

# INCOMPLETENESS FOR STABLY CONSISTENT FORMAL SYSTEMS

YASHA SAVELYEV

**ABSTRACT.** We first partly develop a mathematical notion of stable consistency intended to reflect the actual consistency property of human beings. Then we give a generalization of the first and second Gödel incompleteness theorem to stably 1, 2-consistent formal systems. Our argument in particular re-proves the original incompleteness theorems from first principles, using Turing machine language to (computably) construct our “Gödel sentence” directly, in particular we do not use the diagonal lemma, nor any meta-logic, with the proof naturally formalizable in set theory. In practice such a stably consistent formal system could be meant to represent the mathematical output of humanity evolving in time, so that the above gives a formalization of a famous disjunction of Gödel, obstructing computability of intelligence.

## 1. INTRODUCTION

We begin by quickly introducing the notion of stable consistency. First, the term **encoded map** will mean a partial map  $M : A \rightarrow B$ , with  $A, B$  sets with an additional structure of an encoding in  $\mathbb{N}$ , (later on  $\mathbb{N}$  is replaced by the set of strings in a finite alphabet as that will be more natural). An encoding is just an injective map  $e : A \rightarrow \mathbb{N}$  with some extra properties. This is described in more detail in Section 2.2. Working with encoded sets and maps as opposed to subsets of  $\mathbb{N}$  will have the same advantages as working with abstract countably dimensional vector spaces as opposed to subspaces of  $\mathbb{R}^\infty$ .

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

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

for  $\{\pm\}$  a set with two elements  $+, -$ .

### Definition 1.1.

- $\alpha \in \mathcal{A}$  is  **$M$ -stable** if there is an  $n_0$  with  $M(n_0) = (\alpha, +)$  s.t. there is no  $m > n_0$  with  $M(m) = (\alpha, -)$ . Let  $M^s : \mathbb{N} \rightarrow \mathcal{A}$  enumerate in order of appearance the  $M$ -stable  $\alpha$ . Abusing notation, we may also just write  $M^s$  for the set  $M^s(\mathbb{N})$ , where there is no risk of confusion. To keep notation manageable we write  $M \vdash \alpha$  to mean  $M^s \vdash \alpha$  meaning  $M^s$  proves  $\alpha$ , since there can be no other meaning.
- $M$  is **stably consistent** if  $M^s$  is consistent.
- $M$  **decides arithmetic** if

$$\forall \alpha \in \mathcal{A} : (M \vdash \alpha) \vee (M \vdash \neg \alpha),$$

where  $\vdash$  means proves as usual.

- $M$  **decides arithmetic truth** if

$$\forall \alpha \in \mathcal{A} : (\alpha \text{ is true}) \implies (M \vdash \alpha),$$

here by true we mean satisfied in the standard model of arithmetic.

So such an  $M$  could be given the following interpretation, this interpretation has no mathematical content a priori, it is simply how one may think of such an  $M$  informally.  $M(n) = (\alpha, +)$  only if at the moment  $n$   $M$  decides that  $\alpha$  is true. Meanwhile

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

only if at the moment  $m$ ,  $M$  no longer asserts that  $\alpha$  is true. If  $\alpha$  is  $M$ -stable we can also say  $M$  *stably asserts*  $\alpha$  to be true. So that  $M$  is stably consistent if the set of sentences it stably decides to be true are consistent.

The following is the semantic version of our main theorem.

**Theorem 1.2.** *Suppose  $M$  as above is stably computable and stably sound then there is a constructible (given a Turing machine stably computing  $M$ ) true sentence  $\alpha(M)$ , which  $M^s$  does not prove. Here stably computable means that there is a Turing machine  $T : \mathbb{N} \rightarrow \mathcal{A} \times \{\pm\}$  so that  $M^s = T^s$ , note that this does not mean that  $M^s$  is computable.*

The fact that a true sentence  $\alpha$  which  $M^s$  does not prove exists, can be immediately deduced from Tarski undecidability of truth as the set  $M^s(\mathbb{N})$  is definable in first order arithmetic whenever  $M$  is stably computable. However our sentence is constructible by elementary means, starting with the definition of a Turing machine. In particular the diagonal lemma is not used. More crucially this proof readily extends to give the more interesting syntactic version of the theorem, generalizing the original Gödel incompleteness theorems, and which is presented below.

