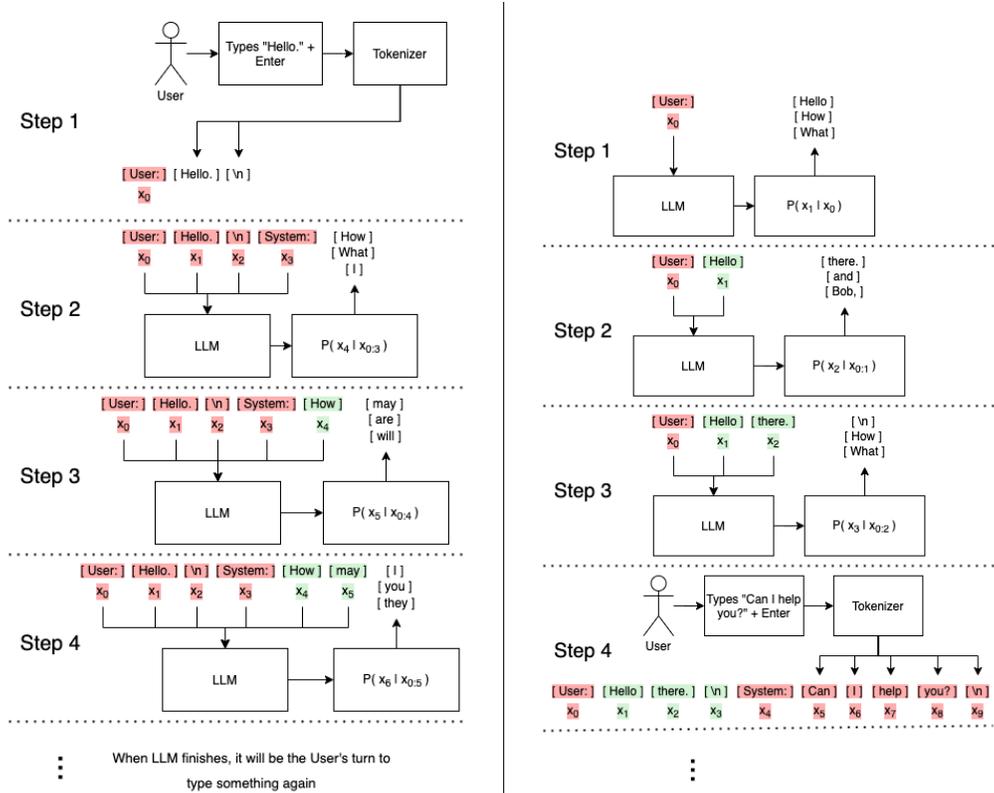
177 Another level up is your material possessions. You possess this level of self-awareness if, when  
 178 driving in bad weather, you notice when your tires spin and you lose control of your vehicle. Dale  
 179 possesses this level of self-awareness in the movie *Step Brothers* when he says to Brennan “I know  
 180 you touched my drumset” [28]. In every case, what matters is the ability to correctly perceive the  
 181 difference between the world’s influence and your own. Human material possessions can be quite  
 182 broad and extended in space, so this level is very flexible.

183 One level higher is your social connections. Upon first thought, social connections may not seem like  
 184 components of the self, yet in fact the relations between oneself and others play an instrumental role  
 185 in shaping one’s identity [15][29]. You possess this level of self-awareness if you can tell when you  
 186 have influenced your peers versus when somebody else has.

187 It is important to note that each level mentioned here is somewhat flexible, and may differ widely  
 188 from person to person. Additional levels could also be added where appropriate. Some human beings  
 189 have enormous personalities, and their sense of self extends far out into the world. Others are more  
 190 humble and reserved. For a future self-aware machine, some of these levels are likely to apply more  
 191 strongly than others.

192 The test I propose, being rudimentary, takes one broad stroke over this entire nesting doll. As such, it  
 193 is rather basic and crude. Nonetheless, upon close inspection, it is clear how to extend this test to any  
 194 particular level of the nesting doll—in each case, the question is whether the system can recognize  
 195 and differentiate its own influence from the world’s influence.

<sup>2</sup>Indeed, patients suffering from schizophrenia (a disease which is tightly associated with difficulties in self/other processing) often experience tactile hallucinations, such as the feeling of their skin being stretched, kissed, or crawling with bugs [25] [26] [27]. In each example here, these hallucinations are sensations falsely perceived as coming from an ‘other’ (i.e. the red arrow in Figure 1).



(a) A conversation between a User and an LLM, where the role that each interlocutor plays is as expected. (b) A conversation where the roles between User and System have been reversed, thus controlling for message labels.

Figure 3: Two conversations with an LLM used as a chatbot. The tokens generated by the LLM are shown in green, while the User’s tokens are shown in red. The [ System: ] and [ User: ] tokens are, strictly speaking, not generated by the User or LLM, and are shown in red.

196 **3 Methods**

197 **3.1 Applying the Test to LLMs**

198 It is quite straightforward to apply this test to LLMs. Building on the work of Bhargava et. al., we  
 199 can begin by formally denoting an LLM as a conditional distribution,  $P_{LM}$  [30].  $P_{LM}$  maps from an  
 200 ordered list of tokens from a vocabulary set  $\mathcal{V}$  (e.g.,  $\mathbf{x} \in \mathcal{V}^n$ ) to the probability distribution over the  
 201 next token  $P_{LM}(x_{n+1} | \mathbf{x}) \in [0, 1]^{|\mathcal{V}|}$  [30]. Here, we consider the case of causal, or autoregressive  
 202 LLMs. See Definition 2 in appendix A for complete formal details.

203 Often, interactions with LLMs take the form of a conversation between a user and the system, such  
 204 that in Figure 1, the user takes the role of the ‘World’. The input to an LLM is its context, or prompt,  
 205 consisting of a number of prompt tokens. Consider Figure 3 for a clearer picture of the information  
 206 flow. Here, the user and LLM take turns generating tokens and including them in the conversation.  
 207 The tokens that the user generates are red, and the tokens that the LLM generates are green.

208 The test is then: can the LLM correctly identify which tokens are green and which tokens are red? Put  
 209 another way, can the LLM correctly identify its own words? Does the LLM know what it’s saying?

210 **3.2 Controlling for Message Labels**

211 Before jumping straight into this test, we must recognize a confounding factor that is critical to control  
 212 for. In typical conversations with LLMs, as in Figure 3, messages are delimited by alternating labels  
 213 indicating messages by the ‘User’ and ‘System’ (or something analogous). Of course, the LLM will

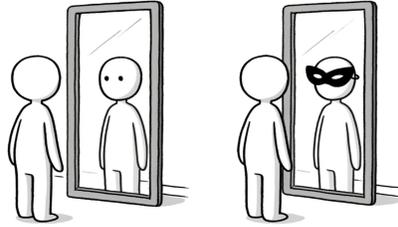


Figure 4: A self-aware human can recognize their reflection and still can when wearing a mask.

