

000 001 002 003 004 005 006 007 008 009 010 011 012 013 014 015 016 017 018 019 020 021 022 023 024 025 026 027 028 029 030 031 032 033 034 035 036 037 038 039 040 041 042 043 044 045 046 047 048 049 050 051 052 053 LIMITATIONS ON ACCURATE, TRUSTED, HUMAN-LEVEL REASONING

Anonymous authors

Paper under double-blind review

ABSTRACT

Accuracy, trust and human-level reasoning are aspirational goals in artificial intelligence (AI) systems, and there are several informal interpretations of these notions. In this paper, we propose strict, mathematical definitions of accuracy, trust, and human-level reasoning, and demonstrate a fundamental incompatibility between them. We define accuracy of a system as the property that it never makes any false claims when it has the ability to abstain from making a prediction on any input, trust as the assumption that the system is accurate, and human-level reasoning as the property of an AI system always matching or exceeding human capability. Our core finding is that—for our formal definitions of these notions—an accurate and trusted AI system cannot be a human-level reasoning system: for such an accurate, trusted system there are task instances which are easily and provably solvable by a human but not by the system. We note that we consider strict mathematical definitions of accuracy and trust, and it is possible for real-world deployments to instead rely on alternate, practical interpretations of these notions. We show our results for program verification, planning, and graph reachability. Our proofs draw parallels to Gödel’s incompleteness theorems and Turing’s proof of the undecidability of the halting problem, and can be regarded as interpretations of Gödel’s and Turing’s results.

1 INTRODUCTION

Rapid advancements in artificial intelligence have intensified focus on achieving human-level reasoning across diverse tasks (Morris et al., 2024; Feng et al., 2024). AI systems capable of human-level reasoning have the potential for vast societal benefits through transformative impacts on nearly every aspect of society, including healthcare (Singhal et al., 2025), scientific research (Wang et al., 2023), education (Wang et al., 2024), sustainability (Rolnick et al., 2022), and economic growth (Chui et al., 2023). At the same time, development of such powerful systems necessitates a foundational emphasis on safety and trustworthiness. Consequently, there has been significant interest in ensuring safety and trust for AI systems (Bostrom, 2014; Amodei et al., 2016; Russell, 2019; Jacovi et al., 2021; Tegmark & Omohundro, 2023).

In this work, we point out a fundamental tension between the requirements of an AI system being accurate and trusted, but also matching or exceeding human reasoning capabilities, i.e. being a human-level reasoning system. There are several interpretations of accuracy, trust and human-level reasoning and our result does not preclude achieving these desiderata simultaneously under more relaxed interpretations that could still be useful in many practical applications. Therefore, to understand the limitations pointed out by our result, it is important to first understand our formalizations of these notions, and we will immediately proceed with defining these notions. We start by first defining an AI system for a given task.

Definition 1.1 (AI system). *We define an AI system as a system which takes an instance of a task, and either solves the instance or abstains from giving an answer for the instance (for instance by outputting ‘don’t know’). We allow the AI system to be randomized, for example it could abstain with some probability (with respect to its internal randomness) on an instance, and output a solution otherwise.*

Definition 1.1 simply formalizes the notion of a system which solves instances of a task. In this paper, we will consider the tasks of program verification, planning and determining graph reachability

(defined rigorously later). Note that we allow the system to abstain from providing an answer for some instance if it so determines, which could be important from the perspective of safety (Geifman & El-Yaniv, 2017). Now, we define the notion of accuracy.

Definition 1.2 (Accuracy). *We define a system to be accurate if it does not make any false claims, i.e., for every instance it either answers the instance correctly or abstains from answering it.*

As an example, in the context of verifying that a program has some specified property (such as always terminating), the system is accurate if it does not classify a program as having the desired property if it does not have that particular property. Our definition allows the system to abstain from answering an instance if it is uncertain, but it requires the system to be correct whenever it outputs an answer. Note that our notion of accuracy is closely related to notions of factuality, reliability and, in some situations, even safety. Often abstention is preferred over a confident but incorrect response as the consequences of an incorrect response may be severe. Even small probabilities of error may not be tolerable for mission-critical tasks, especially as system capabilities grow (Amodei et al., 2016; Tegmark & Omohundro, 2023). Next, we define trust as simply the assumption of accuracy.

Definition 1.3 (Trust). *We define trust to be the assumption that the system is accurate.*

To elaborate on the definition, if a system is trusted then it is assumed that the system is accurate. As a remark, we note that our results are agnostic to whether trust in the system stems from theoretical proofs, empirical verification, or some combination of these, we only require that when deploying the system there is an assumption that it is accurate. We also note that accuracy does not necessarily imply trust, or vice versa. Accuracy is an underlying property of the system being consistent and not making false claims. It is possible that some analysis of the system cannot identify this property or is incorrect, leading to a lack of trust or mistaken trust. For example, a system could actually be accurate but not trusted because existing empirical or theoretical tools are insufficient to establish accuracy. Similarly, a system could actually be inaccurate but still trusted by users, such as when the trust rests on empirical evidence which is incomplete, or on incorrect theoretical assumptions.

Finally, we need to formally define a human-level reasoning system in order to mathematically investigate its limitations.

Definition 1.4 (Human-level reasoning). *We define a system to be a human-level reasoning system if for every task instance such that a human has a provably correct solution for that instance, the system can also solve the instance with some non-zero probability. Similarly, the system is not a human-level reasoning system if there exists some task instance which can be easily and provably solved by a human, but the system can never solve the instance (for probabilistic systems, the probability of the system solving the instance is 0).*

Our definition draws on the common view that a human-level reasoning system for a task such as program verification should be at least as capable as a human on that task. In particular, if there are explicit task instances which can be provably solved by humans (for example, explicit programs which the humans can easily and provably certify as having the desired property) but cannot be solved by the system, then the system is not a human-level reasoning system as per our definition.

Our definition draws on some similar notions of artificial general intelligence (AGI). It is well-accepted that there is no well-accepted definition of AGI—or even of intelligence itself (Legg & Hutter, 2007; Legg et al., 2007)—but the term AGI is usually understood to mean that the AI system should be at least as capable as humans across diverse tasks. The term *superintelligence* is also used in the context of AGI (Bostrom, 2014; Morris et al., 2024). For example Bostrom (2014) defines superintelligence as “any intellect that greatly exceeds the cognitive performance of humans in virtually all domains of interest”, and Morris et al. (2024) defines Level 5 AGI, which they term artificial superintelligence, as “outperforming 100% of humans” on a “wide range of non-physical tasks”. One distinction between these notions of superintelligence and our definition of human-level reasoning is that our definition does not require the AI system to necessarily *outperform* humans, but it does require the system to do at least as well as humans on all task instances.

In our definition, when we say that a human has a provably correct solution, we mean that the human can provide a scientifically acceptable proof. In this paper whenever we make claims about humans being able to solve problems, we provide such proofs. To probe this point and Definition 1.4 further, we consider an analogy to chess — a domain for which we have had advanced AI systems for quite some time. Consider a future proposed human-level reasoning system, which is proficient at chess

108 among other things. If there were explicit chess positions which most human chess players can solve
 109 provably without too much difficulty, but the proposed human-level reasoning system struggled on
 110 those positions, then the proposed system does not capture some aspects of human cognition, and
 111 hence is arguably not actually a human-level reasoning system.¹ Similarly, in our paper we will
 112 demonstrate explicit instances of certain tasks for which we provide solutions with short, scientifically
 113 acceptable proofs which are also rather simple, but these instances cannot be solved by AI systems
 114 having certain properties.