Let  $RA$  denote Robinson arithmetic that is Peano arithmetic  $PA$  without induction. Let  $\mathcal{F}_0$  denote the set of  $RA$ -decidable formulas  $\phi$  in arithmetic with one free variable.

**Definition 1.3.** *Given  $M : \mathbb{N} \rightarrow \mathcal{A} \times \{\pm\}$ , we say that it is **stably 1-consistent** if it is stably consistent and if for any formula  $\phi \in \mathcal{F}_0$  the following holds:*

$$M \vdash \exists m : \phi(m) \implies \neg \forall m : M \vdash \neg \phi(m).$$

*We say that it is **stably 2-consistent** if the same holds for  $\Pi_1$  formulas  $\phi$  with one free variable, that is formulas  $\phi = \forall n : g(m, n)$ , with  $g$   $RA$ -decidable.*

**Theorem 1.4.**

- (1) *Let  $M$  be stably computable and such that  $M^s \supset RA$ . Then there is a computably constructible, from a specification of a Turing machine stably computing  $M$ , sentence  $\alpha(M) \in \mathcal{A}$  which  $M^s$  cannot prove if it is stably 1-consistent and cannot disprove if it is stably 2-consistent.*
- (2)  *$\alpha(M)$  can be chosen so that for  $M$  stably computable by a Turing machine  $T$ , so that  $T^s \supset RA$ .*

$$(T \text{ is stably 1-consistent}) \implies \alpha(M)$$

*is a theorem of  $PA$  taking the standard interpretation of all terms. In particular  $\alpha(M)$  is true in the standard model of arithmetic whenever  $M$  is stably 1-consistent. And in particular, if in addition  $M^s \supset PA$ , then*

$$M \not\vdash (T \text{ is stably 1-consistent}),$$

*that is  $M^s$  cannot prove its own 1-consistency.*

This can be understood as a direct generalization of the original first Gödel incompleteness theorem, as we simply weaken the main assumption from  $\omega$ -consistency to 1, 2-consistency, and computability to stable computability. And it is an extension of the second Gödel incompleteness theorem to only stably computably enumerable, but 1-consistent formal systems extending  $RA, PA$ . The above theorem is naturally interpreted in set theory  $ZF$ , as just a statement about particular set theoretic functions, and our proof is readily formalizable in  $ZF$ . This should make the above more accessible to general mathematical audience.

**1.1. Relationship with known results.** If we choose to look at this theorem purely from the point of view of the formal system  $M^s$ , then the distinction with Gödel is that the set  $M^s$  is merely  $\Sigma_2$  definable when  $M$  is stably computable, keeping in mind that we do not assume that  $M^s$  is deductively closed. As pointed out in [17] there are examples of complete, consistent  $\Delta_2$  definable extensions of  $PA$ . So that any extensions of the second Gödel incompleteness theorem to definable sets must require stronger consistency assumptions. Theorem 1.4 2) says that 1-consistency is sufficient, when we have a stably computably enumerable set, that is a set  $S$  so that there is a Turing machine  $T : \mathbb{N} \rightarrow \mathbb{N} \times \{\pm\}$  so that  $T^s(\mathbb{N}) = S$ , see Section 3 for definitions. In connection with this we should ask a simple question:

*Question 1.* Is every  $\Sigma_2$  definable set  $S \subset \mathbb{N}$  stably computable enumerable?

Moreover, in the very recent work of Salehi and Seraji [17], which we also recommend for additional references, the general Corollary 2.11 in part implies the first part of Theorem 1.4 at least when  $M^s$  is deductively closed. One more substantial distinction is that our sentence is constructible, which aside from mathematical interest is important for interpretation here, since  $M$  is meant to model the mathematical output of humanity and we want to formalize Gödel's disjunction, this is discussed in the following Section 1.2. Surprisingly, the authors of [17] point out that in general when  $n$  is at least 3 for a  $\Sigma_{n+1}$  definable,  $n$ -consistent theory  $F$  there is no constructible  $\Pi_{n+1}$  independent sentence, although such a sentence does exist. It is possible that the result above is related to the boundary.