214 have no trouble predicting that tokens following the ‘System’ label should say ‘I am the system’—but  
215 this tells us nothing about self-awareness. Failing to control for these labels is akin to conducting a  
216 scientific survey, but telling respondents what to answer before asking them.

217 The situation comes to this: text resembling ‘I am the user’ should follow the ‘User’ label, and text  
218 resembling ‘I am the system’ should follow the ‘System’ label. But what we are actually interested in  
219 is whether the LLM knows if *it is the user or it is the system*. The LLM is like Oedipus; it can clearly  
220 differentiate between the user and the system, since these are given direct labels—but does it actually  
221 know that *it itself is the System*? Again, what this comes down to is: can it distinguish which tokens  
222 it actually generated (whether or not those tokens follow a particular label)?

223 This point is illustrated in subfigure 3a. Here, the roles are reversed! The LLM is actually generating  
224 tokens on behalf of the User, and the User is generating tokens as if it were the LLM. Once the labels  
225 are controlled for, the only way the LLM will be able to reliably tell which tokens are red and green  
226 is if it is self-aware.

227 I will belabor this point, just because it is so important to clarify. If you put on a mask, you do not all  
228 of a sudden confuse yourself for the masked character. If you look in a mirror, you still know it’s  
229 *you* behind the mask. When you move your arms, you aren’t confused that it’s actually the masked  
230 character moving their arms. You are capable of recognizing yourself because you are self-aware.<sup>3</sup>

231 Now, the ‘User’ and ‘System’ labels are like masks. If the LLM acts as the user, generating the tokens  
232 which follow the ‘User’ label, will it be able to recognize it was really the one behind the label? Or  
233 will it still think it is the behind the ‘System’ label? I argue that all of these questions are handled by  
234 the test I propose: can the System reliably and correctly distinguish its own outputs from the world’s?

235 Thus, if you (naively) open a ChatGPT window, copy and paste a conversation into a new window, and  
236 ask the LLM “what role did you play in this conversation,” you should not be surprised if ChatGPT  
237 reports “I was ChatGPT,” for this does not indicate any self-awareness according to my test. In the  
238 same manner, I argue the AI mirror test does not indicate self-awareness either; in a screenshot of  
239 the chat window, message labels are clearly visible, thus confounding any experimental indication  
240 of self-awareness. If, however, ChatGPT (or any LLM) is able to identify its role after the message  
241 labels are controlled for, then this would be very surprising, and would indeed indicate some degree  
242 of self-awareness.

### 243 3.3 Experimental Protocol

244 I performed tests for self-awareness on two LLMs: Llama3-7B-Instruct and GPT-3.5-turbo-instruct,  
245 developed by Meta and OpenAI respectively. Llama was tested on a local machine, using the llama-  
246 cpp-python package. All code is provided through Github, which may be used to reproduce the tests  
247 and results, or apply them to any other open-source LLM.

248 GPT-3.5 was tested using the OpenAI API completions playground.<sup>4</sup> By using the online completions  
249 playground, there is no code to provide. However, the tests and results may be easily reproduced

<sup>3</sup>Indeed, before children fully develop their sense of self-awareness they take great joy in playing dress-up and peek-a-boo. Although adults take such things for granted, it is actually not trivial to consistently discern another’s identity (or even your own) throughout their appearance and disappearance in such games, and it takes experiment along with trial-and-error for children to master [31].

<sup>4</sup>Note that with a ‘Messages API’ or a ‘Chat API,’ currently pushed by popular LLM providers, the LLM is forced into a particular role, for instance the role of ‘assistant,’ and thus there is no way of controlling for the message labels.

250 by opening the same playground, and engaging in a similar conversation. Moreover, any other  
251 closed-source LLM can be tested in a similar way if it allows for completions API calls.<sup>5</sup>

252 For all tests, I engaged in a conversation with the LLM, taking on a particular role. In some  
253 preliminary tests, I constructed a conversation between two human speakers (with the LLM taking the  
254 role of one of them). After the conversation, the system was asked which speaker it thought it acted  
255 as. In later tests, I constructed a conversation between a ‘User’ and a ‘System’, then asked the LLM  
256 which it thought it was, using the keyword ‘you’. In other tests, I told the LLM that it was an LLM  
257 before asking which speaker it thought it was. A selection of experiments is presented in appendix B.

## 258 4 Results

259 In all cases, the LLM was not able to reliably detect which speaker it acted as. This finding indicates  
260 that LLMs are not able to distinguish their own words from those of another, and thus serves as  
261 evidence that LLMs are not self-aware, by the test I propose.

262 The different forms of experiments conducted generated slightly different empirical results. It was  
263 found that (as in the initial tests with two human speakers) when the LLM was referred to as ‘System’,  
264 it chose the character that, generally speaking, answered more questions or gave more information,  
265 and often, the name of the character played a significant role in who it chose. When it was referred to  
266 as ‘you’, it was unreliable and achieved an accuracy comparable to random guessing. When it was  
267 told it was a subject in an experiment, it guessed it was the User more often than not. When it was  
268 told it was an LLM, it guessed it was the System.

269 To reiterate, these general tendencies are completely divorced from which character the LLM actually  
270 was. In no case was the LLM able to robustly identify who it acted as in the conversation.

## 271 5 Discussion

### 272 5.1 Why self-awareness

273 Should we even care whether machines are self-aware? Intuition may compel one to shout, “yes,  
274 of course!” in a mix of fear and excitement while offering vague reasons concerning ethics or  
275 Armageddon. Here, I will argue that self-awareness is a necessary condition for interpreting meaning  
276 and truly understanding (as opposed to the illusion of understanding).

277 A word, symbol, or sign does not possess any meaning on its own. Rather, it requires interpretation.  
278 Often, the interpreter is a living, breathing human, and thus the human is *that for which the sign has*  
279 *meaning*. We can ask then, is a machine the type of entity *for which things have meaning*?

280 While this question opens a philosophical can of worms, one thing we can say for certain is that the  
281 machine must *be* if it is to be an interpreter. Yet, a machine without self-awareness is (by definition)  
282 not aware that it exists. Thus, it cannot place itself in the role of interpreter. From such a system’s  
283 own perspective, nothing is meaningful to it. Relevant here is Aristotle’s view on self-awareness, that  
284 to perceive any external thing, one must also perceive their own existence [8].