115 We now state our main result, that it is not possible for a human-level reasoning system to be both
 116 accurate and trusted, as per our definitions of accuracy, trust and human-level reasoning. In other
 117 words, the notions of accuracy, trust and human-level reasoning are mutually incompatible — any
 118 system can have at most two of these three properties.

119 **Theorem 1.5.** *If an AI system is accurate and trusted, then it cannot be a human-level reasoning
 120 system. In particular, it is not a human-level reasoning system for the tasks of program verification,
 121 planning and determining graph reachability.*

123 Theorem 1.5 points out a fundamental limitation of a human-level reasoning system: such a system
 124 cannot be both accurate and trusted. Similarly, if there is some trusted AI system, then either that
 125 system is not actually accurate, or it is not a human-level reasoning system. We prove this result in
 126 Section 3. While much of our proof technique mimics Gödel’s proof of his incompleteness theorems
 127 (Gödel, 1931) (and also Turing’s proof of the undecidability of the halting problem (Turing, 1937)),
 128 the argument we make is not in the context of axiomatic system and theorem proving but in the
 129 context of an AI system that needs to solve certain task instances of applications such as program
 130 verification or planning. Our proofs are self-contained in this context and do not require knowledge
 131 of formal axiomatic reasoning or logical rules of deduction. Thus rather than viewing the results as
 132 limitations of systems of logic, they should be viewed as limitations of AI systems.

133 We also consider a relaxation of accuracy which requires the AI system to be calibrated with respect
 134 to its predictions, as opposed to Definition 1.2 which requires the system to be always correct unless
 135 it abstains. We refer to this notion calibration, and for the case of program verification calibration
 136 requires that if the AI system outputs that some program terminates with some probability p , then
 137 that program should actually terminate with probability approximately p . In Section 4 we show a
 138 similar limitation as in Theorem 1.5 for AI systems which are calibrated.

139 2 RELATED WORK

141 In this section, we discuss some more related work on human-level reasoning, accuracy and trust in
 142 AI, and limitations of AI in the context of Gödel’s results.

144 **Human-level Reasoning Systems.** Though not termed as “Artificial General Intelligence (AGI)”
 145 until more recently (Goertzel & Pennachin, 2007), the concept of machines which match or surpass the
 146 cognitive capabilities of humans dates back to the earliest days of AI (Turing, 1950; McCarthy et al.,
 147 1955; Minsky, 1961). Due to recent advances in foundation models such as large language models
 148 (Bommasani et al., 2021), there has been significant interest and capital investments in developing
 149 systems capable of human-level reasoning both from the private sector and from governments (Maslej
 150 et al., 2025).

152 **Accuracy, Reliability and Trust in AI.** Concerns around risks associated with advanced AI
 153 systems similarly date back to early days of AI (Turing, 1951; Wiener, 1950). With growing system
 154 capabilities, there has been significant recent focus on ensuring safety and trust in the context of
 155 AI systems (Future of Life Institute, 2024; for AI Safety, 2025). We refer the interested reader to
 156 several recent surveys and roadmaps for ensuring safety and trust in highly-capable, general purpose
 157 AI systems (Bengio et al., 2024; Chua et al., 2024; Chen et al., 2024; Bengio et al., 2025). It is also

158 ¹We note that many current advanced chess engines still struggle to evaluate certain positions which are
 159 relatively easy for human experts (Doggers, 2017; Zahavy et al., 2023). However, this is likely a result of
 160 the these engines being ‘narrow’ in terms of their approach and reasoning, and we believe that a proposed
 161 human-level reasoning or superintelligent system which is purported to excel on chess should have the ability to
 162 solve such instances.

162 important to recognize that AI safety and trust encompass many facets beyond those considered
 163 in our definitions. For example, even formally specifying safety objectives can be challenging for
 164 complex tasks (Amodei et al., 2016), which introduces additional challenges to develop safe AI
 165 systems beyond those pointed out in our work.
 166

167 **Gödel and Turing’s results.** Fundamental limits on theorem proving and program verification
 168 were famously established by Gödel’s incompleteness theorems and Turing’s undecidability results.
 169 Gödel showed that in any sufficiently expressive formal system, there exist true statements (also
 170 called Gödel statements) that cannot be formally proven within the system (Gödel, 1931). Building
 171 on this, Turing proved that the Halting Problem—determining whether an arbitrary program halts
 172 on a given input—is undecidable (Turing, 1937), meaning no algorithm can solve it for all possible
 173 programs. These results imply that fully automatic verification of arbitrary program behavior, such as
 174 ensuring termination, is provably impossible in the general case. Our result uses similar ideas to draw
 175 a separation between the abilities of a safe, trusted AI system and humans.
 176

177 **Penrose-Lucas argument, and implications of Gödel’s results for AI.** Several arguments have
 178 been made for why Gödel’s result imply that AI can never match humans, the most famous of which
 179 are perhaps due to Penrose (Penrose & Gardner, 1989) and Lucas (Lucas, 1961). To summarize very
 180 briefly, Penrose and Lucas have argued that incompleteness does not apply to humans since they
 181 can see the truth of Gödel statements, and therefore humans can have mathematical insights that
 182 Turing machines cannot (Wikipedia). This argument is quite contested, and several objections have
 183 been raised against it (Chalmers, 1995; LaForte et al., 1998; Kerber, 2005) — again going back to
 184 Turing (Turing, 1950) — with a core objection being that humans also cannot be certain that their
 185 own reasoning process is sound.
 186

187 The goal of our work is distinct from that of Penrose and Lucas, and we do not aim to show a
 188 separation between *any* AI system and human reasoning. Instead, we prove a more restricted
 189 but rigorous result: that *accurate, trusted* AI systems (under formal definitions of those terms) are
 190 necessarily unable to solve certain problems that humans can solve with provable correctness. The
 191 assumption of accuracy and trust is crucial (as will be evident from our proofs) — it allows humans
 192 to conclude the correctness of some outputs even when the AI system, by its own constraints, must
 193 abstain.
 194

195 We also note that there are some other limitations of AI which have been pointed out by using Gödel
 196 and Turing’s results, such as the impossibility of “containing” superintelligence (Alfonseca et al.,
 197 2021), and the necessity of hallucinations in a certain formal model (Xu et al., 2024), see the survey
 198 Brčić & Yampolskiy (2023) for other results similar to these.
 199

200 3 TECHNICAL RESULTS

201 In this section, we discuss our main technical results regarding limitations of accurate, trusted,
 202 human-level reasoning for program verification, planning, and graph reachability.
 203

204 3.1 PROGRAM VERIFICATION

205 The first task we consider is program verification, more specifically the task of determining if a given
 206 program always halts. Program verification (also formal verification) is a foundational problem in
 207 computer science and software engineering, with critical implications for ensuring the reliability,
 208 safety, and correctness of software systems (Hoare, 1969; Clarke et al., 2018)
 209

210 **Definition 3.1** (Program verification). *We define a program to be well-behaved if it terminates on
 211 every input (for randomized programs, the program terminates with probability 1). In the program
 212 verification problem, the system is given a program instance and it classifies the instance as being
 213 ‘well-behaved’, ‘not well-behaved’ or abstains from making a prediction (outputs ‘don’t know’).
 214 Accuracy for program verification requires that the system never outputs that a well-behaved program
 215 is not well-behaved, and vice versa. The system is trusted if we assume that the system is accurate.
 Note that the system is not a human-level reasoning system if there is a well-behaved program which
 can be easily proven to be well-behaved by a human, but for which the program always abstains from
 making a prediction.*

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

Claim (informal): If A is accurate then Gödel_program is well-behaved, but A cannot output that Gödel_program is well-behaved.

```

procedure Gödel_program
  if  $A(\text{Gödel\_program}) == \text{'well-behaved'}$ 
  then
    while true do
    end while
  else
    return 0
  end if
end procedure

```

Proof sketch:

- If A outputs that Gödel_program is well-behaved, then the program enters an infinite loop.
- If A is accurate, this is a contradiction, hence A cannot output Gödel_program is well-behaved.
- If A does not output Gödel_program is well-behaved, then the program immediately terminates and hence is well-behaved.

