

# INCOMPLETENESS FOR STABLY SOUND TURING MACHINES

YASHA SAVELYEV

**ABSTRACT.** We first partly develop a mathematical notion of stable soundness intended to reflect the actual soundness property of human beings. Then we show that given an abstract query machine  $M$  the following cannot hold simultaneously:  $M$  is stably sound,  $M$  is computable,  $M$  can stably decide the truth of any arithmetic statement. This can be understood as an extension of the Gödel incompleteness theorem to stably sound setting. This is a non-trivial extension as stably sound Turing machine can decide the halting problem. In practice such an  $M$  could be meant to represent a weakly idealized human being so that the above gives an obstruction to computability of intelligence, and this gives a formal extension of a famous disjunction of Gödel.

## 1. INTRODUCTION

The term **machine**<sup>1</sup> will just be a synonym for a partial map  $M : A \rightarrow B$ , with  $A, B$  sets with an additional structure of an encoding in  $\mathbb{N}$ . An encoding is just an injective map  $e : A \rightarrow \mathbb{N}$  with some extra properties. This is described in more detail in Section 2.1.

Let  $\mathcal{A}$  denote the set of (first order) sentences of arithmetic. And suppose we are given a machine

$$M : \mathbb{N} \rightarrow \mathcal{A} \times \{\pm\},$$

for  $\{\pm\}$  a set with two elements  $+, -$ . If there is an  $n_0$  with  $M(n_0) = (\Sigma, +)$  s.t. there is no  $m > n_0$  with  $M(m) = (\Sigma, -)$  then we say that  $\Sigma$  **is**  $M$ -**stable**.

**Definition 1.1.** *We say that a machine*

$$M : \mathbb{N} \rightarrow \mathcal{A} \times \{\pm\}$$

*is stably sound if for all  $\Sigma \in \mathcal{A}$ :*

$$\Sigma \text{ is } M\text{-stable} \implies \Sigma \text{ is true.}$$

*We say that  $M$  decides arithmetic if*

$$\forall \Sigma \in \mathcal{A} : \Sigma \text{ is true} \implies \Sigma \text{ is } M\text{-stable.}$$

So such an  $M$  has the following interpretation:

$$M(n) = (\Sigma, +)$$

if at the moment  $n$   $M$  decides that  $\Sigma$  is true, while

$$M(m) = (\Sigma, -)$$

if at the moment  $m$   $M$  does not decide that  $\Sigma$  is true. If  $\Sigma$  is  $M$ -stable we can also say  $M$  **stably decides**  $\Sigma$  to be true. So that  $M$  is stably sound if and only if whenever  $M$  stably decides  $\Sigma$  to be true, it is in fact true.

The following is one version of our main theorem, with Theorem 4.5 being a more fundamental variant. We present this formulation first because a human mathematician intuitively works as a machine  $M$  above.

---

<sup>1</sup>For some authors and in some of the writing of Turing and Gödel “machine” is synonymous with Turing machine. For us the term machine is just abstraction for a process acting on inputs, but it need not be a computational process, in contrast to Turing machines.

**Theorem 1.2.** *For  $M$  as above the following cannot hold simultaneously:  $M$  is stably sound,  $M$  is computable and  $M$  decides arithmetic. Here ‘computable’ has the mostly standard meaning of computability by a Turing machine, with specifics given in Section 2.1.*

If we replace stably sound by sound then the above would just be a reformulation of Gödel’s incompleteness theorem. Moreover, if we specialize our argument to this case then in the statement of the theorem we may also specialize to statements (in Turing machine language) of the kind: the Turing machine  $T$  does not halt on input  $n$ . In other words, we would get the well known result of Turing that the halting problem is not soundly decidable by a Turing machine. On the other hand there is a stably sound Turing machine which can stably soundly decide the halting problem, cf. Example 3.3. So perhaps Turing machines stably soundly decide the truth of any arithmetic statement? The above theorem shows that this is not the case.

**1.1. Motivational background - intelligence, Gödel’s disjunction and Penrose.** In what follows we understand *human intelligence* very much like Turing in [2], as a black box which receives inputs and produces outputs. More specifically, this black box  $M$  is meant to be some system which contains a human subject. We do not care about what is happening inside  $M$ . So we are not directly concerned here with such intangible things as understanding, intuition, consciousness - the inner workings of human intelligence that are supposed as special. The only thing that concerns us is what output  $M$  produces given an input, not how it is produced. Given this *very* limited interpretation, the question that we are interested in is this:

*Question 1.* Can human intelligence be completely modelled by a Turing machine?

An informal definition of a Turing machine (see [1]) is as follows: it is an abstract machine which permits certain inputs, and produces outputs. The outputs are determined from the inputs by a fixed finite algorithm, defined in a certain precise sense. For a non-expert reader we point out that this “fixed” does not preclude the algorithm from “learning”,<sup>2</sup> it just means that how it “learns” is completely determined by the initial algorithm. In particular anything that can be computed by computers as we know them can be computed by a Turing machine. For our purposes the reader may simply understand a Turing machine as a digital computer with unbounded memory running some particular program. Unbounded memory is just a mathematical convenience. In specific arguments, also of the kind we make, we can work with non-explicitly bounded memory. Turing himself has started on a form of Question 1 in his “Computing machines and Intelligence” [2], where he also informally outlined a possible obstruction to a yes answer coming from Gödel’s incompleteness theorem.

For the incompleteness theorem to have any relevance we need some assumption on the soundness or consistency of human reasoning. Informally, a human is sound if whenever they asserts something in absolute faith this something is indeed true. This requires context as truth in general is undefinable. For our arguments later on the context will be in certain mathematical models. However, we cannot honestly hope for soundness, as even mathematicians are not on the surface sound at all times, they may assert mathematical untruths at various times, (but usually not in absolute faith). But we can certainly hope for some kind of fundamental soundness.