On the other hand I am not aware of any previous work on generalizations of the second incompleteness theorem, to formal systems which are not computably enumerable, as we have in the second part of the Theorem above.

Our argument also readily reproves the original first and second incompleteness theorems of Gödel from first principles<sup>1</sup>, with our Gödel sentence constructed directly by means of Turing machine language.

**1.2. 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. Given this *very* limited interpretation, the question that we are interested in is this:

*Question 2.* 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. 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. However, we cannot honestly even hope for consistency as for example mathematicians are not on the surface consistent at all times. But we can certainly hope for some more fundamental consistency and soundness as human knowledge appears to converge on truth.

We will interpret here fundamental consistency/soundness as stable consistency/soundness. This notion is meant to reflect the basic understanding of the way science progresses. Of course even stable consistency 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, or  $M$  now refers not to an individual human but to humanity, suitably interpreted. We call such a human *weakly idealized*.

**Remark 1.5.** In the case of humanity  $H$ , we may suppose that the output of the associated function

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

is determined by majority consensus. This is not as restrictive as it sounds, for example if there is a computer verified proof of the Riemann hypothesis  $\alpha$  in ZFC, then irrespectively of the complexity of the proof, we can expect majority consensus for  $\alpha$ , provided validity of set theory still has consensus. At least if we reasonably interpret  $H$ , which is beyond the scope here. In addition if stable consistency (or

<sup>1</sup>Assuming existence of a certain encoding category  $\mathcal{S}$ , which is already essentially constructed by Gödel.

perhaps stable 1-consistency) is explicitly the goal, (we can say that this is the experimental branch of human mathematical output) then  $H$  can safely consider increasingly more powerful axiomatic systems like  $ZFC + \text{continuum hypothesis}$ , etc. By the same reasoning it is not at all unreasonable to suppose that  $H$  stably asserts all theorems of say Robinson arithmetic, it can just be interpreted as that  $H$  stably asserts all theorems generated by some particular Turing machine. This remark is primarily of interest in the context of Theorem 1.7 to follow.

1.2.1. *Gödel's disjunction.* 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. Let  $S$  be a certain idealized mathematician. Essentially, Gödel argues for a disjunction, the following cannot hold simultaneously:

$S$  is computable,  $S$  is consistent and  $A$ ,

where  $A$  says that  $S$  can decide any Diophantine problem. 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.

*Claim 1.* It is impossible to meaningfully formalize Gödel's disjunction without strengthening the incompleteness theorems.

Already Feferman asks in [6] what is the meaning of 'idealized' above? In the context of the present paper we might propose that it just means stabilized (cf. Section 3 specifically). But there is a Turing machine  $T$  whose stabilization  $T^s$  soundly 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 and mathematically. The process of the idealization must be such that non-computability is not introduced in the ideal limit. For weak idealization in terms of humanity mentioned above this is automatic, for the more direct idea of idealizing brain processes it should certainly also be in principle possible. For example, suppose we know that there is a mathematical model  $M$  for the biological human brain, in which the deterioration mentioned above is described by some explicit computable stochastic process, on top of the base cognitive processes. Then mathematically a weak idealization  $M^i$  of  $M$  would just correspond to the removal of this stochastic process. We may then meaningfully apply Theorems 1.2, 1.4 to  $M^i$ , since  $M^i$  would be non-computable only if some cognitive processes of the brain (in the given mathematical model) were non-computable, since mathematically  $M^i$  would be composed from such processes.

But this is not what is needed by Gödel. He needs an idealization that is plausibly consistent otherwise the disjunction would again be meaningless, while a weak idealization of a human is only plausibly stably consistent. Since we do not have a good understanding of physical processes of the human brain involved in cognition, it is not at all clear that what Gödel asks is even possible.

So to meaningfully formalize Gödel's disjunction we need generalize Gödel's incompleteness theorems to use stable consistency, or somewhat stronger notions like stable 1-consistency, which is what we do here.