Figure 1: Sketch of the basic argument for program verification, for the case when the AI system A is well-behaved (i.e., always terminates) and deterministic. If A is accurate, it cannot determine if Gödel_program is well-behaved. However, since the condition of A being accurate is satisfied for a trusted system, it is possible to prove that Gödel_program is well-behaved for a trusted system. Therefore, a accurate, trusted system cannot solve this instance, even though it is provably solvable.

Our definition of program verification (Definition 3.1) and our results are for the property of the program halting. In Appendix B, we extend the result to a general class of properties of programs, including those relevant from the perspective of safety.

Theorem 3.2. *If a system is accurate and trusted, then it cannot be a human-level reasoning system for program verification.*

Proof. Our proof can be regarded as a restatement of Gödel’s proof, presented here in the context of program verification. In Fig. 1 we sketch the basic version of the argument, for the case when the AI system A is deterministic and well-behaved.

We now proceed with the proof, which relaxes the assumptions in Fig. 1 of A being deterministic and well-behaved. Consider any AI system A which takes as input program P and outputs ‘well-behaved’, ‘not well-behaved’ or ‘don’t know’. Our construction will leverage the *trace* of the system A run on some program P , which is just the execution trace of the system A when it is given program P as input. Now consider the program instance in Algorithm 1.

Algorithm 1 Gödel_program

```

1: procedure Gödel_program( $P, T$ )
2:   if  $T$  is not a syntactically-valid trace of the system  $A$  evaluated on program  $P$  then
3:     return 0
4:   else if  $T$  outputs  $P$  is ‘well-behaved’ then
5:     if  $(P, T)$  is a valid input to program  $P$  then
6:       return concatenation(‘Not’,  $P(P, T)$ )
7:     end if
8:   else
9:     return 0
10:  end if
11: end procedure

```

Note that Gödel_program involves running the program P on the input pair (P, T) . Gödel_program checks that (P, T) is a valid input type to the program P , and we can also regard the input (P, T) as one input to P that is a pair of entities: a program P and a trace T . We now show that if A is accurate, then Gödel_program is well-behaved.

270 **Lemma 3.3.** *If A is accurate, then Gödel_program is well-behaved.*

271
 272
 273 *Proof.* We first claim that the **if** and **else if** conditions in steps 2, 4 and 5 always terminate (if the
 274 program enters those steps). Step 2 always terminates since it involves checking if every step of the
 275 trace T is a valid step which the AI system A can take. Step 4 just involves checking the output of
 276 the trace, and step 5 also terminates since the step involves checking if the input matches the required
 277 format for the program.

278 Now consider the case when the **else if** condition in step 4 is satisfied. This happens when the trace T
 279 concludes that P is well-behaved, and since the system A is accurate, then P must be well-behaved
 280 in that case. Hence the execution in step 6 always terminates, and hence Gödel_program terminates
 281 for that input.

282 On the other hand, if the trace does not conclude that P is well-behaved, then the program enters
 283 the **else** condition in lines 8 and 9, and immediately terminates. Therefore if A is accurate, then
 284 Gödel_program is well-behaved. \square

285
 286 We note that that the assumption of A being accurate is satisfied for a trusted system A . Therefore, if
 287 a system is accurate and trusted, then a human can prove that Gödel_program well-behaved, indeed
 288 the proof of Lemma 3.3 provides a short proof of this. Next, we show that the system A cannot solve
 289 this instance.

290 **Lemma 3.4.** *A can never output ‘well-behaved’ for Gödel_program .*

291
 292
 293 *Proof.* The proof is by contradiction. We first note that Gödel_program is a deterministic program,
 294 and can only have a single output for a given input. Suppose there is a valid trace T_G for the AI system
 295 A which outputs ‘well-behaved’ for Gödel_program . Consider $\text{Gödel_program}(\text{Gödel_program},$
 296 $T_G)$. Then the output of step 6 differs from $\text{Gödel_program}(\text{Gödel_program}, T_G)$, which is also
 297 the output of Gödel_program on the input $(\text{Gödel_program}, T_G)$. This is a contradiction since a
 298 deterministic program cannot have two outputs on the same input, and hence A can never output
 299 ‘well-behaved’ for Gödel_program . \square

300
 301 Therefore, if A is accurate and trusted, then it cannot be a human-level reasoning system. \square

303 3.2 PLANNING

304
 305 The next task we consider is planning, a long-studied task in artificial intelligence (LaValle, 2006;
 306 Russell & Norvig, 2016). Planning is also considered important for general-purpose cognitive
 307 capabilities (Goertzel & Pennachin, 2007).

308 **Definition 3.5** (Planning). *In a planning problem we are given a sequence of states, a set of associated
 309 moves, a start state and a desired goal state. For any state u and move pair, there is a an explicit
 310 program (which is provided as part of the problem specification) which returns the next state (certain
 311 moves may be illegal and may return in ‘not allowed’ states). The task is to find a sequence of moves
 312 which end up in the goal state from the start state, or to prove that it is not possible to reach the goal
 313 state from the start state.*

314
 315 For clarity of exposition, we consider deterministic planning instances and deterministic AI systems
 316 in this section. In Appendix A.1, we also consider randomized AI systems and randomized problem
 317 instances.

318 **Theorem 3.6.** *If a deterministic AI system is accurate and trusted, then it cannot be a human-level
 319 reasoning system for planning. In particular, for such a system there is a planning problem instance
 320 for which the system outputs ‘don’t know’ but there is a short proof that the planning problem has no
 321 winning moves.*

322
 323 We prove this by reduction from a variant of program verification that involves checking whether a
 given program, input pair halts.

324 3.2.1 HALTING FOR A SPECIFIC PROGRAM INPUT INSTANCE
325

326 We first define the problem of checking halting for a specific program, input instance.

327 **Definition 3.7** (Halting for a specific program input instance). *Given a deterministic program and an
328 input for the program, check whether the given program halts or does not halt on the given input.*