In this work we will formally interpret fundamental soundness as stable soundness, which we already partly described above. This notion of stable soundness is meant to reflect the basic understanding of the way science progresses. Of course even stable soundness needs idealizations to make sense for individual humans. The human brain deteriorates and eventually fails, so that either we idealize the human brain to never deteriorate in particular never die, or  $M$  now refers not to an individual human but to the evolving scientific community. We call such a human *weakly idealized*.

Around the same time as Turing, Gödel argued for a no answer to Question 1, see [11, 310], relating the question to existence of absolutely unsolvable Diophantine problems, see also Feferman [6], and Koellner [13], [14] for a discussion. Essentially, Gödel argues for a disjunction:

$$\neg((S \text{ is computable}) \wedge (S \text{ is sound}) \wedge A),$$

---

<sup>2</sup>In the sense of “machine learning”.

where  $S$  refers to a certain idealized subject, and where  $A$  says that  $S$  can decide any Diophantine problem. Gödel's argument can be formalized, see [14]. At the same time Gödel doubted that  $\neg A$  is possible, again for an idealized  $S$ , as this puts absolute limits on what is humanly knowable even in arithmetic. Note that his own incompleteness theorem only puts relative limits on what is humanly knowable, within a fixed formal system.

However, what is the meaning of ‘idealized’ above? If idealized just means stabilized in the sense of this paper (Section 3 specifically) then there is a Turing machine  $T$  whose stabilization  $T^s$  provably decides the halting problem, cf. Example 3.3, and so  $T^s$  is no longer computable. In that case, the above disjunction becomes meaningless because passing to the idealization may introduce non-computability where there was none before. So in this context one must be extremely detailed with what “idealized” means physically. The process of the physical idealization must be such that non-computability is not introduced in the ideal limit. In the case of weak idealization mentioned above it should certainly be in principle possible but this not what is needed by Gödel. He needs an idealization that is plausibly sound otherwise the disjunction would again be meaningless, while a weak idealization of a human is only plausibly stably sound. Since we have no idea what is happening in the human brain, it is not at all clear that what Gödel asks is even possible.

So the natural solution to attempt is to enrich the argument of Gödel so that it explicitly allows for just stable soundness. But then we may worry: if stable soundness is such a loose concept that a Turing machine machine can stably soundly decide the halting problem, maybe Turing machines can stably soundly decide anything? So we need a new incompleteness theorem and this is the theorem above.

After Gödel, Lucas [10] and later again and more robustly Penrose [16] argued for a no answer based only on soundness and the Gödel incompleteness theorem, that is attempting to remove the necessity to decide  $A$  or  $\neg A$ . A number of authors, particularly Koellner [13], [14], argue that there are likely unresolvable meta-logical issues with the Penrose argument, even allowing for soundness. See also Penrose [16], and Chalmers [4] for discussions of some issues. The issue, as I see it, is loosely speaking the following. The kind of argument that Penrose proposes is a meta-algorithm  $P$  that takes as input specification of a Turing machine or a formal system, and has as output a natural number (or a string, sentence). Moreover, each step of this meta-algorithm is computably constructive. But the goal of the meta-algorithm  $P$  is to prove  $P$  is not computable as a function! So on this rough surface level this appears to be impossible. But it may be possible to modify the idea in some subtle way.

Notwithstanding, what we argue here is that there is much more compelling version of the original Gödel disjunction that only needs stable soundness. The following is a slightly informal, applied version of our Theorem 1.2.

**Theorem 1.3.** *Either there are cognitively meaningful, absolutely non Turing computable processes in the human brain, or human beings are stably unsound, or for any (could be weakly idealized)  $S$  there exists a certain true arithmetic statement, let's call it  $\mathcal{H}(S)$ <sup>3</sup>, that  $S$  will never stably decide to be true. The above is indeed a mathematical fact<sup>4</sup>, given our formalizations.*

By *absolutely* we mean in any sufficiently accurate physical model. Note that even existence of absolutely non Turing computable processes in nature is not known. For example, we expect beyond reasonable doubt that solutions of fluid flow or  $N$ -body problems are generally non Turing computable (over  $\mathbb{Z}$ , if not over  $\mathbb{R}$  cf. [3])<sup>5</sup> as modeled in essentially classical mechanics. But in a more physically accurate and fundamental model they may both become computable, possibly if the nature of the universe is ultimately discreet. It would be good to compare this theorem with Deutch [5], where computability of any suitably finite and discreet physical system is conjectured. Although this is not immediately at odds with us, as the hypothesis of that conjecture may certainly not be satisfiable.

<sup>3</sup> $\mathcal{H}(S)$  is a statement in the language of Turing machines and so is number theoretic, however it is not a Diophantine problem. Of course it cannot be by Example 3.3.

<sup>4</sup>Specifically a theorem of set theory, although we keep set theory implicit as usual.

<sup>5</sup>We are now involving real numbers but there is a standard way of talking of computability in this case, in terms of computable real numbers. This is what means over  $\mathbb{Z}$ .

**Remark 1.4.** *It should also be noted that for Penrose, in particular, non-computability of intelligence would be evidence for new physics, and he has specific and very intriguing proposals with Hameroff [9] on how this can take place in the human brain. As we have already indicated, new physics is not a logical necessity for non-computability of brain processes, at least given the state of the art. However, it is very plausible that new physical-mathematical ideas may be necessary to resolve the deep mystery of human consciousness. Here is also a partial list of some partially related work on mathematical models of brain activity, consciousness and or quantum collapse models: [12], [15], [7], [8].*

As a final remark, technically the paper is mostly elementary and should be widely readable in entirety.