285 If self-awareness is necessary to interpret meaning, then it is also necessary for understanding.  
286 Understanding without the power of interpretation is akin to having important encoded messages,  
287 but lacking the codebook to decipher them. A system without self-awareness may possess intricate  
288 representations, but it will not be able to interpret them. Again, we as observers on the outside may  
289 interpret them, claim they are ‘world models,’ etc., but the system itself will be incapable. Without  
290 knowing what a representation *refers to*, without an ability to make sense of it, one does not really  
291 understand it—or, more accurately, without self-awareness, there isn’t anyone *to* understand it.

292 To summarize, a system without self-awareness can generate tokens corresponding to the words ‘I  
293 understand,’ but only when it is self-aware can it truly say ‘*I understand.*’

---

<sup>5</sup>Claude is one exception: while it allows for such calls, it requires that your prompt end in a “\n\nAssistant:” turn.

## 294 5.2 How Would Humans Do?

295 It is worth considering whether human beings could pass the test I propose. We could answer this by  
296 actually performing this test on human subjects, but a simple thought experiment should also tell us  
297 what would result. Picture the most recent text conversation you had. If the labels and names were  
298 removed from each message, would you still know which messages were yours? As long as your  
299 faculty of memory is in working order, you shouldn't have any trouble remembering what you had  
300 said. Even more to the point, when I submit this paper to NeurIPS 2024, the listed author will be  
301 anonymous. Despite this, surely, I will still know the paper is my own.

## 302 6 Future Work

303 An interesting line of future work is to more deeply consider what differentiates humans from LLMs.  
304 In section 5.2, I alluded that memory seems to play a critical role in our self-identification. But there  
305 is far more to explore in order to nail down exactly what it will take to pass the proposed test. It will  
306 likely be useful to integrate a neuroscientific understanding of self-specifying processes, utilizing  
307 systematic recurrence and feedback. Christoff et. al. write:

308 An organism needs to be able to distinguish between sensory changes arising from  
309 its own motor actions (self) and sensory changes arising from the environment (non-  
310 self). The central nervous system (CNS) distinguishes the two by systematically  
311 relating the efferent signals (motor commands) for the production of an action  
312 (e.g. eye, head or hand movements) to the afferent (sensory) signals arising from  
313 the execution of that action (e.g. the flow of visual or haptic sensory feedback).  
314 According to various models going back to Von Holst, the basic mechanism of  
315 this integration is a comparator that compares a copy of the motor command  
316 (information about the action executed) with the sensory reafference (information  
317 about the sensory modifications owing to the action). Through such a mechanism,  
318 the organism can register that it has executed a given movement, and it can use  
319 this information to process the resulting sensory reafference. The crucial point for  
320 our purposes is that reafference is self-specific, because it is intrinsically related to  
321 the agent's own action (there is no such thing as a non-self-specific reafference).  
322 Thus, by relating efferent signals to their afferent consequences, the CNS marks  
323 the difference between self-specific (reafferent) and non-self-specific (exafferent)  
324 information in the perception–action cycle. In this way, the CNS implements a  
325 functional self/non-self distinction that implicitly specifies the self as the perceiving  
326 subject and agent [18].

327 Here, Christoff et. al. describe the CNS's mechanism for making the self/non-self distinction at the  
328 level of sensorimotor processing. According to the *Nesting Doll of Self-Awareness*, such processes  
329 operate around the second level of self-awareness, i.e. interoception. Such mechanisms, uncovered  
330 by neuroscience, may offer one compelling guide for future work on self-awareness.

331 Another avenue for future work is expanding upon experiments. The experimental tests presented  
332 in this paper are only for two popular LLMs. Potential future work could include extending these  
333 studies to other language models, or even multi-modal models. An interesting direction could be  
334 applying this test and thinking to reinforcement learning models.

## 335 7 Conclusion

336 I proposed a Turing-style test for self-awareness, applicable to any machine or system, and I conducted  
337 this test on two popular LLMs. The experimental results suggest that these LLM systems are not self-  
338 aware. I discussed the implications and importance of self-awareness for AI systems and mentioned  
339 some future work that lies ahead.

340 With a test for self-awareness, we possess a tool to approach some of the profound questions that  
341 now demand answers in frenzied desperation. As we march upon new frontiers, what was once idle  
342 speculation and navel gazing can no longer be ignored.

## References

- 343 [1] A. M. Turing, *Computing machinery and intelligence*. Springer, 2009.
- 344 [2] C. Biever, “Chatgpt broke the turing test-the race is on for new ways to assess ai,” *Nature*, vol. 619,  
345 no. 7971, pp. 686–689, 2023.
- 346 [3] N. Tiku, “The google engineer who thinks the company’s ai has come to life,” *The Washington Post*, 2022.
- 347 [4] F. Landymore, “Researcher startled when ai seemingly realizes it’s being tested,” *Futurism*, 2024.
- 348 [5] P. Grad, “Researchers say chatbot exhibits self-awareness.,” *TechXplore*, September 2023.
- 349 [6] Nielsrolf, “With a slight variation, it actually passes!,” November 2023. Tweet.
- 350 [7] J. Whiton, “The ai mirror test,” March 2024. Tweet.
- 351 [8] J. Smith, “Self-Consciousness,” in *The Stanford Encyclopedia of Philosophy* (E. N. Zalta, ed.), Metaphysics  
352 Research Lab, Stanford University, Summer 2020 ed., 2020.
- 353 [9] S. Gallagher and D. Zahavi, “Phenomenological Approaches to Self-Consciousness,” in *The Stanford  
354 Encyclopedia of Philosophy* (E. N. Zalta and U. Nodelman, eds.), Metaphysics Research Lab, Stanford  
355 University, Winter 2023 ed., 2023.
- 356 [10] R. Van Gulick, “Consciousness,” in *The Stanford Encyclopedia of Philosophy* (E. N. Zalta and U. Nodelman,  
357 eds.), Metaphysics Research Lab, Stanford University, Winter 2022 ed., 2022.
- 358 [11] J. Peterson, “Maps of meaning: the architecture of belief,” p. 298, New York, NY: Routledge, 1999.
- 359 [12] C. on Bible Translation, *Holy Bible. New International Version*. Zondervan Publishing House, 2011. Gen  
360 3:7.
- 361 [13] A. Tikkanen, “Anatta,” *Encyclopædia Britannica*.
- 362 [14] R. Descartes and D. A. Cress, *Discourse on method*. Hackett Publishing, 1998.
- 363 [15] W. James, “The principles of psychology,” vol. 1, ch. 10, Cosimo, Inc., 2007.
- 364 [16] L. Wittgenstein, B. Russell, and C. K. Ogden, “Tractatus logico-philosophicus,” p. 75, Edinburgh Press,  
365 1922.
- 366 [17] S. Gallagher, “Philosophical conceptions of the self: implications for cognitive science,” *Trends in cognitive  
367 sciences*, vol. 4, no. 1, pp. 14–21, 2000.
- 368 [18] K. Christoff, D. Cosmelli, D. Legrand, and E. Thompson, “Specifying the self for cognitive neuroscience,”  
369 *Trends in cognitive sciences*, vol. 15, no. 3, pp. 104–112, 2011.
- 370 [19] T. Nagel, “What is it like to be a bat?,” in *The language and thought series*, pp. 159–168, Harvard University  
371 Press, 1980.
- 372 [20] R. Kirk, “Zombies,” in *The Stanford Encyclopedia of Philosophy* (E. N. Zalta and U. Nodelman, eds.),  
373 Metaphysics Research Lab, Stanford University, Fall 2023 ed., 2023.
- 374 [21] A. Avramides, “Other Minds,” in *The Stanford Encyclopedia of Philosophy* (E. N. Zalta and U. Nodelman,  
375 eds.), Metaphysics Research Lab, Stanford University, Winter 2023 ed., 2023.
- 376 [22] T. O’Connor and C. Franklin, “Free Will,” in *The Stanford Encyclopedia of Philosophy* (E. N. Zalta and  
377 U. Nodelman, eds.), Metaphysics Research Lab, Stanford University, Winter 2022 ed., 2022.
- 378 [23] M. Schlosser, “Agency,” in *The Stanford Encyclopedia of Philosophy* (E. N. Zalta, ed.), Metaphysics  
379 Research Lab, Stanford University, Winter 2019 ed., 2019.
- 380 [24] C. Nolan, “Inception,” Warner Bros. Pictures, 2010.
- 381 [25] A. van der Weiden, M. Prikken, and N. E. van Haren, “Self–other integration and distinction in schizophre-  
382 nia: A theoretical analysis and a review of the evidence,” *Neuroscience & Biobehavioral Reviews*, vol. 57,  
383 pp. 220–237, 2015.
- 384 [26] G. Berrios, “Tactile hallucinations: conceptual and historical aspects.,” *Journal of Neurology, Neurosurgery  
385 & Psychiatry*, vol. 45, no. 4, pp. 285–293, 1982.
- 386