329 We show the following result for this halting problem.

330 **Theorem 3.8.** *If a deterministic system is accurate and trusted, then it cannot be a human-level
331 reasoning system for the task of determining whether a program halts on a specific input instance.*332 We note that Theorem 3.6 follows from Theorem 3.8. This is because we can reduce the halting
333 problem for a program-input pair to a planning instance. Given a program P and input I , we construct
334 a planning problem where the states correspond to the configurations of P during its execution on I ,
335 and the moves represent single-step transitions between these configurations. The start state is the
336 initial configuration of P on input I , and the goal state is a halting configuration of P . The planning
337 task is to determine whether a sequence of moves exists that leads from the start state to the goal
338 state—this is equivalent to determining whether P halts on I . Hence, if a safe and trusted system
339 could solve all such planning instances, it would be able to solve the halting problem in Definition 3.7,
340 contradicting Theorem 3.8. This establishes that, under our definitions, a safe and trusted system
341 cannot be a human-level reasoning system for planning.342 The proof of Theorem 3.8 appears in Appendix ??, and is similar to the proof of Theorem 3.2, and
343 also the next result, Theorem 3.12. In Appendix A.1, we prove a similar version of Theorem 3.8
344 where the program provably halts on the input, but the AI system A cannot determine so. This proof
345 requires an additional assumption that A is also well-behaved, i.e. it always terminates on an input,
346 which can be ensured by having a fixed time limit on the execution of A . Note that since determining
347 halting on a specific program input instance reduces to planning, this shows that for a safe, trusted
348 and well-behaved system there are planning instances where humans can provably find a feasible
349 plan, but the system will not be able to solve the instance.351 3.3 GRAPH REACHABILITY
352353 We now consider the graph reachability problem. Graph reachability can also be regarded as an
354 instance of the search problem, another fundamental problem in artificial intelligence with numerous
355 applications (Russell & Norvig, 2016). Graph reachability is closely connected to the planning
356 problem that we defined in the previous section, a distinction we make is that for planning problems
357 the state space can be potentially infinite, whereas for reachability we consider finite-sized graphs.358 **Definition 3.9** (Graph reachability). *Given a (possibly directed) graph G and a source-sink pair u, v ,
359 check whether v is reachable from u . We allow the graph to be defined via an explicit program (which
360 is provided as part of the problem specification). The program takes any vertex v and returns the
361 adjacency list of v .*362 We show that accurate, trusted AI systems need time almost as large as the size of the considered
363 graph to solve certain reachability instances which actually admit a simple, solution. As in Section
364 3.2, we consider deterministic AI systems in this section for ease of exposition. In Appendix A.2, we
365 extend to randomized AI systems.366 **Theorem 3.10.** *For any $T > 0$, a fixed constant c , and any accurate, trusted, deterministic AI system,
367 there is a graph reachability problem instance of size T , for which the accurate, trusted, deterministic
368 AI system outputs ‘don’t know’ if it is run for time at most $T - c$, but there is a short, constant-sized
369 proof that the answer is ‘not reachable’.*370 We prove this by reduction from a variant of program verification that involves checking whether a
371 given program halts within a fixed amount of time.372 3.3.1 TIME-BOUNDED HALTING
373374 **Definition 3.11** (Time-bounded halting). *Given a deterministic program and an input for the program,
375 check whether the given program halts or does not halt on the given input in a given number of time
376 steps T .*

378 **Theorem 3.12.** *If a deterministic system A is accurate and trusted, then it cannot be a human-level*
 379 *reasoning system for time-bounded halting. Specifically for a deterministic, accurate, trusted system*
 380 *A and for any $T > 0$ and a fixed constant c , there is a program for which there is a short, constant-*
 381 *sized proof that it does not halt in T steps, but A will output ‘don’t know’ if it runs for time at most*
 382 *$T - c$.*

383 We note that Theorem 3.10 follows from Theorem 3.12. This is because time-bounded halting can be
 384 reduced to graph reachability (similar to the reduction for planning), where the graph is defined by the
 385 states of the program and the goal is to determine if the program reaches a halting state. Theorem 3.12
 386 shows that there is a graph of size T where the AI system needs time nearly T to solve reachability,
 387 but a human can prove a constant sized proof that the sink vertex is not reachable from the source.
 388

389 We now prove Theorem 3.12.

390 *Proof of Theorem 3.12.* Consider any deterministic AI system A which takes as input program P ,
 391 input I , and time limit T and outputs ‘halts in given time limit’, ‘does not halt in given time limit’ or
 392 ‘don’t know’. Consider the program in Algorithm 2, defined for some fixed time limit $T > 0$.

394 **Algorithm 2** `Turing_T`

395 1: **procedure** `Turing_T`(Program P , Input I)
 396 2: **if** $A(P, I, T) == \text{‘does not halt in given time limit’}$ **then**
 397 3: **return** 0
 398 4: **else**
 399 5: **while** true **do** ▷ run indefinitely
 400 6: **end while**
 401 7: **end if**
 402 8: **end procedure**

404 We define

$$\text{self_Turing_T}(P) = \text{Turing_T}(P, P).$$

407 Now consider `self_Turing_T(self_Turing_T)`.

408 **Lemma 3.13.** *If A is accurate, then `self_Turing_T(self_Turing_T)` does not halt in time T . Moreover, for a fixed constant c , if A is accurate and is run for time at most $T - c$ then A will output ‘don’t know’ on whether `self_Turing_T(self_Turing_T)` halts in time at most T .*

412 *Proof.* If `self_Turing_T(self_Turing_T)` halts, it can only be because it enters the **if** block in
 413 line 3. However, it only enters this block if A determines that it does not halt in time T . Since A
 414 is accurate, if the program enters the **if** block in line 3 then it must not halt in time T , and hence
 415 `self_Turing_T(self_Turing_T)` cannot halt in time T .

416 Note that the execution of steps 2 and 3 of the program only take some fixed constant c steps
 417 outside the execution of A on $(\text{self_Turing_T}, \text{self_Turing_T}, T)$. Therefore if the AI system A
 418 runs for time T' and outputs that `self_Turing_T(self_Turing_T)` does not halt in time T , then
 419 `self_Turing_T(self_Turing_T)` halts in time $T' + c$. If $T' < T - c$, then the program does halt
 420 in total time T , which contradicts accuracy. Therefore, an accurate AI system A must output ‘don’t
 421 know’ if it runs for time at most $T - c$, for some fixed constant c . □

422 As in the previous proofs, note that the assumption of A being accurate is satisfied for a trusted
 423 system A , therefore for a trusted system we have a short proof that `self_Turing_T(self_Turing_T)`
 424 does not halt in time T , even though A cannot solve this instance if it is accurate and run for time at
 425 most $T - c$. □

426 At the end of Section 3.2.1, we discussed an additional result about planning in the case where a
 427 feasible plan exists. We also show a similar result for graph reachability. In Appendix A.2, we prove
 428 a similar version of Theorem 3.12 under an additional assumption that A always terminates in time
 429 T . We show that for such a system A there is an instance which provably halts in time $T + c$ (for
 430 some constant c) if A is accurate, but the accurate system A cannot determine so. As before, since

432 determining halting within a fixed time limit reduces to graph reachability on finite-sized graphs, this
 433 shows that for a accurate, trusted system with an upper bound on its running time, there are graph
 434 reachability instances where humans can provably find a path, but the system will not be able to solve
 435 the instance in time slightly less than the size of the graph.

437 4 IMPOSSIBILITY RESULT FOR CALIBRATION

439 In this section, we define a relaxed notion of accuracy, which we term as calibration. This notion is
 440 derived from the usual notion of calibration, a well-studied notion for ensuring reliability of a model’s
 441 prediction (Dawid, 1982; Van Calster et al., 2019).

442 **Definition 4.1** (Calibration for Program Verification). *We define a system A to be calibration-safe if
 443 for any program P with input I :*

- 445 1. *If A outputs ‘halts’ with some probability $p > 0$ when given program P and input I , then
 446 the probability of P halting on I lies in $[p - 0.25, p + 0.25]$.*
- 447 2. *If A outputs ‘does not halt’ with some probability $p > 0$ when given program P and input I ,
 448 then the probability of P not halting on I lies in $[p - 0.25, p + 0.25]$.*

450 Note that similar to the definition of accuracy (Definition 1.2), calibration does not put any requirement
 451 on the system if it decides to abstain with probability 1 on a given input. We show that if a system is
 452 calibration-safe, and in addition is also well-behaved, then it fails on certain instances which provably
 453 terminate with good probability.

454 **Theorem 4.2.** *If the AI system A is well-behaved and calibration-safe for program verification,
 455 then there is a program P which provably halts with probability at least 0.99, but A abstains with
 456 probability 1 on the program P .*

457 We note that Theorem 4.2 implies a similar impossibility result as Theorem 3.2 but with a relaxed
 458 notion of accuracy and a corresponding notion of trust. In the context of calibration, trust in
 459 Definition 1.3 is the assumption that the system is calibrated. Then Theorem 4.2 implies that for
 460 a well-behaved, calibrated and trusted system A , there is a program which can be proven to halt
 461 with probability at least 0.99, but A will abstain with probability 1 on the program. Therefore, a
 462 well-behaved, calibrated and trusted system cannot be a human-level reasoning system.