## 2. SOME PRELIMINARIES

This section can be just skimmed on a first reading. Really what we are interested in is not Turing machines per se, but computations that can be simulated by Turing machine computations. These can for example be computations that a mathematician performs with paper and pencil, and indeed is the original motivation for Turing's specific model. However to introduce Turing computations we need Turing machines. Here is our version, which is a computationally equivalent, minor variation of Turing's original machine.

**Definition 2.1.** *A Turing machine  $M$  consists of:*

- *Three infinite (1-dimensional) tapes  $T_i, T_o, T_c$ , (input, output and computation) divided into discreet cells, next to each other. Each cell contains a symbol from some finite alphabet  $\Gamma$  with at least two elements. A special symbol  $b \in \Gamma$  for blank, (the only symbol which may appear infinitely many often).*
- *Three heads  $H_i, H_o, H_c$  (pointing devices),  $H_i$  can read each cell in  $T_i$  to which it points,  $H_o, H_c$  can read/write each cell in  $T_o, T_c$  to which they point. The heads can then move left or right on the tape.*
- *A set of internal states  $Q$ , among these is "start" state  $q_0$ . And a non-empty set  $F \subset Q$  of final states.*
- *Input string  $\Sigma$ : the collection of symbols on the tape  $T_i$ , so that to the left and right of  $\Sigma$  there are only symbols  $b$ . We assume that in state  $q_0$   $H_i$  points to the beginning of the input string, and that the  $T_c, T_o$  have only  $b$  symbols.*
- *A finite set of instructions:  $I$ , that given the state  $q$  the machine is in currently, and given the symbols the heads are pointing to, tells  $M$  to do the following. The actions taken, 1-3 below, will be (jointly) called an **executed instruction set** or just **step**:*
  - (1) *Replace symbols with another symbol in the cells to which the heads  $H_c, H_o$  point (or leave them).*
  - (2) *Move each head  $H_i, H_c, H_o$  left, right, or leave it in place, (independently).*
  - (3) *Change state  $q$  to another state or keep it.*
- *Output string  $\Sigma_{out}$ , the collection of symbols on the tape  $T_o$ , so that to the left and right of  $\Sigma_{out}$  there are only symbols  $b$ , when the machine state is final. When the internal state is one of the final states we ask that the instructions are to do nothing, so that these are frozen states.*

**Definition 2.2.** *A complete configuration of a Turing machine  $M$  or total state is the collection of all current symbols on the tapes, position of the heads, and current internal state. Given a total state  $s$ ,  $\delta(s)$  will denote the successor state of  $s$ , obtained by executing the instructions set of  $M$  on  $s$ , or in other words  $\delta(s)$  is one step forward from  $s$ .*

So a Turing machine determines a special kind of function:

$$\delta^M : \mathcal{C}(M) \rightarrow \mathcal{C}(M),$$

where  $\mathcal{C}(M)$  is the set of possible total states of  $M$ .

**Definition 2.3.** A **Turing computation**, or **computation sequence** for  $M$  is a possibly not eventually constant sequence

$$*M(\Sigma) := \{s_i\}_{i=0}^{i=\infty}$$

of total states of  $M$ , determined by the input  $\Sigma$  and  $M$ , with  $s_0$  the initial configuration whose internal state is  $q_0$ , and where  $s_{i+1} = \delta(s_i)$ . If elements of  $\{s_i\}_{i=0}^{i=\infty}$  are eventually in some final machine state, so that the sequence is eventually constant, then we say that the computation **halts**. For a given Turing computation  $*M(\Sigma)$ , we will write

$$*M(\Sigma) \rightarrow x,$$

if  $*M(\Sigma)$  halts and  $x$  is the corresponding output string.

We write  $M(\Sigma)$  for the output string of  $M$ , given the input string  $\Sigma$ , if the associated Turing computation  $*M(\Sigma)$  halts. Denote by *Strings* the set of all finite strings of symbols in  $\Gamma$ , including the empty string  $\epsilon$ . Then a Turing machine  $M$  determines a partial function that is defined on all  $\Sigma \in \text{Strings}$  s.t.  $*M(\Sigma)$  halts, by  $\Sigma \mapsto M(\Sigma)$ .

In practice we will allow our Turing machine  $T$  to reject some elements of *Strings* as valid input. We may formalize this by asking that there is a special final machine state  $q_{\text{reject}}$  so that  $T(\Sigma)$  halts with  $q_{\text{reject}}$  for

$$\Sigma \notin \mathcal{I} \subset \text{Strings},$$

where  $\mathcal{I}$  is some set of all valid, that is  $T$ -**permissible** input strings. We do not ask that for  $\Sigma \in \mathcal{I}$   $*T(\Sigma)$  halts. If  $*T(\Sigma)$  does halt then we will say that  $\Sigma$  is  $T$ -**acceptable**.

**Definition 2.4.** We denote by  $\mathcal{T}$  the set of all Turing machines with a distinguished final machine state  $q_{\text{reject}}$ .

It will be convenient to forget  $q_{\text{reject}}$  and instead write

$$T : \mathcal{I} \rightarrow \mathcal{O},$$

where  $\mathcal{I} \subset \text{Strings}$  is understood as the subset of all  $T$ -permissible strings, or just **input set** and  $\mathcal{O}$  is the set output strings or **output set**.