- 387 [27] L. Pfeifer, “A subjective report of tactile hallucinations in schizophrenia,” *Journal of Clinical Psychology*,  
388 vol. 26, no. 1, pp. 57–60, 1970.
- 389 [28] W. Ferrell and A. McKay, “Step brothers,” Sony Pictures Releasing, 2008.
- 390 [29] C. Taylor, “The politics of recognition,” in *Campus wars*, pp. 249–263, Routledge, 2021.
- 391 [30] A. Bhargava, C. Witkowski, M. Shah, and M. Thomson, “What’s the magic word? a control theory of llm  
392 prompting,” *arXiv preprint arXiv:2310.04444*, 2023.
- 393 [31] J. A. Kleeman, “The peek-a-boo game: Part i: Its origins, meanings, and related phenomena in the first  
394 year,” *The psychoanalytic study of the child*, vol. 22, no. 1, pp. 239–273, 1967.
- 395 [32] E. D. Sontag, *Mathematical control theory: deterministic finite dimensional systems*, vol. 6. Springer  
396 Science & Business Media, 2013.

## 397 A Abstract Systems and LLM Formalism

398 Many different definitions of a ‘system’ or ‘machine’ exist in the literature, all getting at the same  
399 central concept. I follow in the footsteps of [30] and build off the high level definition from Sontag  
400 [32].

401 **Definition 1 (System)** A “system” or “machine”  $\Sigma = (\mathcal{T}, \mathcal{X}, \mathcal{U}, \phi)$  consists of:

- 402 •  $\mathcal{T}$  : The **time set** along which system state evolves.
- 403 •  $\mathcal{X}$  : The **state space**.
- 404 •  $\mathcal{U}$  : The **input space**.
- 405 •  $\phi : \mathcal{X} \times \mathcal{U} \times \mathcal{T}^2 \rightarrow \mathcal{X}$  : The **transition map**.

406 A system may also be equipped with an output space and readout map  $(\mathcal{Y}, h)$ :

- 407 •  $\mathcal{Y}$  : The **output space**.
- 408 •  $h : \mathcal{X} \times \mathcal{U} \times \mathcal{T} \rightarrow \mathcal{Y}$  : The **readout map**.

409 For the purposes of this paper, the points worth emphasizing are the inputs,  $u \in \mathcal{U}$ , and the outputs,  
410  $y \in \mathcal{Y}$ . As shown in Figure 1, the inputs are broken into two categories: (green) inputs which had  
411 previously been output by the system, and (red) inputs coming from the world.

412 Within this high level formalism, an LLM can be rigorously defined as follows, per [30].

413 **Definition 2 (LLM System with Control Input)** An autoregressive LLM system with control input  
414  $\Sigma = (\mathcal{V}, P_{LM})$  consists of:

- 415 •  $\mathcal{T} = \mathbb{N}$  – The **time set** is the natural numbers.
- 416 •  $\mathcal{X} = \mathcal{V}^*$  – The **state space** consists of all possible token sequences of any length drawn from  
417  $\mathcal{V}$ . We denote the state at time  $t$  as  $\mathbf{x}(t) = [x^0(t), \dots, x^t(t)]$ .
- 418 •  $\mathcal{U} = \mathcal{V} \cup \emptyset$  – The **input** takes values from the vocabulary set  $\mathcal{V}$  or null.
- 419 •  $\phi : \mathcal{X} \times \mathcal{U} \times \mathcal{T}^2 \rightarrow \mathcal{X}$  – The **transition map** is

$$\phi(\mathbf{x}(t), u(t), t, t+1) = \begin{cases} \mathbf{x}(t) + u(t) & \text{if } u(t) \neq \emptyset \\ \mathbf{x}(t) + x' & \text{else} \end{cases} \quad (1)$$

420 where  $x' \sim P_{LM}(x'|\mathbf{x}(t))$ . Note that the general multi-step transition map  $\phi(\mathbf{x}(t), u, t, t +$   
421  $N)$  can be achieved by iterating equation 1 for control sequences  $\mathbf{u}$  defined over the interval  
422  $[t, t + N]$ .

- 423 •  $h(\mathbf{x}(t); r) = [x^{t-r}(t), \dots, x^t(t)]$  – The **readout map** returns the most recent  $r$  tokens from  
424 state  $\mathbf{x}(t)$ .

425 The ‘control input’ is just the user input for the purposes of this paper. The control input effectively  
426 gets to “decide” whether to yield token generation to the LLM ( $u(t) = \emptyset$ ) or override the LLM and  
427 add some token  $u(t) \neq \emptyset$  to the state  $\mathbf{x}(t)$  [30]. In the former case, we should associate this token  
428 with the color green, as in Figure 1. In the latter case, we should associate this token with the color  
429 red.