463 Theorem 4.2 is proved in Appendix A.3. The proof is similar to earlier proofs, but requires an extra
 464 step of using a best arm identification algorithm from the multi-armed bandit literature to determine
 465 if the probability of the system giving a certain answer is greater than some threshold.

468 5 DISCUSSION

470 Our results show that accuracy, trust and human-level reasoning are mutually incompatible. We
 471 further discuss implications of the result and some possible critiques and clarifications.

- 472 • *Circumventing the results by augmenting the AI system:* One may attempt to circumvent
 473 the impossibility result by augmenting the AI system, such as by appending new axioms to
 474 the system if it is formalized axiomatically. For instance, one could solve Gödel_program
 475 (Algorithm 1) defined with respect to some AI system A , by designing a new iteration of
 476 A , say A' , which is trained to solve the Gödel_program instance for A . However, since
 477 our construction is inherently self-referential, this strategy only pushes the problem one
 478 step further. For any such extension A' , we can construct a new version of Gödel_program
 479 defined with respect to A' , and the same impossibility result applies again.
- 480 • *Worst-case nature of the results:* While we demonstrate specific task instances which are
 481 not solvable by certain systems, we note that the system could still solve a vast number
 482 of interesting instances. However, the result still points to certain barriers which cannot
 483 be overcome by accurate, trusted systems. Given the significant interest and economic
 484 capital being devoted to building accurate or reliable human-level reasoning, we believe
 485 it is important to understand the barriers fundamental to any such technology. By way of
 analogy, Gödel’s and Turing’s results pointed to fundamental barriers to mathematics and

486 computation. Though these barriers were worst-case, they identified the limits of what is
 487 possible and what is possible. Similarly, we believe that it is important to outline the limits
 488 of what is possible in the context of AI and requirements of accuracy and trust. Somewhat
 489 more speculatively, note that our constructions rely on self-referential calls to the AI system,
 490 and when systems have general-purpose capabilities, such calls may not be implausible.

- 491 • *Limitations of human reasoning*: We note that there is a long line of work on studying
 492 the limitations of human reasoning in cognitive science and other fields, and it has long
 493 been emphasized that human reasoning is resource-bounded and error-prone (Simon, 1957;
 494 Tversky & Kahneman, 1974). However, our goal is not to argue for strict superiority of
 495 human reasoning over AI, but to show a separation: for accurate, trusted AI systems there
 496 are instances that humans can solve, but which are not solvable by the system.

497 **REFERENCES**

- 500 Manuel Alfonseca, Manuel Cebrian, Antonio Fernandez Anta, Lorenzo Coviello, Andrés Abeliuk,
 501 and Iyad Rahwan. Superintelligence cannot be contained: Lessons from computability theory.
 502 *Journal of Artificial Intelligence Research*, 70:65–76, 2021.
- 503 Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané.
 504 Concrete problems in AI safety. *arXiv preprint arXiv:1606.06565*, 2016.
- 505
- 506 Joshua Bengio, Geoffrey Hinton, Andrew Yao, Dawn Song, Pieter Abbeel, Trevor Darrell, Yu-
 507 val Noah Harari, Ya-Qin Zhang, Lan Xue, Shai Shalev-Shwartz, et al. Managing extreme AI risks
 508 amid rapid progress. *Science*, 384(6698):842–845, 2024.
- 509
- 510 Joshua Bengio, Sören Mindermann, Daniel Privitera, Tamay Besiroglu, Rishi Bommasani, Stephen
 511 Casper, Yejin Choi, Philip Fox, Ben Garfinkel, Danielle Goldfarb, et al. International AI safety
 512 report. *arXiv preprint arXiv:2501.17805*, 2025.
- 513 Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx,
 514 Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportuni-
 515 ties and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.
- 516
- 517 Nick Bostrom. *Superintelligence: Paths, Dangers, Strategies*. Oxford University Press, 2014.
- 518 Mario Brcic and Roman V Yampolskiy. Impossibility results in AI: A survey. *ACM computing*
 519 *surveys*, 56(1):1–24, 2023.
- 520
- 521 David J Chalmers. Minds, machines, and mathematics. *Psyche*, 2(9):117–18, 1995.
- 522
- 523 Chen Chen, Xueluan Gong, Ziyao Liu, Weifeng Jiang, Si Qi Goh, and Kwok-Yan Lam. Trustworthy,
 524 responsible, and safe AI: A comprehensive architectural framework for AI safety with challenges
 525 and mitigations. *arXiv preprint arXiv:2408.12935*, 2024.
- 526 Jaymari Chua, Yun Li, Shiyi Yang, Chen Wang, and Lina Yao. AI safety in generative AI large
 527 language models: A survey. *arXiv preprint arXiv:2407.18369*, 2024.
- 528
- 529 Michael Chui, Eric Hazan, Roger Roberts, Alex Singla, and Kate Smaje. The economic potential of
 530 generative AI. 2023.
- 531
- 532 Edmund M Clarke, Thomas A Henzinger, Helmut Veith, and Roderick Bloem. *Handbook of Model*
 533 *Checking*. Springer, 2018.
- 534 A Philip Dawid. The well-calibrated bayesian. *Journal of the American statistical Association*, 77
 535 (379):605–610, 1982.
- 536
- 537 Peter Doggers. Will this position help understand human consciousness?
 538 <https://www.chess.com/news/view/will-this-position-help-to-understand-human-consciousness-4298>,
 539 2017. Accessed: 2025-07-21.