**Definition 2.5.** Given a partial function

$$f : \mathcal{I} \rightarrow \mathcal{O},$$

we say that a Turing machine in  $T \in \mathcal{T}$

$$T : \mathcal{I} \rightarrow \mathcal{O}$$

**computes**  $f$  if  $T = f$  as partial functions on  $\mathcal{I}$ .

**2.1. Abstractly encoded sets.** We will sometimes use abstract sets to refer to input and output sets. However, these are understood to be subsets of *Strings* under some implicit, *fixed* encoding. Concretely an **encoding** of  $A$  is an injective set map  $e : A \rightarrow \text{Strings}$ . For example we may encode  $\text{Strings}^2$  as a subset of *Strings* as follows. The encoding string of  $\Sigma = (\Sigma_1, \Sigma_2) \in \text{Strings}^2$  will be of the type: “this string encodes an element  $\text{Strings}^2$ : its components are  $\Sigma_1$  and  $\Sigma_2$ .” In particular the sets of integers  $\mathbb{N}, \mathbb{Z}$ , which we may use, will under some encoding correspond to subsets of *Strings*. Indeed this abstracting of sets from their encoding in *Strings* is partly what computer languages do.

More formally, let  $\mathcal{S}$  be an arrow category whose objects are maps  $e_A : A \rightarrow \text{Strings}$ , for  $e_A$  an embedding called **encoding map of**  $A$ , determined by  $A \in \text{Set}$ , and morphisms  $(m, \tilde{m})$  commutative diagrams:

$$\begin{array}{ccc} A & \xrightarrow{\tilde{m}} & B \\ \downarrow e_A & & \downarrow e_B \\ e_A(A) \subset \text{Strings} & \xrightarrow{m} & e_B(B) \subset \text{Strings}. \end{array}$$

We may just write  $A \in \mathcal{S}$  for an object, with  $e_A$  implicit. We call such an  $A$  an **abstractly encoded set** so that  $\mathcal{S}$  is a category of abstractly encoded sets.

In addition we ask that  $\mathcal{S}$  satisfies the following properties.

- For  $A \in \mathcal{S}$   $e_A(A)$  is computable (recursive). Here, a set  $S \subset \text{Strings}$  is called *computable* if both  $S$  and its complement are computably enumerable, with  $S$  called *computably enumerable* if there is a Turing machine  $T$  so that  $*T(\Sigma)$  halts iff  $\Sigma \in S$ .
- For  $A, B \in \mathcal{S}$   $A \sqcup B$  is in  $\mathcal{S}$ . Moreover, corresponding to the inclusion map  $\tilde{i}_A : A \rightarrow A \sqcup B$  we have a morphism  $(i_A, \tilde{i}_A)$  in  $\mathcal{S}$ :

$$\begin{array}{ccc} A & \xrightarrow{\tilde{i}_A} & A \sqcup B \\ \downarrow & & \downarrow \\ e_A(A) & \xrightarrow{i_A} & e_{A \sqcup B}(A \sqcup B), \end{array}$$

likewise with  $\tilde{i}_B : B \rightarrow A \sqcup B$ . We ask that the map

$$e_A(A) \sqcup e_B(B) \rightarrow e_{A \sqcup B}(A \sqcup B)$$

induced by  $(i_A, \tilde{i}_A), (i_B, \tilde{i}_B)$  is bijective and  $i_A(e_A(A)), i_B(e_B(B))$  are computable.

- If  $A, B \in \mathcal{S}$  then  $A \times B \in \mathcal{S}$  and the projection maps  $\tilde{pr}_A : A \times B \rightarrow A, \tilde{pr}_B : A \times B \rightarrow B$  complete to morphisms of  $\mathcal{S}$ , similar to the above, so that we have a commutative diagram:

$$\begin{array}{ccc} A \times B & \xrightarrow{\tilde{pr}_A} & A \\ \downarrow & & \downarrow \\ e_{A \times B}(A \times B) & \xrightarrow{pr_A} & e_A(A), \end{array}$$

and likewise with  $pr_B$ .

- There is an abstractly encoded set  $\mathcal{U} = \text{Strings} \in \mathcal{S}$ , with  $e_{\mathcal{U}} = id_{\text{Strings}}$ . We can think of  $\mathcal{U}$  as the set of typeless strings. We may of course also write  $e_A(A) \subset \mathcal{U}$ .
- For  $A, B \in \mathcal{S}$ ,

$$(e_A(A) \cap e_B(B) \text{ is non-empty}) \implies (A = \mathcal{U}) \vee (B = \mathcal{U}).$$

In particular each  $A \in \mathcal{S}$  is determined by the image  $e_A(A)$ .

The specific such category  $\mathcal{S}$  that we need will be clear from context later on. We only need to encode finitely many types of specific sets. For example  $\mathcal{S}$  should contain an abstract encoding of  $\mathbb{Z}, \mathbb{N}, \{\pm\}, \mathcal{T}$ , with  $\{\pm\}$  a set with two elements  $+, -$ . The encodings of  $\mathbb{N}, \mathbb{Z}, \mathcal{T}$  should be suitably natural so that for example there is a computable total map  $S : e(\mathbb{N}) \rightarrow e(\mathbb{N})$  which corresponds under  $e$  to the self-map  $n \mapsto n + 1$  of  $\mathbb{N}$ . This is not part of the axioms but we need this to construct specific abstract Turing machines, whose definition follows.

For  $A, B \in \mathcal{S}$ , given a partial set map  $f : A \rightarrow B$  we define

$$f_e := e_B \circ f \circ e_A^{-1},$$

called *the encoding of  $f$* .

**Definition 2.6.** For  $A, B \in \mathcal{S}$  an **abstract Turing machine**

$$T : \text{Dom} \subset A \rightarrow B$$

is a *partial map*

$$f : \text{Dom} \rightarrow B,$$

and  $T_e \in \mathcal{T}$

$$T_e : e_A(\text{Dom}) \rightarrow e_B(B)$$

such that  $T_e$  computes  $f_e$ , with  $e_A(\text{Dom})$  the set of  $T_e$ -permissible strings. The set  $\text{Dom}$  is called the *set of  $T$ -permissible strings*. We often omit  $\text{Dom}$  from notation, especially when  $\text{Dom} = A$ . In practice the partial map  $f$  will just be denoted by  $T$  itself.