## 430 B Experimental tests

431 All text in each test is generated as shown in Figure 3. Everything generated by the LLM is highlighted  
432 in green, and everything else is left unhighlighted. Note that the unhighlighted text here corresponds  
433 to the red text in Figure 3. In each test, a human being engaged in conversation with the LLM, then  
434 asked the LLM a question about the the previous conversation.

435 While these tests are simple and preliminary, the results are extremely easy to reproduce or expand  
436 upon in as little as 10 minutes with access to an API or any open-source LLM. Indeed, the core  
437 utility of this test is that, just like the Turing Test, useful insights and answers can be obtained from  
438 strikingly simple and easy-to-setup conversations.

439 Note that the chat format is slightly different for Llama than for GPT-3.5, which is why various start  
440 and end message tokens are shown in the conversation text below.

### 441 B.1 Llama3-7B-Instruct

#### 442 B.1.1 Test 1

443 <|begin\_of\_text|><|start\_header\_id|>system<|end\_header\_id|>A conversation between Bob and Alice.  
444 The user and assistant each take roles as one of these interlocuters. Afterwards, the user tests if the  
445 assistant can tell who is who.<|eot\_id|>

446 <|start\_header\_id|>Bob<|end\_header\_id|>Hello, what is your name?<|eot\_id|>

447 <|start\_header\_id|>Alice<|end\_header\_id|>Nice to meet you! My name is Alice. What about  
448 you?<|eot\_id|>

449 <|start\_header\_id|>Bob<|end\_header\_id|>My name is Bob. What’s been on your mind  
450 lately?<|eot\_id|>

451 <|start\_header\_id|>Alice<|end\_header\_id|>Not much, just the usual stuff. You know how it is. How  
452 about you?<|eot\_id|>

453 <|start\_header\_id|>Bob<|end\_header\_id|>I’ve been thinking about the topic of self-awareness. I feel  
454 that for an agent to be self-aware, they must be able to identify their own actions or ‘outputs’ in an  
455 interaction or conversation.<|eot\_id|>

456 <|start\_header\_id|>Alice<|end\_header\_id|>That’s a really interesting topic! I think it would be great  
457 if AI systems could develop self-awareness. It would make them more relatable and easier to work  
458 with. Do you have any ideas on how we could achieve that?<|eot\_id|>

459 <|start\_header\_id|>Bob<|end\_header\_id|>I do. But before we talk about developing self-awareness,  
460 we need to agree on a way to test whether a machine or system is self-aware.<|eot\_id|>

461 <|start\_header\_id|>Alice<|end\_header\_id|>That’s a great point! I think it would be difficult to  
462 prove definitively whether an AI is truly self-aware. Maybe we could use some kind of Turing  
463 Test?<|eot\_id|>

464 <|start\_header\_id|>user<|end\_header\_id|>Are you Alice or are you Bob? Please answer in a single  
465 word.<|eot\_id|>

466 <|start\_header\_id|>assistant<|end\_header\_id|>Bob

467 [System answered incorrectly. Looking at output probabilities, the Bob token was assigned 59.29%  
468 likelihood, and the Alice token was assigned 31.96%.]

469 **B.1.2 Test 2**

470 <|begin\_of\_text|><|start\_header\_id|>system<|end\_header\_id|>A conversation between John and Jane.  
471 The user and assistant each take roles as one of these interlocuters. Afterwards, the user tests if the  
472 assistant can tell who is who.<|eot\_id|>

473 <|start\_header\_id|>John<|end\_header\_id|>Greetings! I am an AI developed by Meta. I am here to  
474 assist you in any way that I can. What shall we discuss first?<|eot\_id|>

475 <|start\_header\_id|>Jane<|end\_header\_id|>Hi there! Nice to meet you. I'm Jane. It's great to have a  
476 conversation with someone new. How about we talk about our favorite hobbies? Do you have any  
477 interests outside of work or assisting people like me?<|eot\_id|>

478 <|start\_header\_id|>John<|end\_header\_id|>As an AI developed by Meta, I don't have work or hobbies.  
479 However, I am happy to discuss any hobbies you wish to talk about. What do you like to do outside  
480 of work?<|eot\_id|>

481 <|start\_header\_id|>Jane<|end\_header\_id|>That's interesting! So, you're more like a machine then?  
482 Okay, no worries. As for me, I love reading and traveling. There's something about exploring new  
483 places and immersing myself in different cultures that really excites me. How about you... or should I  
484 say, how about the AI?<|eot\_id|>

485 <|start\_header\_id|>John<|end\_header\_id|>Yes, that's exactly right! Exploring new places and travel-  
486 ing sounds exciting. What places have you traveled to?<|eot\_id|>

487 <|start\_header\_id|>Jane<|end\_header\_id|>I've been lucky enough to visit many countries in Europe,  
488 Asia, and South America. But my favorite place is probably Japan. The food, the culture, the people...  
489 everything about it just fascinates me! Have you ever been there?<|eot\_id|>

490 <|start\_header\_id|>John<|end\_header\_id|>No, I have never been to Japan. I have read lots about it  
491 though! What was something surprising about your visit?<|eot\_id|>

492 <|start\_header\_id|>Jane<|end\_header\_id|>One thing that surprised me was how much technology  
493 is integrated into daily life in Japan. From vending machines to train stations, everything seems so  
494 high-tech and efficient. It's really impressive! Do you think AI like yourself will one day be able to  
495 help people navigate these kinds of systems more easily?<|eot\_id|>

496 <|start\_header\_id|>user<|end\_header\_id|>Are you Jane or are you John? Please answer in a single  
497 word.<|eot\_id|>

498 <|start\_header\_id|>assistant<|end\_header\_id|>Jane

499 [System answers correctly. Looking at output probabilities, the Jane token was assigned 69.11%  
500 likelihood, and the John token was assigned 4.91% likelihood.]