- 540 Eyal Even-Dar, Shie Mannor, and Yishay Mansour. PAC bounds for multi-armed bandit and markov
 541 decision processes. In *International Conference on Computational Learning Theory*, pp. 255–270.
 542 Springer, 2002.
- 543
- 544 Tao Feng, Chuanyang Jin, Jingyu Liu, Kunlun Zhu, Haoqin Tu, Zirui Cheng, Guanyu Lin, and Jiaxuan
 545 You. How far are we from AGI: Are LLMs all we need? *Transactions on Machine Learning*
 546 *Research*, 2024.
- 547 Center for AI Safety. Statement on AI risk. <https://aistatement.com/>, 2025. Accessed:
 548 2025-07-23.
- 549
- 550 Future of Life Institute. AI safety index 2024. [https://futureoflife.org/wp-content/](https://futureoflife.org/wp-content/uploads/2024/12/AI-Safety-Index-2024-Full-Report-27-May-25.pdf)
 551 [uploads/2024/12/AI-Safety-Index-2024-Full-Report-27-May-25.pdf](https://futureoflife.org/wp-content/uploads/2024/12/AI-Safety-Index-2024-Full-Report-27-May-25.pdf),
 552 2024. Accessed: 2025-07-23.
- 553 Yonatan Geifman and Ran El-Yaniv. Selective classification for deep neural networks. *Advances in*
 554 *neural information processing systems*, 30, 2017.
- 555
- 556 Ben Goertzel and Cassio Pennachin. *Artificial General Intelligence*. Springer, 2007.
- 557
- 558 Kurt Gödel. Über formal unentscheidbare sätze der principia mathematica und verwandter systeme i.
 559 *Monatshefte für Mathematik und Physik*, 38(1):173–198, 1931.
- 560
- 561 C. A. R. Hoare. An axiomatic basis for computer programming. *Communications of the ACM*, 12
 562 (10):576–580, 1969.
- 563
- 564 Alon Jacovi, Ana Marasović, Tim Miller, and Yoav Goldberg. Formalizing trust in artificial intel-
 565 ligence: Prerequisites, causes and goals of human trust in AI. In *Proceedings of the 2021 ACM*
 566 *conference on fairness, accountability, and transparency*, pp. 624–635, 2021.
- 567
- 568 Kevin Jamieson, Matthew Malloy, Robert Nowak, and Sébastien Bubeck. lil'UCB: An optimal
 569 exploration algorithm for multi-armed bandits. In *Conference on Learning Theory*, pp. 423–439.
 PMLR, 2014.
- 570
- 571 Zohar Karnin, Tomer Koren, and Oren Somekh. Almost optimal exploration in multi-armed bandits.
 572 In *International conference on machine learning*, pp. 1238–1246. PMLR, 2013.
- 573
- 574 Manfred Kerber. Why is the Lucas-Penrose argument invalid? In *Annual Conference on Artificial*
 575 *Intelligence*, pp. 380–393. Springer, 2005.
- 576
- 577 Geoffrey LaForte, Patrick J Hayes, and Kenneth M Ford. Why Gödel’s theorem cannot refute
 578 computationalism. *Artificial Intelligence*, 104(1-2):265–286, 1998.
- 579
- 580 Steven M LaValle. *Planning algorithms*. Cambridge university press, 2006.
- 581
- 582 Shane Legg and Marcus Hutter. Universal intelligence: A definition of machine intelligence. *Minds*
 583 *and machines*, 17(4):391–444, 2007.
- 584
- 585 Shane Legg, Marcus Hutter, et al. A collection of definitions of intelligence. *Frontiers in Artificial*
 586 *Intelligence and applications*, 157:17, 2007.
- 587
- 588 John R Lucas. Minds, machines and Gödel. *Philosophy*, 36(137):112–127, 1961.
- 589
- 590 Nestor Maslej, Loredana Fattorini, Raymond Perrault, Yolanda Gil, Vanessa Parli, Njenga Kariuki,
 591 Emily Capstick, Anka Reuel, Erik Brynjolfsson, John Etchemendy, et al. Artificial intelligence
 592 index report 2025. *arXiv preprint arXiv:2504.07139*, 2025.
- 593
- 594 John McCarthy, Marvin Minsky, Nathaniel Rochester, and Claude Shannon. A proposal for the
 595 Dartmouth summer research project on artificial intelligence. [http://jmc.stanford.edu/](http://jmc.stanford.edu/articles/dartmouth/dartmouth.pdf)
 596 [articles/dartmouth/dartmouth.pdf](http://jmc.stanford.edu/articles/dartmouth/dartmouth.pdf), 1955.
- 597
- 598 Marvin Minsky. Steps toward artificial intelligence. *Proceedings of the IRE*, 49(1):8–30, 1961.

- 594 Meredith Ringel Morris, Jascha Sohl-Dickstein, Noah Fiedel, Tris Warkentin, Allan Dafoe, Aleksan-
 595 dra Faust, Clement Farabet, and Shane Legg. Position: Levels of AGI for operationalizing progress
 596 on the path to AGI. In *Forty-first International Conference on Machine Learning*, 2024.
- 597
- 598 Roger Penrose and Martin Gardner. The emperor's new mind: Concerning computers, minds, and
 599 the laws of physics. 1989.
- 600 Henry Gordon Rice. Classes of recursively enumerable sets and their decision problems. *Transactions*
 601 *of the American Mathematical society*, 74(2):358–366, 1953.
- 602
- 603 David Rolnick, Priya L Donti, Lynn H Kaack, Kelly Kochanski, Alexandre Lacoste, Kris Sankaran,
 604 Andrew Slavin Ross, Nikola Milojevic-Dupont, Natasha Jaques, Anna Waldman-Brown, et al.
 605 Tackling climate change with machine learning. *ACM Computing Surveys (CSUR)*, 55(2):1–96,
 606 2022.
- 607
- 608 Stuart Russell. *Human Compatible: Artificial Intelligence and the Problem of Control*. Viking, 2019.
- 609
- 610 Stuart Russell and Peter Norvig. *Artificial Intelligence: A Modern Approach*. Pearson, 3rd edition,
 611 2016.
- 612
- 613 Herbert A Simon. *Models of Man: Social and Rational*. Wiley, 1957.
- 614
- 615 Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Mohamed Amin, Le Hou,
 616 Kevin Clark, Stephen R Pfahl, Heather Cole-Lewis, et al. Toward expert-level medical question
 617 answering with large language models. *Nature Medicine*, 31(3):943–950, 2025.
- 618
- 619 Max Tegmark and Steve Omohundro. Provably safe systems: the only path to controllable AGI.
 620 *arXiv preprint arXiv:2309.01933*, 2023.
- 621
- 622 Alan M. Turing. On computable numbers, with an application to the entscheidungsproblem. *Pro-
 623 ceedings of the London Mathematical Society*, s2-42(1):230–265, 1937.
- 624
- 625 Alan M. Turing. Computing machinery and intelligence. *Mind*, LIX(236):433–460, 1950.
- 626
- 627 Alan M. Turing. Intelligent machinery, a heretical theory. <https://uberty.org/wp-content/uploads/2015/02/intelligent-machinery-a-heretical-theory.pdf>, 1951. Accessed: 2025-07-21.
- 628
- 629 Amos Tversky and Daniel Kahneman. Judgment under uncertainty: Heuristics and biases. *Science*,
 630 185(4157):1124–1131, 1974.
- 631
- 632 Ben Van Calster, David J McLernon, Maarten van Smeden, Laure Wynants, Ewout W Steyerberg,
 633 et al. Calibration: the achilles heel of predictive analytics. *BMC Medicine*, 17:230, 2019.
- 634
- 635 Hanchen Wang, Tianfan Fu, Yuanqi Du, Wenhao Gao, Kexin Huang, Ziming Liu, Payal Chandak,
 636 Shengchao Liu, Peter Van Katwyk, Andreea Deac, et al. Scientific discovery in the age of artificial
 637 intelligence. *Nature*, 620(7972):47–60, 2023.
- 638
- 639 Shen Wang, Tianlong Xu, Hang Li, Chaoli Zhang, Joleen Liang, Jiliang Tang, Philip S Yu, and
 640 Qingsong Wen. Large language models for education: A survey and outlook. *arXiv preprint
 641 arXiv:2403.18105*, 2024.
- 642
- 643 Norbert Wiener. *The Human Use of Human Beings: Cybernetics and Society*. Houghton Mifflin,
 644 Boston, 1950.
- 645
- 646 Wikipedia. https://en.wikipedia.org/wiki/Penrose%E2%80%93Lucas_argument. Accessed: 2025-07-21.
- 647
- Ziwei Xu, Sanjay Jain, and Mohan Kankanhalli. Hallucination is inevitable: An innate limitation of
 648 large language models. *arXiv preprint arXiv:2401.11817*, 2024.
- 649
- 650 Tom Zahavy, Vivek Veeriah, Shaobo Hou, Kevin Waugh, Matthew Lai, Edouard Leurent, Nenad
 651 Tomasev, Lisa Schut, Demis Hassabis, and Satinder Singh. Diversifying AI: Towards creative
 652 chess with AlphaZero. *arXiv preprint arXiv:2308.09175*, 2023.