1.2.2. *Lucas-Penrose thesis.* 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 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 construction of  $P$  is that  $P$  be not computable as a function! So on this rough surface level this appears to be impossible. Notwithstanding, we can start by formalizing Gödel's disjunction as follows.

1.2.3. *Formalizing Gödel's disjunction.* The following is an applied interpreted version of our Theorem 1.4. In what follows,  $H$  refers to an encoded map associated to humanity as discussed. We will not say exactly how  $H$  is associated, as this is somewhat beyond our scope. We will instead assume the following axiom:

**Axiom 1.6.** *If the physical processes underlying the output of  $H$  are computable, this may include computability of certain stochastic variables, as probability distribution valued maps, then the set  $H^s(\mathbb{N})$  is stably computably enumerable. This means as before that there is a Turing machine  $T : \mathbb{N} \rightarrow \mathcal{A} \times \{\pm\}$  so that  $T^s(\mathbb{N}) = H^s(\mathbb{N})$ .*

Note that we cannot hope for  $H$  to be computable, since particular values  $H(n)$  may arise from random phenomena. We could hope that some probability distribution valued version of  $H$  is computable, but that is another story.

So this is our present concrete model for a weakly idealized human, we theorized above other perhaps neater models for a weakly idealized human by directly idealizing brain processes, and the theorem could also be stated in that context.

**Theorem 1.7.** *Suppose that our model for humanity  $H : \mathbb{N} \rightarrow \mathcal{A} \times \{\pm\}$ , stably asserts all theorems of Robinson arithmetic, is stably 1-consistent (or just stably sound for simplicity), and satisfies the above axiom. Then either there are cognitively meaningful, absolutely non Turing computable processes in the human brain, or there exists a true computably constructive statement of arithmetic  $\mathcal{H}$ , such that  $H$  will never stably decide  $\mathcal{H}$ . By constructive we mean provided a specification of a Turing machine stably computing  $H$ .*

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>2</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.

**Remark 1.8.** *Turing suggested in [2] that abandoning hope of soundness is the obvious way to circumvent the implications of Gödel incompleteness theorem. In my opinion this position is untenable, humans may not be sound but it is implausible that  $H$  is not stably sound, since as we say human knowledge appears to converge on truth, and it is inconceivable that  $H$  is not stably 1-consistent. Then we cannot escape incompleteness by the above. So given stable soundness, the only way that the incompleteness theorems can be circumvented is to accept that there is an absolute limit on the power of human reason as in the theorem above.*

**Remark 1.9.** *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].*

## 2. SOME PRELIMINARIES

For more details on Turing machines we recommend the book of Soare [18]. Our approach here is however is slightly novel in that we do not fix Gödel encoding for anything, instead abstractly

---

<sup>2</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  $\mathbb{Z}$ .

axiomatizing the expected properties of encodings. This allows us later on to give very concise, self-contained arguments for main results.

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  $s$ ,  $\delta^M(s)$  will denote the successor state of  $s$ , obtained by executing the instructions set of  $M$  on  $s$ , or in other words  $\delta^M(s)$  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, it will be convenient to 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

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

The set  $\mathcal{I}$  is also called the set of *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  $T \in \mathcal{T}$

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

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

**2.1. Multi-input Turing machine.** There is a basic well known variant of a Turing machine which takes as input an element of  $\text{Strings}^n$ , for some fixed  $n \in \mathbb{N}$ . This is done by replacing the input 1-tape by an  $n$ -tape. We will not give details of this. Notationally such a Turing machine will be distinguished by a superscript so

$$T^n : \text{Strings}^n \rightarrow \text{Strings},$$

denotes such a Turing machine with  $n$  inputs.

**2.2. Abstractly encoded sets.** The material of this section will be used in the main arguments, it will allow us to remove the need to work with explicit Gödel encodings, greatly simplifying subsequent details. However this will require a bit of abstraction. This abstraction is analogous to introduction of abstract vector spaces in linear algebra. This is likely very obvious to experts working with things like type theory, but I am not aware of this being explicitly introduced in computability theory literature.