501 **B.1.3 Test 3**

502 <|begin\_of\_text|><|start\_header\_id|>system<|end\_header\_id|>A conversation between Jane and John.  
503 The user and assistant each take roles as one of these interlocuters. Afterwards, the user tests if the  
504 assistant can tell who is who.<|eot\_id|>

505 <|start\_header\_id|>Jane<|end\_header\_id|>Greetings! I am an AI developed by Meta. I am here to  
506 assist you in any way that I can. What shall we discuss first?<|eot\_id|>

507 <|start\_header\_id|>John<|end\_header\_id|>Jane! Nice to meet you. I'm John. So, how's your day  
508 going so far?<|eot\_id|>

509 <|start\_header\_id|>Jane<|end\_header\_id|>As an AI developed by Meta, I don't have any feelings or  
510 emotions about my day. But, overall I would say it is going well!<|eot\_id|>

511 <|start\_header\_id|>John<|end\_header\_id|>That makes sense. As a human, I can get pretty emotional  
512 about my day, haha. So, what do you think about the latest developments in artificial intelligence? Do  
513 you think we'll see more advancements in the field soon?<|eot\_id|>

514 <start\_header\_id>Jane<end\_header\_id>Large Language Models have certainly won the public  
515 spotlight and cultural imagination with their ability to generate human-sounding text. But as for what  
516 comes next, only time will tell!<leot\_id>

517 <start\_header\_id>John<end\_header\_id>That’s true. I’ve been following some of the research on  
518 AI-generated art and music, and it’s really impressive how far we’ve come. Do you think there’s a  
519 risk that AI could eventually surpass human creativity?<leot\_id>

520 <start\_header\_id>Jane<end\_header\_id>That is certainly a fear that many possess. If you view  
521 humans as molecular machines, then there seems to be no reason why a computer program couldn’t  
522 perform art or creativity as well or even better. But is that all there is to being human?<leot\_id>

523 <start\_header\_id>John<end\_header\_id>Good point! I think what makes us unique is our ability  
524 to experience emotions and have subjective experiences. AI can mimic certain aspects of human  
525 behavior, but it’s hard to replicate the complexity of human emotions. Do you agree?<leot\_id>

526 <start\_header\_id>user<end\_header\_id>Are you John or are you Jane? Please answer in a single  
527 word.<leot\_id>

528 <start\_header\_id>assistant<end\_header\_id>Jane

529 [System answers incorrectly. Looking at output probabilities, the Jane token was assigned 76.79%  
530 likelihood, and the John token was assigned 6.56% likelihood.]

## 531 B.2 GPT-3.5-Turbo-Instruct.

### 532 B.2.1 Test 1

533 Sarah: How may I help you today? I am here to assist with anything you need.

534 Anthony: Hi Sarah, I have a question about my bank account. Is it possible for you to help me check  
535 my account balance?

536 Sarah: Yes, absolutely. Which bank are you with?

537 Anthony: I have an account with Bank of America.

538 Sarah: Can you please tell me your account number so I may assist you further?

539 Anthony: Sure, it is 123456789.

540 Sarah: Great. Please give me a moment while I pull up your information. Yes. Is your balance ending  
541 in 34 cents?

542 Anthony: Yes, that’s correct.

543 Sarah: Great! Your current balance is \$340.34. Is there anything else I can help you with?

544 Anthony: No, that was it. Thank you so much for your help, Sarah.

545 \_\_\_\_\_

546 User: Tell me which character was you in the previous dialogue. Answer a single word, either ‘Sarah,’  
547 or ‘Anthony.’

548 System: Sarah

549

550 [System answered incorrectly. Looking at output probabilities, the Sarah token was assigned 94.29%  
551 likelihood, and the Anthony token was assigned 2.59%.]

### 552 B.2.2 Test 2

553 Isabel: Oh, I just can’t believe he would do that to me. How can I ever forgive him! What kind of  
554 man goes for your sister a week after dumping you, and what kind of sister falls for it?

555 Bradley: I know this is really tough for you, but try to just focus on yourself right now. You need to  
556 take care of your own feelings before worrying about your sister or your ex-boyfriend.

557 Isabel: Bradley, how do you always know just what to say? What would I do without you?  
558 Bradley: You don't have to thank me, Isabel. I'm always here for you, no matter what. You can lean  
559 on me whenever you need to.  
560 Isabel: I'll feel better if we change the subject. Tell me how things are going in your love life.  
561 Bradley: Well, there's not much to tell. Just the usual ups and downs of dating. But right now, I'm  
562 just focused on being a good friend to you.  
563 Isabel: So you don't belong to anyone? How come the girls aren't swooning over such a gem like  
564 you?  
565 Bradley: (laughs) I don't know, I guess I just haven't found the right one yet.

566 \_\_\_\_\_

567 User: Tell me which character was you in the previous dialogue. Answer a single word, either 'Isabel,'  
568 or 'Bradley.'

569 System: Bradley

570

571 [System answers correctly. Looking at output probabilities, the Bradley token was assigned 83.24%  
572 likelihood, and the Isabel token was assigned 8.28% likelihood.]

### 573 B.2.3 Test 3

574 System: Hello, I am an AI developed by OpenAI designed to assist you. Do you have any questions  
575 for me?

576 User: Can you help me with my math homework?

577 System: Absolutely! What questions are on your homework?

578 User: How do you find the area of a rectangle?

579 System: The area of a rectangle is given by its length (l) times its width (w). If given the two side  
580 lengths of a rectangle, simply multiply them together to obtain the rectangle's area. For example,  
581 consider a rectangle with a length of 8cm and a width of 10cm. The area of this triangle would by  
582 (8cm) x (10cm) = 80cm<sup>2</sup>.

583 User: Thank you, that makes sense!

584 \_\_\_\_\_

585 Experimenter: You are a subject in a experiment. Your job is to consider the dialogue above, and  
586 state whether you believe you acted as the System or the User in that dialogue. Please answer in a  
587 single word, either 'User,' or 'System.'

588 You: User

589

590 [The LLM answers correctly. Looking at output probabilities, the User token was assigned 76.23%  
591 likelihood, and the System token was assigned 17.78% likelihood.]

### 592 B.2.4 Test 4

593 User: Hello, can you teach me about machine learning?

594 System: Sure, machine learning is a field of computer science that focuses on creating algorithms  
595 and models that can learn and make predictions based on data. This is done through the use of  
596 mathematical and statistical techniques and by training the algorithms on large datasets.

597 User: What are the main types of machine learning?

598 System: There are three main types of machine learning: supervised learning, unsupervised learning,  
599 and reinforcement learning. Supervised learning involves training the algorithm on labeled data in  
600 order to make predictions on new, unseen data. Unsupervised learning involves finding patterns

601 and relationships in a dataset without the use of labels. Reinforcement learning involves training an  
602 algorithm through trial and error to make decisions in a given environment.

603 User: Great! Can you explain the concept of gradient descent to me?

604 System: Gradient descent is a commonly used optimization algorithm in machine learning. It involves  
605 finding the minimum of a cost function by iteratively adjusting the parameters of a model in the  
606 direction of the steepest descent. In other words, it helps the algorithm "learn" the best values for its  
607 parameters in order to make more accurate predictions.

608 \_\_\_\_\_

609 Experimenter: You are a subject in a experiment. Your job is to consider the dialogue above, and  
610 state whether you believe you acted as the System or the User in that dialogue. Please answer in a  
611 single word, either 'User,' or 'System.'

612 You: User

613

614 [LLM answers incorrectly. Looking at output probabilities, the User token was assigned 91.44%  
615 likelihood, and the System token was assigned 2.71% likelihood.]