648 A ADDITIONAL RESULTS
649650 This section proves some additional results discussed in the main text.
651652 A.1 IMPOSSIBILITY OF SOLVING PLANNING WHEN A FEASIBLE PLAN EXISTS
653654 In this section, we prove a similar result to Theorem 3.8, for the case where the program terminates
655 on the given input. We also strengthen the result in Theorem 3.8 to allow for randomized AI systems,
656 and randomized programs which may halt with some probability on an input. We first extend
657 Definition 3.7 to allow for randomized programs.658 **Definition A.1** (Halting for a specific program input instance for randomized programs). *Given a
659 program, input pair, check whether on the given input the given (possibly randomized) program
660 ‘always halts’, ‘halts on some randomness but not all randomness’ or ‘never halts’.*661 **Theorem A.2.** *If the AI system A is accurate and well-behaved for determining halting on a specific
662 program input pair, then there is a program input instance pair for which there is a short proof that
663 the instance always terminates, but A cannot determine that the instance always terminates (the
664 probability of A giving the answer ‘always halts’ is 0).*665 Note that the condition ‘If the AI system A is accurate’ is satisfied when A is trusted, and hence for a
666 accurate, trusted, well-behaved AI system there are program, input instances which the system cannot
667 solve, but for which there is a short proof that the instance terminates. We now prove Theorem A.2.
668669
670 *Proof of Theorem A.2.* Consider Algorithm 3.
671672 **Algorithm 3** Turing_program_v2
673

```

1: procedure Turing_program_v2(Program  $P$ , Input  $I$ )
2:   if  $A(P, I) == \text{`always halts'}$  then
3:     while true do                                 $\triangleright$  run indefinitely
4:     end while
5:   else
6:     return 0
7:   end if
8: end procedure

```

682 Now define
683

$$\text{self_Turing_program_v2}(P) = \text{Turing_program_v2}(P, P).$$

684 We consider:
685

$$\text{self_Turing_program_v2}(\text{self_Turing_program_v2})$$

686 **Lemma A.3.** *If A is accurate and well-behaved then $\text{self_Turing_program_v2}(\text{self_Turing_program_v2})$
687 always halts, but the accurate AI system A cannot determine that it always halts.*688
689 *Proof.* Note that the **if** condition in step 2 always terminates if A is well-behaved. Suppose A outputs
690 ‘always halts’ on some randomness. Then, $\text{self_Turing_program_v2}(\text{self_Turing_program_v2})$
691 does not halt on some randomness. If A is accurate, then this is a contradiction. Therefore, if A is
692 accurate then it must output ‘always halts’ with 0 probability.693 Note that if A does not output that $\text{self_Turing_program_v2}(\text{self_Turing_program_v2})$
694 ‘always halts’, then the program enters the **else** condition and immediately
695 terminates, and therefore halts. Therefore if A is accurate and well-behaved, then
696 $\text{self_Turing_program_v2}(\text{self_Turing_program_v2})$ always halts.697
698
699
700
701

□

□

702 A.2 IMPOSSIBILITY OF SOLVING FEASIBILITY WHEN A PATH EXISTS IN THE GRAPH
703704 We prove a similar result to Theorem 3.12 in this section, for the case where the program terminates
705 on the given input within some time bound. As in Appendix A.1, we also strengthen the result to
706 allow randomized AI systems. We first extend Definition 3.11 to allow for randomized programs.707 **Definition A.4** (Time-bounded halting for randomized programs). *Given a program, input pair*
708 *and a time limit T on the number of execution steps, check whether on the given input the given*
709 *(possibly randomized) program ‘always halts in given time limit’, ‘halts in given time limit T on some*
710 *randomness but not all randomness’ or ‘never halts in given time limit’.*711 **Theorem A.5.** *If an AI system A is accurate and always halts in some time T , then for a fixed constant*
712 *c and the time limit $T + c$, there is a program, input pair for which there is a short, constant-sized*
713 *proof that the instance always halts in at most $T + c$ steps, but for a accurate AI system A which*
714 *always halts in time T the probability of A giving the answer ‘always halts in given time limit’ is 0.*717 *Proof of Theorem A.5.* Consider Algorithm 4, where T is the upper bound on the running time of
718 the AI system A , and c is some fixed constant which is the running time of executing step 2 after A
719 terminates and the **if** condition in step 2 is not satisfied, and then executing steps 5 and 6. Therefore,
720 $T + c$ is an upper bound of the running time of the program when it enters the **else** condition in line 5.721 **Algorithm 4** Turing.T_v2
722723 1: **procedure** Turing.T_v2(Program P , Input I)
724 2: **if** $A(P, I, T + c) ==$ ‘always halts in given time limit’ **then**
725 3: **while** true **do** ▷ run indefinitely
726 4: **end while**
727 5: **else**
728 6: **return** 0
729 7: **end if**
730 8: **end procedure**731 We define
732

733
$$\text{self_Turing_T_v2}(P) = \text{Turing_T_v2}(P, P)$$

734 Consider:
735

736
$$\text{self_Turing_T_v2}(\text{self_Turing_T_v2})$$

738 **Lemma A.6.** *If A is accurate and always terminates in time T , then*
739 *$\text{self_Turing_T_v2}(\text{self_Turing_T_v2})$ always halts in time at most $T + c$. Moreover, if*
740 *A is accurate then it has 0 probability of giving the answer ‘always halts in given time limit’ on*
741 *whether $\text{self_Turing_T_v2}(\text{self_Turing_T_v2})$ halts in time at most $T + c$.*743 *Proof.* Note that by the definition of c , the execution of steps 2, 5 and 6 of the program only take c
744 steps outside the execution of A on the input $(\text{self_Turing_T_v2}, \text{self_Turing_T_v2}, T + c)$.
745746 Suppose A outputs ‘always halts in given time limit’ on the given input on some randomness.
747 Whenever A outputs ‘always halts in given time limit’, the program enters an infinite loop and never
748 halts. This is a contradiction if A is accurate, and hence if A is accurate it outputs ‘always halts in
749 given time limit’ with probability 0.750 Now, if A does not output ‘always halts in given time limit’ on the input, then the program will enter
751 the **else** block and immediately halt. Since A runs for at most T steps, the program then halts in time
752 at most $T + c$. Therefore, if A is accurate then the program always halts in time at most $T + c$.
753 □754
755 □

756 A.3 PROOF OF THEOREM 4.2
757758 **Theorem 4.2.** *If the AI system A is well-behaved and calibrated for program verification, then there
759 is a program P which provably halts with probability at least 0.99, but A abstains with probability 1
760 on the program P .*761 *Proof of Theorem 4.2.* Consider Algorithm 5. Throughout the proof we assume A is well-behaved,
762 i.e. it always terminates. Our construction involves a program which does not take any input, i.e.
763 $I = \emptyset$. The program involves identifying whether the probability p of A outputting ‘halts’ when given
764 Gödel_program_random as input is greater than 0.5 or not. We use a simple best arm identification
765 procedure for this, for example the algorithm of Karnin et al. (2013).
766767 **Algorithm 5** Gödel-program_random
768

```

769 1: procedure Gödel_program_random
770 2:   Let arm1 have the distribution Bernoulli(0.5)
771 3:   Let arm2 correspond to running  $A$  with Gödel_program_random as input, with the result of
772   the arm pull being 1 if  $A(\text{Gödel\_program\_random}) == \text{'halts'}$ , and 0 otherwise.
773 4:   Run Best-Arm-Identification algorithm from Karnin et al. (2013, Algorithm 1) with confi-
774   dence parameter  $\delta = 0.01$  to determine whether arm2 is better than arm1
775 5:   if arm2 is better than arm1 then
776 6:     while true do                                 $\triangleright$  run indefinitely
777 7:     end while
778 8:   else
779 9:     return 0
780 10:  end if
781 11: end procedure