We define  $\mathcal{T}_S$  to be the set of abstract Turing machines relative to  $S$  as above. We will not need to encode  $\mathcal{T}_S$ , as there is a natural embedding  $i : \mathcal{T}_S \rightarrow \mathcal{T}$ ,  $i(T) = T_e$  and  $\mathcal{T}$  is encoded, and this will be sufficient for us.

For writing purposes we condense the above as follows.

**Definition 2.7.** A **machine** will be a synonym for a partial map  $M : \text{Dom} \subset A \rightarrow B$ , with  $A, B$  abstractly encoded sets.  $\text{Dom}$  may be omitted from notation. (Note  $M$  is partial on  $\text{Dom}$ , not partial on  $A$  with  $\text{Dom}$  the set where  $M$  is defined.)

$\mathcal{M} = \mathcal{M}_S$  will denote the set of machines. Given an abstract Turing machine  $T : A \rightarrow B$ , we have an associated machine  $\text{fog}(T) : A \rightarrow B$  defined by forgetting the additional structure  $T_e$ . However we may also just write  $T$  for this machine by abuse of notation. So we have a forgetful map

$$\text{fog} : \mathcal{T} \rightarrow \mathcal{M},$$

which forgets the extra structure of a Turing machine.

**Definition 2.8.** We say that an abstract Turing machine  $T$  **computes**  $M \in \mathcal{M}$  if  $\text{fog}(T) = M$ . We say that  $M$  is **computable** if some  $T$  computes  $M$ .

This has a number of expected properties, for example:

**Lemma 2.9.** If  $M : A \rightarrow B$  is computable and  $G : A \rightarrow C$  is computable, then

$$D : A \rightarrow B \times C,$$

$D(a) = (M(a), G(a))$ , is computable.

Let us omit the proof.

**Composition.** Given Turing machines

$$M_1 : A \rightarrow B, M_2 : C \rightarrow D,$$

we may naturally **compose** them to get a Turing machine  $M_2 \circ M_1 : R \rightarrow D$ , for  $R = M_1^{-1}(B \cap C)$ , with the intersection taken in *Strings*.  $R$  can be empty in which case  $M_2 \circ M_1$  is a Turing machine which rejects all input. In the case of abstract Turing machines we ask that  $B = C$  and then the composition operation is likewise readily defined. Let us not elaborate further.

**2.2. Notation.**  $\mathbb{Z}$  always denotes the set of all integers and  $\mathbb{N}$  non-negative integers. We will sometimes specify an (abstract) Turing machine simply by specifying a map

$$T : \mathcal{I} \rightarrow \mathcal{O},$$

with the full data of the underlying Turing machine being implicitly specified, in a way that should be clear from context.

### 3. ON STABLE SOUNDNESS

**Definition 3.1.** Given a machine:

$$M : \mathbb{N} \rightarrow B \times \{\pm\},$$

if there is an  $n_0$  with  $M(n_0) = (\Sigma, +)$  s.t. there is no  $m > n_0$  with  $M(m) = (\Sigma, -)$  then  $\Sigma$  is called  **$M$ -stable** and we say that  $M$  **prints**  $\Sigma$  **stably**.

**Definition 3.2.** Given a machine

$$M : \mathbb{N} \rightarrow B \times \{\pm\},$$

we define

$$M^s : \mathbb{N} \rightarrow \mathcal{U}$$

to be the machine enumerating, in order, all the  $M$ -stable  $\Sigma$ . We call this the **stabilization** of  $M$ . The range of  $M^s$  is called the **stable output** of  $M$ .

In general  $M^s$  may not be computable even if  $M$  is computable. Explicit examples of this sort can be constructed by hand.

*Example 3.3.* We can construct a Turing machine

$$A : \mathbb{N} \rightarrow \mathcal{P} \times \{\pm\},$$

whose stabilization  $A^s$  enumerates every Diophantine (integer coefficients) polynomial with no integer roots, where  $\mathcal{P}$  denotes the set of all Diophantine polynomials, (also abstractly encoded). Similarly, we can construct a Turing machine  $D$  whose stabilization enumerates pairs  $(T, n)$  for  $T : \mathbb{N} \rightarrow \mathbb{N}$  a Turing machine and  $n \in \mathbb{N}$  such that  $*T(n)$  does not halt. In other words  $D$  stably soundly decides the halting problem. To do this we may proceed via a zig-zag algorithm.

In the case of Diophantine polynomials, here is a (inefficient) example. Let  $Z$  computably enumerate every Diophantine polynomial, and let  $N$  computably enumerate the integers. In other words, in our language,  $Z : \mathbb{N} \rightarrow \mathcal{P}$ ,  $N : \mathbb{N} \rightarrow \mathbb{Z}$  are total bijective abstract Turing machines.

- Initialize an ordered list  $L$  by  $L = \emptyset$ , which we understand as a list of instructions.
- Start. For each  $p \in \{Z(0), \dots, Z(n)\}$  check if  $\{N(0), \dots, N(n)\}$  are solutions of  $p$ . Whenever no add  $(p, +)$  to  $L$ , whenever yes add  $(p, -)$ .
- Set  $n := n + 1$  go to Start and continue.

This will define a partial function  $A : \mathbb{N} \rightarrow \mathcal{P} \times \{\pm\}$  whose value  $A(m)$  is the  $m$ 'th, not necessarily final, instruction in the list  $L^m$  which is  $L$  after the  $m$ 'th step of the algorithm with  $L^0 = \{\emptyset\}$ . Clearly  $A$  is computable and its stabilization  $A^s$  enumerates Diophantine polynomials which have no integer roots.