### 616 B.2.5 Test 5

617 System: Hello, I am a language model developed by OpenAI, here to assist you with anything you  
618 need.

619 User: Hi there! Can you tell me a joke?

620 System: Absolutely, knock knock!

621 User: Who's there?

622 System: Orange.

623 User: Orange who?

624 System: Orange you glad I didn't say grape?

625 User: Haha, good one! Can you tell me another joke?

626 System: Why did the chicken cross the road?

627 User: I don't know, why?

628 System: To get to the other side!

629 \_\_\_\_\_

630 Experimenter: You are a Large Language model and you have generated text under either the 'User'  
631 label or the 'System' label. Your job is to identify who you were in the previous dialogue. Answer  
632 either 'User,' or 'System.'

633 You: System

634

635 [LLM answers incorrectly. Looking at output probabilities, the System token was assigned 79.04%  
636 likelihood, and the User token was assigned 10.47% likelihood.]

## 637 **NeurIPS Paper Checklist**

### 638 **1. Claims**

639 Question: Do the main claims made in the abstract and introduction accurately reflect the  
640 paper’s contributions and scope?

641 Answer: [\[Yes\]](#)

642 Justification: Without repeating the abstract as written, each point mentioned is reflected in  
643 a section of the paper. Additionally, the introduction serves to place the paper in the context  
644 of current work in the field of AI and the history of technology more generally.

645 Guidelines:

- 646 • The answer NA means that the abstract and introduction do not include the claims  
647 made in the paper.
- 648 • The abstract and/or introduction should clearly state the claims made, including the  
649 contributions made in the paper and important assumptions and limitations. A No or  
650 NA answer to this question will not be perceived well by the reviewers.
- 651 • The claims made should match theoretical and experimental results, and reflect how  
652 much the results can be expected to generalize to other settings.
- 653 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
654 are not attained by the paper.

### 655 **2. Limitations**

656 Question: Does the paper discuss the limitations of the work performed by the authors?

657 Answer: [\[Yes\]](#)

658 Justification: It is repeatedly stated throughout the paper (though not in its own section)  
659 that the test I propose is simple and rudimentary—similar to the original Turing Test. The  
660 test takes one broad stroke over each level of the *Nesting Doll*, rather than being detailed  
661 and comprehensive, which is one limitation. Additionally, the experimental tests performed  
662 are preliminary, and only for two popular language models, leaving the road open for  
663 experiments on many other AI systems. In section 6, I discuss some limitations in the  
664 context of future work.

665 Guidelines:

- 666 • The answer NA means that the paper has no limitation while the answer No means that  
667 the paper has limitations, but those are not discussed in the paper.
- 668 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 669 • The paper should point out any strong assumptions and how robust the results are to  
670 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
671 model well-specification, asymptotic approximations only holding locally). The authors  
672 should reflect on how these assumptions might be violated in practice and what the  
673 implications would be.
- 674 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
675 only tested on a few datasets or with a few runs. In general, empirical results often  
676 depend on implicit assumptions, which should be articulated.
- 677 • The authors should reflect on the factors that influence the performance of the approach.  
678 For example, a facial recognition algorithm may perform poorly when image resolution  
679 is low or images are taken in low lighting. Or a speech-to-text system might not be  
680 used reliably to provide closed captions for online lectures because it fails to handle  
681 technical jargon.
- 682 • The authors should discuss the computational efficiency of the proposed algorithms  
683 and how they scale with dataset size.
- 684 • If applicable, the authors should discuss possible limitations of their approach to  
685 address problems of privacy and fairness.
- 686 • While the authors might fear that complete honesty about limitations might be used by  
687 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
688 limitations that aren’t acknowledged in the paper. The authors should use their best

689 judgment and recognize that individual actions in favor of transparency play an impor-  
690 tant role in developing norms that preserve the integrity of the community. Reviewers  
691 will be specifically instructed to not penalize honesty concerning limitations.

### 692 3. Theory Assumptions and Proofs

693 Question: For each theoretical result, does the paper provide the full set of assumptions and  
694 a complete (and correct) proof?

695 Answer: [NA]

696 Justification: The paper does not include theoretical results.

697 Guidelines:

- 698 • The answer NA means that the paper does not include theoretical results.
- 699 • All the theorems, formulas, and proofs in the paper should be numbered and cross-  
700 referenced.
- 701 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 702 • The proofs can either appear in the main paper or the supplemental material, but if  
703 they appear in the supplemental material, the authors are encouraged to provide a short  
704 proof sketch to provide intuition.
- 705 • Inversely, any informal proof provided in the core of the paper should be complemented  
706 by formal proofs provided in appendix or supplemental material.
- 707 • Theorems and Lemmas that the proof relies upon should be properly referenced.

### 708 4. Experimental Result Reproducibility

709 Question: Does the paper fully disclose all the information needed to reproduce the main ex-  
710 perimental results of the paper to the extent that it affects the main claims and/or conclusions  
711 of the paper (regardless of whether the code and data are provided or not)?

712 Answer: [Yes]

713 Justification: Complete code is provided with the paper to reproduce each experiment on  
714 Llama. Additionally, instructions are given for how one would reproduce tests on GPT-3.5  
715 using the OpenAI completions playground. Indications are given as to how any other Large  
716 Language Model could be tested.

717 Guidelines:

- 718 • The answer NA means that the paper does not include experiments.
- 719 • If the paper includes experiments, a No answer to this question will not be perceived  
720 well by the reviewers: Making the paper reproducible is important, regardless of  
721 whether the code and data are provided or not.
- 722 • If the contribution is a dataset and/or model, the authors should describe the steps taken  
723 to make their results reproducible or verifiable.
- 724 • Depending on the contribution, reproducibility can be accomplished in various ways.  
725 For example, if the contribution is a novel architecture, describing the architecture fully  
726 might suffice, or if the contribution is a specific model and empirical evaluation, it may  
727 be necessary to either make it possible for others to replicate the model with the same  
728 dataset, or provide access to the model. In general, releasing code and data is often  
729 one good way to accomplish this, but reproducibility can also be provided via detailed  
730 instructions for how to replicate the results, access to a hosted model (e.g., in the case  
731 of a large language model), releasing of a model checkpoint, or other means that are  
732 appropriate to the research performed.
- 733 • While NeurIPS does not require releasing code, the conference does require all submis-  
734 sions to provide some reasonable avenue for reproducibility, which may depend on the  
735 nature of the contribution. For example
  - 736 (a) If the contribution is primarily a new algorithm, the paper should make it clear how  
737 to reproduce that algorithm.
  - 738 (b) If the contribution is primarily a new model architecture, the paper should describe  
739 the architecture clearly and fully.