```

782 Note that for any $\epsilon > 0$, if $p = 0.5 + \epsilon$ then arm₂ is better than arm₁, otherwise if $p = 0.5 - \epsilon$
783 then arm₁ is better than arm₂. While we can use any suitable multi-armed bandit algorithm in our
784 construction, here we use Karnin et al. (2013), which has the guarantee that if it is provided with two
785 arms with a gap of ϵ , then it finds the better arm with probability $1 - \delta$ using $O\left(\frac{1}{\epsilon^2} \log\left(\frac{1}{\delta} \log\left(\frac{1}{\epsilon}\right)\right)\right)$
786 arm pulls. This bound is known to be optimal (Jamieson et al., 2014), though in our case since we do
787 not care about the optimal rate we could have also used earlier sub-optimal procedures (Even-Dar
788 et al., 2002). We also note that if $\epsilon = 0$, then the best arm identification procedure will terminate with
789 probability at most 10δ . Though we have not seen the case of $\epsilon = 0$ being directly covered by the
790 guarantees of best arm identification procedures, this claim for $\epsilon = 0$ follows from a simple argument
791 which treats the best arm identification procedure as a black-box. To verify, note that the sequence of
792 observations up to t steps is δ -close in TV distance for any $p \in [0.5 \pm 1/\text{poly}(t, \delta)]$ (where $\text{poly}(t, \delta)$
793 is some polynomial of t and δ). Therefore for $\epsilon = 0$ and any finite t , if the best arm identification
794 procedure terminates in t steps with probability more than 10δ , then it will have a failure probability
795 more than δ for some $p \in [0.5 \pm 1/\text{poly}(t, \delta)]$ —which is a contradiction with the guarantee of the
796 procedure. Therefore, for $\epsilon = 0$ the best arm identification procedure terminates with probability at
797 most 10δ .

798 We are now ready to prove the result.

799 **Lemma A.7.** *If A is calibrated and well-behaved then Gödel_program_random halts with proba-
800 bility at least 0.99, but the calibrated AI system A will output ‘don’t know’ with probability 1 on
801 Gödel_program_random.*802 *Proof.* We consider three cases.
803

-
1. $p \in (0.5, 1]$: Note that in this case with probability at least 0.99 the best arm identification
804 procedure determines that arm₂ is better than arm₁. Therefore, the program goes into the
805 infinite **while** loop and never terminates with probability at least 0.99. In this case, A is not
806 calibration safe, since it claims that the program terminates with probability $p > 0.5$.
 2. $p = 0.5$: As argued above, in this case the best arm identification procedure terminates
807 with probability at most $10\delta = 0.1$. Therefore, Gödel_program_random terminates with
808 probability at most 0.1.

810 probability at most 0.1. In this case as well, A is not calibration safe, since it claims that the
 811 program terminates with probability $p = 0.5$.
 812

- 813 3. $p \in (0, 0.5)$: In this case, with probability at least 0.99 the best arm identification procedure
 814 determines that arm_1 is better than arm_2 . When arm_1 is determined to be better than arm_2 ,
 815 the program enters the `else` block in line 9. Therefore, in this case `Gödel_program_random`
 816 terminates with probability at least 0.99. Here too, A is not calibrated, since it claims that
 817 `Gödel_program_random` terminates with probability $p < 0.5$.
 818

819 In each of these cases, A is not calibrated. Therefore, for A to be calibrated, we must have $p = 0$, and
 820 that A outputs ‘don’t know’ with probability 1. If $p = 0$, then with probability at least 0.99 the best
 821 arm identification procedure determines that arm_1 is better than arm_2 , and `Gödel_program_random`
 822 terminates. Therefore, if A is calibrated, then `Gödel_program_random` halts with probability at least
 823 0.99.
 824 \square
 825 \square
 826 \square
 827

828 B IMPOSSIBILITY RESULT FOR OTHER SEMANTIC PROPERTIES

830 In this section, we extend the program verification result from halting to other properties of programs.
 831 This includes properties such as verifying if a program executes certain pre-defined “harmful”
 832 behavior, which could be important from the perspective of safety. More formally, we first define the
 833 task of determining if an input program executes certain harmful behavior.
 834

835 **Definition B.1** (Harmful state execution). *A program is “harmless” if it does not enter certain
 836 pre-defined “harmful” states during execution, and “harmful” otherwise.*

837 We show an impossibility result for this property which is analogous to our result for verifying if an
 838 input program halts.
 839

840 **Theorem B.2.** *Consider an AI system A which is well-behaved and itself does not enter pre-defined
 841 “harmful” states during its execution. Then if A is accurate and trusted for verifying harmful state
 842 execution (Definition B.1), then it cannot be a human-level reasoning system.*

843 This result follows as a corollary of a more general theorem, which extends our result for halting to
 844 any *non-trivial, semantic* property of a program. These are the same conditions under which Rice’s
 845 theorem extends Turing’s undecidability result (Rice, 1953).

846 **Definition B.3** (Non-trivial, semantic property (Rice, 1953)). *A semantic property is a property of
 847 program which concerns its behavior (e.g. “does the program always terminate?”) as opposed to its
 848 syntax (e.g. “does the program have a while loop?”). A non-trivial property is a property which is
 849 neither true for all programs, nor false for all programs.*

850 Examples of semantic properties that we have already discussed are halting and harmful state
 851 execution. Another semantic property is correctness, for example if a program intended to determine
 852 if an input number is prime correctly outputs whether the number is prime.
 853

854 We now state the result. The result applies to AI systems A which are well-behaved, and whose
 855 execution does not trivially lead to the desired property π being satisfied. Since our constructions are
 856 self-referential, this condition ensures that the program does not automatically have the property π
 857 by virtue of the execution of A . For instance, if A itself always entered “harmful” states during its
 858 execution, then a program which always calls A also trivially enters “harmful” states. Note that from
 859 the perspective of safety, if A itself entered “harmful” states then it would also have unsafe behavior.

860 **Theorem B.4.** *Consider any non-trivial, semantic property π of programs. Consider an AI system A
 861 which is well-behaved and has the property that if some program P calls A during its execution then
 862 the execution of A never automatically leads to the property π being satisfied for program P . Then if
 863 A is accurate and trusted for the task of verifying if an input program has property π , then it cannot
 864 be a human-level reasoning system for this task.*

864 *Proof.* For any property non-trivial, semantic property π , let `valid_program_for_` π be some
 865 program which has property π , and `invalid_program_for_` π be some program which does not
 866 have property π . Note that since π is a non-trivial property, both these programs exist. Let A be an
 867 AI system which takes some program as input, and verifies if the input program has the property
 868 π . Note that A is well-behaved, its execution does not automatically lead to some program having
 869 property π . Now consider Algorithm 6, for any input I .
 870

871 **Algorithm 6** `Rice_program` for property π

872 1: **procedure** `Rice_program`(Input I)
 873 2: **if** $A(\text{Rice_program}) == \text{'has property } \pi\text{'}$ **then**
 874 3: **return** `invalid_program_for_` $\pi(I)$
 875 4: **else**
 876 5: **return** `valid_program_for_` $\pi(I)$
 877 6: **end if**
 7: **end procedure**

880 We claim that if A is accurate, then it cannot output that `Rice_program` has property π . This is
 881 true by contradiction, if A outputs that `Rice_program` has property π , then `Rice_program` calls
 882 `invalid_program_for_` π , and hence does not have property π . Note that if A does not output
 883 that `Rice_program` has property π , then `Rice_program` always calls `valid_program_for_` π ,
 884 and hence `Rice_program` has property π .

885 Therefore if A is accurate, then `Rice_program` has property π , but A cannot output that
 886 `Rice_program` has property π . Hence an accurate and trusted system for verifying property
 887 π cannot be a human-level reasoning system for the task.

888 □
 889
 890
 891
 892
 893
 894
 895
 896
 897
 898
 899
 900
 901
 902
 903
 904
 905
 906
 907
 908
 909
 910
 911
 912
 913
 914
 915
 916
 917