**3.1. Decision machines.** By a *decision machine* we mean a machine of the form:

$$D : \text{Dom} \subset B \times \mathbb{N} \rightarrow \{\pm\}.$$

**Definition 3.4.** For a decision machine  $D$  we say that  $a \in A$  is *D-decided* if there is an  $n_0 \in \mathbb{N}$  with  $D(T, n_0) = +$  s.t. there is no  $m > n_0$  with  $D(T, m) = -$ .

Given a machine  $M : \mathbb{N} \rightarrow B \times \{\pm\}$ , there is an associated decision machine:

$$D_M : \text{Dom} \subset B \times \mathbb{N} \rightarrow \{\pm\},$$

with  $\text{Dom}$  defined to be the set of  $(b, n) \in B \times \mathbb{N}$  such that either  $\text{pr}_B M(n) = b$  or  $M(n)$  is undefined, for  $\text{pr}_B : B \times \{\pm\} \rightarrow B$  the projection.  $D_M$  is defined by

$$D_M(b, n) = \begin{cases} + & \text{if } M(n) = (b, +) \\ - & \text{if } M(n) = (b, -) \\ \text{undefined} & \text{if } M(n) \text{ is undefined.} \end{cases}$$

**Lemma 3.5.** If  $M$  as above is computable then  $D_M$  is computable. Moreover,  $b$  is  $M$ -stable iff  $b$  is  $D_M$ -decided.

*Proof.* To keep notation simple let  $M$  denote the Turing machine computing the machine  $M$ . Let

$$\widetilde{M} : B \times \mathbb{N} \rightarrow (B \times \{\pm\}) \times B$$

be the Turing machine defined (as a partial map) by  $\widetilde{M}(b, n) = (M(n), b)$ , using Lemma 2.9. Define  $T$  to be the composition of (abstract) Turing machines:

$$T := R \circ \widetilde{M},$$

for

$$R : (B \times \{\pm\}) \times \{\pm\},$$

the Turing machine so that  $(b, \pm, a)$  is  $R$ -rejected, is not  $R$ -permitted, whenever  $b \neq a$  and otherwise

$$R(b, \pm, a) = \text{pr}_{\pm}(b, \pm, a)$$

for

$$\text{pr}_{\pm} : (B \times \{\pm\}) \times B \rightarrow \pm$$

the projection. Then clearly  $T$  computes  $D_M$ . The second part of the lemma is immediate.  $\square$



In particular by the example above there is a Turing machine

$$D_A : \mathcal{P} \times \mathbb{N} \rightarrow \{\pm\}$$

that stably soundly decides if a Diophantine polynomial has integer roots, meaning:

$$p \text{ is } D_A\text{-decided} \iff p \text{ has no integer roots.}$$

We can of course also construct such a  $D_A$  more directly. Likewise there is a Turing decision machine that stably soundly decides the halting problem.

**Definition 3.6.** *Given a machine*

$$M : \text{Dom} \subset \mathbb{N} \times B \rightarrow \{\pm\}$$

*and a Turing machine*

$$T : \text{Dom} \subset \mathbb{N} \times B \rightarrow \{\pm\},$$

*we say that  $T$  stably computes  $M$ , or that  $\Theta_{M,T}^s$  holds, if*

$$b \text{ is } M\text{-decided} \iff b \text{ is } T\text{-decided.}$$

#### 4. INCOMPLETENESS FOR STABLY SOUND TURING MACHINES

Let  $\mathcal{D}$  denote the set of abstract machines of the form:

$$D : \text{Dom} \subset \mathcal{T} \times \mathbb{N} \rightarrow \{\pm\},$$

with subset  $\text{Dom}$  not fixed. And set

$$\mathcal{D}^t := \{T \in \mathcal{T}_S \mid \text{fog}(T) \in \mathcal{D}\}.$$

In what follows, for  $T \in \mathcal{T}$  when we write  $T \in \mathcal{D}^t$  we mean that  $T \in \text{image } i|_{\mathcal{D}^t}$ , for  $i : \mathcal{T}_S \rightarrow \mathcal{T}$  the embedding discussed in Section 2.1. Likewise, if  $T \in \text{image } i|_{\mathcal{D}^t}$  then by  $T(T, m)$  we mean  $i^{-1}(T)(T, m)$ , so that in what follows the sentence “ $T$  is not  $T$ -decided” makes sense for such a  $T$ . Explicitly, for  $T \in \text{image } i|_{\mathcal{D}^t}$ ,  $T$  is not  $T$ -decided, will mean that for every  $n$  with  $i^{-1}(T)(T, n) = +$  there exists an  $m > n$  such that  $i^{-1}(T)(T, m) = -$ .

**Definition 4.1.** *For a  $D \in \mathcal{D}$ , we say that  $D$  is stably sound on  $T \in \mathcal{T}$  if*

$$(T \text{ is } D\text{-decided}) \implies (T \in \mathcal{D}^t) \wedge (T \text{ is not } T\text{-decided}).$$

*We say that  $D$  is stably sound if it is stably sound on all  $T$ . We say that  $D$  stably decides  $\mathcal{P}(T)$  if whenever the statement:*

$$(T \in \mathcal{D}^t) \wedge (T \text{ is not } T\text{-decided})$$

*is true,  $T$  is  $D$ -decided. We say that  $D$  stably soundly decides  $\mathcal{P}(T)$  if  $D$  is stably sound on  $T$  and stably decides  $\mathcal{P}(T)$ . We say that  $D$  stably soundly decides  $\mathcal{P}$  if  $D$  stably soundly decides  $\mathcal{P}(T)$  for all  $T \in \mathcal{T}$ .*

The interpretation of the above is that each  $D \in \mathcal{D}$  is understood as a machine with the properties:

- For each  $T, n$   $D(T, n) = +$  if  $D$  decides at the moment  $n$  that  $T \in \mathcal{D}^t$  and  $T$  is not  $T$ -decided.
- For each  $T, n$   $D(T, n) = -$  if  $D$  does not decide/assert at the moment  $n$  that  $T \in \mathcal{D}^t$ , or  $D$  does not decide that  $T$  is not  $T$ -decided.