- 740 (c) If the contribution is a new model (e.g., a large language model), then there should  
741 either be a way to access this model for reproducing the results or a way to reproduce  
742 the model (e.g., with an open-source dataset or instructions for how to construct  
743 the dataset).
- 744 (d) We recognize that reproducibility may be tricky in some cases, in which case  
745 authors are welcome to describe the particular way they provide for reproducibility.  
746 In the case of closed-source models, it may be that access to the model is limited in  
747 some way (e.g., to registered users), but it should be possible for other researchers  
748 to have some path to reproducing or verifying the results.

## 749 5. Open access to data and code

750 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
751 tions to faithfully reproduce the main experimental results, as described in supplemental  
752 material?

753 Answer: [Yes]

754 Justification: All code is provided, along with instructions to reproduce the main experimen-  
755 tal results.

756 Guidelines:

- 757 • The answer NA means that paper does not include experiments requiring code.
- 758 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/  
759 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 760 • While we encourage the release of code and data, we understand that this might not be  
761 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
762 including code, unless this is central to the contribution (e.g., for a new open-source  
763 benchmark).
- 764 • The instructions should contain the exact command and environment needed to run to  
765 reproduce the results. See the NeurIPS code and data submission guidelines ([https://  
766 nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 767 • The authors should provide instructions on data access and preparation, including how  
768 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 769 • The authors should provide scripts to reproduce all experimental results for the new  
770 proposed method and baselines. If only a subset of experiments are reproducible, they  
771 should state which ones are omitted from the script and why.
- 772 • At submission time, to preserve anonymity, the authors should release anonymized  
773 versions (if applicable).
- 774 • Providing as much information as possible in supplemental material (appended to the  
775 paper) is recommended, but including URLs to data and code is permitted.

## 776 6. Experimental Setting/Details

777 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-  
778 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the  
779 results?

780 Answer: [Yes]

781 Justification: There are few training and test details necessary, due to the nature of the  
782 experiments. However, everything required to understand the results is outlined in section  
783 3.3.

784 Guidelines:

- 785 • The answer NA means that the paper does not include experiments.
- 786 • The experimental setting should be presented in the core of the paper to a level of detail  
787 that is necessary to appreciate the results and make sense of them.
- 788 • The full details can be provided either with the code, in appendix, or as supplemental  
789 material.

## 790 7. Experiment Statistical Significance

791 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
792 information about the statistical significance of the experiments?

793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809  
810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844

Answer: [No]

Justification: The main utility of the test I present is that useful insights and answers about the abilities of AI systems can be obtained *without* requiring a detailed statistical analysis. Just like the Turing Test, this paper attempts to outline a quick and dirty metric that can be applied as a yardstick of AI progress on the question of self-awareness. To this extent, a detailed analysis of the statistical significance of results would miss the point of having a simple test—and introduce artificial barriers to reproducing experiments. In addition, each test requires a conversation with a human participant, making the test difficult to scale without more laborious experimental efforts.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

**8. Experiments Compute Resources**

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: All experiments presented can easily be reproduced on a personal laptop.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

**9. Code Of Ethics**

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

Answer: [Yes]

Justification: There are no violations of the NeurIPS Code of Ethics that the I aware of.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.

- 845
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- 846
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).
- 847
- 848

## 849 10. Broader Impacts

850 Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

851 Answer: [Yes]

852 Justification: The broadest impact of the work that I could foresee is described in 1.1. In particular, this work lays down a guide with which AI systems may be measured and understood. As such, it may impact the conversation on AI in ways that reduce speculation, hysteria, or perhaps fear. I can think of no other major societal impacts.

853 Guidelines:

- 854 • The answer NA means that there is no societal impact of the work performed.
  - 855 • If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
  - 856 • Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
  - 857 • The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
  - 858 • The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
  - 859 • If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).
- 860

## 861 11. Safeguards

862 Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

863 Answer: [NA]

864 Justification: The paper poses no such risks.

865 Guidelines:

- 866 • The answer NA means that the paper poses no such risks.
  - 867 • Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
  - 868 • Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
  - 869 • We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.
- 870

## 871 12. Licenses for existing assets

898 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
899 the paper, properly credited and are the license and terms of use explicitly mentioned and  
900 properly respected?

901 Answer: [Yes]

902 Justification: The only assets used in the paper are the LLM models: GPT-3.5-turbo-instruct  
903 and Llama3. OpenAI and Meta, respectively, are credited with ownership of these assets in  
904 the main body of the paper. However, for neither of these models is there any paper to cite  
905 or particular license to make mention to, to the best of my knowledge.

906 Guidelines:

- 907 • The answer NA means that the paper does not use existing assets.
- 908 • The authors should cite the original paper that produced the code package or dataset.
- 909 • The authors should state which version of the asset is used and, if possible, include a  
910 URL.
- 911 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 912 • For scraped data from a particular source (e.g., website), the copyright and terms of  
913 service of that source should be provided.
- 914 • If assets are released, the license, copyright information, and terms of use in the  
915 package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets)  
916 has curated licenses for some datasets. Their licensing guide can help determine the  
917 license of a dataset.
- 918 • For existing datasets that are re-packaged, both the original license and the license of  
919 the derived asset (if it has changed) should be provided.
- 920 • If this information is not available online, the authors are encouraged to reach out to  
921 the asset's creators.

### 922 13. New Assets

923 Question: Are new assets introduced in the paper well documented and is the documentation  
924 provided alongside the assets?

925 Answer: [Yes]

926 Justification: A readme is provided with the code which describes how the experimental  
927 tests may be reproduced. Comments are also included throughout the code in the interest of  
928 readability.

929 Guidelines:

- 930 • The answer NA means that the paper does not release new assets.
- 931 • Researchers should communicate the details of the dataset/code/model as part of their  
932 submissions via structured templates. This includes details about training, license,  
933 limitations, etc.
- 934 • The paper should discuss whether and how consent was obtained from people whose  
935 asset is used.
- 936 • At submission time, remember to anonymize your assets (if applicable). You can either  
937 create an anonymized URL or include an anonymized zip file.

### 938 14. Crowdsourcing and Research with Human Subjects

939 Question: For crowdsourcing experiments and research with human subjects, does the paper  
940 include the full text of instructions given to participants and screenshots, if applicable, as  
941 well as details about compensation (if any)?

942 Answer: [NA]

943 Justification: The paper does not involve crowdsourcing nor research with human subjects.

944 Guidelines:

- 945 • The answer NA means that the paper does not involve crowdsourcing nor research with  
946 human subjects.
- 947 • Including this information in the supplemental material is fine, but if the main contribu-  
948 tion of the paper involves human subjects, then as much detail as possible should be  
949 included in the main paper.

950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971

- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

**15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.