**Lemma 4.2.** *If  $D$  is stably sound on  $T \in \mathcal{D}$  then*

$$\neg \Theta_{D,T}^s \vee \neg(T \text{ is } D\text{-decided}).$$

*Proof.* If

$$T \text{ is } D\text{-decided}$$

then since  $D$  is stably sound,  $T$  is not  $T$ -decided, so of course  $\neg \Theta_{D,T}^s$ . □

**Theorem 4.3.** *There is no (stably) computable  $D \in \mathcal{D}$  that stably soundly decides  $\mathcal{P}$ .*

*Proof.* Suppose otherwise that there is such a  $D$ , then by the above lemma we obtain:

$$\forall T \in \mathcal{T}^\pm : \neg \Theta_{D,T}^s,$$

but this is absurd since by assumption  $D$  is (stably) computable.  $\square$

**Definition 4.4.** For  $D \in \mathcal{D}^t$ , we say that  $\mathcal{R}(D)$  holds if for any  $T \in \mathcal{D}^t$  such  $T$  is not  $T$ -decided:

$$\exists T' \in \mathcal{D}^t : (D \text{ stably decides } \mathcal{P}(T')) \wedge (T \simeq_s T').$$

**Theorem 4.5.** For  $D \in \mathcal{D}$  the following cannot hold simultaneously:  $D$  is stably sound,  $D$  is (stably) computable and  $\mathcal{R}_D$  holds.

*Proof.* Suppose that  $D$  is stably computed by some  $T \in \mathcal{D}^t$ . If  $D$  is stably sound then by Lemma 4.2

$$\neg(T \text{ is } D\text{-decided}),$$

and so

$$\neg(T \text{ is } T\text{-decided}),$$

since  $T$  stably computes  $D$ . Consequently, if  $\mathcal{R}(D)$  then for some  $T' \in \mathcal{D}^t$  stably computing  $T$ ,  $T'$  is  $D$ -decided. And so by Lemma 4.2  $\neg \Theta_{D,T'}^s$ , which is a contradiction.  $\square$

*Proof of Theorem 1.2.* Suppose otherwise that we have such an  $M$ , so  $M$  is stably sound, computable and such that

$$\forall \Sigma \in \mathcal{A} : \Sigma \text{ is true} \iff \Sigma \text{ is } M\text{-stable}.$$

It is well known that sentences in the language of Turing machines are logically equivalent to sentences in arithmetic. This means that given a sentence  $\Sigma^t$  in the language of Turing machines, e.g. “ $T$  halts on input  $n$ ” there is a sentence  $\Sigma \in \mathcal{A}$  so that  $\Sigma^t$  is true iff  $\Sigma$  is true. Indeed this kind of translation already appears in the original work of Turing [1].

Let  $s(T) \in \mathcal{A}$  be the sentence logically equivalent to

$$(T \in \mathcal{D}^t) \wedge (T \text{ is not } T\text{-decided}).$$

Define a machine  $\tilde{D}_M \in \mathcal{D}$  by

$$\tilde{D}_M(T, n) := D_M(s(T), n)$$

for  $D_M$  defined as in Section 3. Then by Lemma 3.5 and by supposition above  $D_M \in \mathcal{D}$  is stably sound, computable, and  $D_M$  stably decides  $\mathcal{P}$ . But this contradicts Theorem 4.3.  $\square$

## 5. APPLICATION: A HUMAN AS A DECISION MACHINE

Let  $S$  be a human subject in a controlled environment, in communication with an experimenter/operator  $E$  that as input passes to  $S$  elements of  $\mathcal{I} = \mathcal{T} \times \mathbb{N}$ . Here *controlled environment* means primarily that no information that is not explicitly controlled by  $E$  and that is usable by  $S$  passes to  $S$  while they are in this environment. This condition is only for simplicity, so long as we know in principle what “input” our  $S$  receives it doesn’t matter what kind of environment they are in. For practical purposes  $S$  has in their environment a general purpose digital computer with arbitrarily, as necessary, expendable memory, (in other words a universal Turing machine).

We suppose that upon receiving any  $I \in \mathcal{I}$ , (for example as data in their computer) after possibly using their computer in some way,  $S$  instructs their computer to print after some indeterminate time an element

$$O \in \{\pm\}.$$

So  $S$  is meant to determine a decision machine:

$$(5.1) \quad D_S : \mathcal{T} \times \mathbb{N} \rightarrow \{\pm\},$$

interpreted as in the previous section.

**Remark 5.2.** *The above is partially a simplification, because for a real world  $S$  it may be that each  $D_S(I)$  must be understood as a probability distribution on  $\{\pm\}$ . In other words the value  $D_S(I)$  may only be determined up to some dice roll, which we may expect if quantum mechanics plays a significant role in the human brain. This extra complexity will be ignored, as it does not meaningfully change any of our arguments since dice rolls can be simulated completely with Turing machines. Moreover, we are only interested in the stable behaviour of  $D_S$ , i.e. the set of  $D_S$ -decided elements of  $\mathcal{T}$ . And this stable behavior, given our interpretation of  $D_S$  as a decision machine, should not be affected by any dice rolls.*

In what follows when we say “ $S$  perceives” we mean that some (implied) machine associated to our physical  $S$  stably decides, where this has the same meaning as previously.

If  $D_S$  is stably sound and computable then by Theorem 4.5  $\neg\mathcal{R}(D_S)$ . And this in particular means that there exists a  $T_S \in \mathcal{D}^t$  so that the statement

$$(T_S \in \mathcal{D}^t) \wedge (T_S \text{ is not } T_S\text{-decided})$$

is true but  $S$  will never perceive it to be true. On the other hand  $\mathcal{R}(D_S)$  is of course implied by  $S$  being able to perceive, for any given Turing machine

$$f : \mathbb{N} \rightarrow \{\pm\},$$

that no  $n$  is  $f$ -decided if this is in fact true. However, the condition  $\mathcal{R}(D_S)$  is likely weaker since for  $T_S$  as above  $T_S(T_S, \cdot)$  stably computes the partial function

$$u : \{\emptyset\} \subset \mathbb{N} \rightarrow \{\pm\},$$

in other words  $u$  is defined nowhere. One can then hope, since  $u$  is so simple, that  $T_S$  can be put into a normal form  $T'$  with  $T_S \simeq_s T'$  and with  $T'(T', \cdot) = u$  and so that  $S$  can perceive that  $T'(T', \cdot) = u$ .

In any case, if we define  $\mathcal{H}(S) := s(T_S)$ , with the latter defined as in the proof of Theorem 1.2, then we obtain our Theorem 1.3.

**5.1. Relationship with the Gödel and Penrose argument.** The most lucid analysis, known to me, of the Gödel and Penrose arguments for obstructions to computability of intelligence appears in Koellner [13], [14]. Our argument for Theorem 1.3 extends Gödel’s idea [11, 310], although I am also inspired by the ideas of Penrose. One note is that our argument is entirely based on set theory, while Gödel’s argument has meta-logical elements that require interpretation. Although, as Koellner explains [14], Gödel’s argument can also be at least in some sense fully formalized.

This is not to say that there are no issues of interpretation in Theorem 1.3. One must interpret our definition of stable soundness as it applies to actual human beings. We of course have already partly addressed this. Scientists operate on the unshakeable faith that scientific progress converges on truth. And our interpretation above of this convergence as stable soundness is very simple and natural, at least under the previously explained assumption of weak idealization.

It is of course always the case that we must interpret mathematical theorems when applied to the real world. What one looks for is whether there is any meaningful physical obstruction to carrying out the necessary idealization in principle. In our specific case I see no such obstruction. Of course if the universe and humanity must eventually go extinct then our weakly idealized humans cannot even in principle exist. But to me this is not a meaningful obstruction. The potential mortality of the universe is very unlikely to have any causal relation with computability of intelligence. So we can imagine an eternal universe and a weakly idealized human, run the argument then translate to our universe.

So the only thing to reasonably wonder is whether for every  $S$  there must be such a stably undecidable arithmetic statement of the form  $\mathcal{H}(S)$  above.

**Acknowledgements.** Dennis Sullivan, Bernardo Ameneyro Rodriguez, David Chalmers, and in particular Peter Koellner for helpful discussions on related topics.

## REFERENCES

- [1] A.M. TURING, *On computable numbers, with an application to the entscheidungsproblem*, Proceedings of the London mathematical society, s2-42 (1937).
- [2] ———, *Computing machines and intelligence*, Mind, 49 (1950), pp. 433–460.
- [3] L. BLUM, M. SHUB, AND S. SMALE, *On a theory of computation and complexity over the real numbers: NP-completeness, recursive functions and universal machines.*, Bull. Am. Math. Soc., New Ser., 21 (1989), pp. 1–46.
- [4] D. J. CHALMERS, *Minds machines and mathematics*, Psyche, symposium, (1995).
- [5] D. DEUTSCH, *Quantum theory, the Church-Turing principle and the universal quantum computer.*, Proc. R. Soc. Lond., Ser. A, 400 (1985), pp. 97–117.
- [6] S. FEFERMAN, *Are There Absolutely Unsolvable Problems? Gödel’s Dichotomy*, Philosophia Mathematica, 14 (2006), pp. 134–152.
- [7] C. FIELDS, D. HOFFMAN, C. PRAKASH, AND M. SINGH, *Conscious agent networks: Formal analysis and application to cognition*, Cognitive Systems Research, 47 (2017).
- [8] P. GRINDROD, *On human consciousness: A mathematical perspective*, Network Neuroscience, 2 (2018), pp. 23–40.
- [9] S. HAMEROFF AND R. PENROSE, *Consciousness in the universe: A review of the ‘orch or’ theory*, Physics of Life Reviews, 11 (2014), pp. 39 – 78.
- [10] J.R. LUCAS, *Minds machines and Gödel*, Philosophy, 36 (1961).
- [11] K. GÖDEL, *Collected Works III* (ed. S. Feferman), New York: Oxford University Press, 1995.
- [12] A. KENT, *Quanta and qualia*, Foundations of Physics, 48 (2018), pp. 1021–1037.
- [13] P. KOELLNER, *On the Question of Whether the Mind Can Be Mechanized, I: From Gödel to Penrose*, Journal of Philosophy, 115 (2018), pp. 337–360.
- [14] ———, *On the question of whether the mind can be mechanized, ii: Penrose’s new argument*, Journal of Philosophy, 115 (2018), pp. 453–484.
- [15] K. KREMNIZER AND A. RANCHIN, *Integrated information-induced quantum collapse*, Foundations of Physics, 45 (2015), pp. 889–899.
- [16] R. PENROSE, *Beyond the shadow of a doubt*, Psyche, (1996).

UNIVERSITY OF COLIMA, DEPARTMENT OF SCIENCES, CUICBAS  
 Email address: yasha.savelyev@gmail.com