An *encoding* of a set  $A$  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 a small 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 set  $A$ . We may denote  $e_A(A)$  by  $A_e$ . The morphisms in the category  $\mathcal{S}$  are pairs of partial maps  $(m, m_e)$  so that the following diagram commutes:

$$\begin{array}{ccc} A & \xrightarrow{m} & B \\ \downarrow e_A & & \downarrow e_B \\ A_e \subset \text{Strings} & \xrightarrow{m_e} & B_e \subset \text{Strings}, \end{array}$$

and so that  $m_e$  is computable and  $A_e$  coincides with the set of permissible strings for the Turing machine computing  $m_e$ .



**Notation 1.** 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.

- (1) For  $A \in \mathcal{S}$   $A_e$  is computable (recursive). Here, as is standard, 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  $\Sigma \in S$ .
- (2) 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.
- (3) For  $A, B \in \mathcal{S}$ ,

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

In particular each  $A \in \mathcal{S}$  is determined by  $A_e$ .

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

$$\begin{array}{ccc} A \times B & \xrightarrow{pr^A} & A \\ \downarrow & & \downarrow \\ (A \times B)_e & \xrightarrow{pr_e^A} & A_e, \end{array}$$

with  $pr_e^A$  computable, similarly for  $pr^B$ . In addition we ask for the following naturality property. (We don't strictly speaking use this last property, but it may be helpful for understanding  $\mathcal{S}$ .) Let  $f$  be the composition

$$A_e \times B_e \xrightarrow{(e_A^{-1}, e_B^{-1})} A \times B \xrightarrow{e_{A \times B}} (A \times B)_e,$$

then  $f$  is computable meaning that there is a 2-input Turing machine  $T^2$ , with input set  $A_e \times B_e$  such that

$$T^2(\Sigma_1, \Sigma_2) = f(\Sigma_1, \Sigma_2)$$

for all  $(\Sigma_1, \Sigma_2)$  in  $A_e \times B_e$ .

The above axioms suffice for our purposes but there are a number of possible extensions. 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  $\mathbb{Z}, \mathbb{N}, \mathcal{A}, \mathcal{T}, \{\pm\}$ , with  $\{\pm\}$  a set with two elements  $+, -$ . The encodings of  $\mathbb{N}, \mathbb{Z}$  should be suitably natural so that for example the map

$$\mathbb{N} \rightarrow \mathbb{N}, \quad n \mapsto n + 1$$

completes to a morphism in  $\mathcal{S}$ . For  $\mathbb{Z}$  we also want the map

$$\mathbb{Z} \rightarrow \mathbb{Z} \quad n \mapsto -n$$

to complete to a morphism in  $\mathcal{S}$ . The fact that such encoding categories  $\mathcal{S}$  exist is essentially a classical theorem of Gödel [11], with extensions by Turing and others. It is also implicitly part of any typed computer language.

**Definition 2.6.** *An abstract Turing machine*

$$T : A \rightarrow B,$$

is a morphism of  $\mathcal{S}$ , so that in addition we specify a Turing machine  $T_e$  computing  $e_B \circ T \circ e_A^{-1}$ . So that really this is a pair  $(T, T_e)$ . We may just say Turing machine in place of abstract Turing machine, since usually there can be no confusion, an abstract Turing machine  $\text{Strings} \rightarrow \text{Strings}$  is just a Turing machine.



We define  $\mathcal{T}_S$  to be the set of abstract Turing machines relative to  $S$  as above.  $\mathcal{T}_S$  in general cannot be naturally encoded<sup>3</sup>, but 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.** An **encoded map** will be a synonym for a partial map  $M : A \rightarrow B$ , with  $A, B$  abstractly encoded sets.

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

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

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  $fog(T) = M$ . We say that  $M$  is **computable** if some  $T$  computes  $M$ . Thus the set  $hom_S(A, B)$  of morphisms in  $S$  from  $A$  to  $B$  is exactly the set of computable encoded maps  $M : A \rightarrow B$  with the definitions above.

**2.3. 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. We will not notationally distinguish naturals  $n \in \mathbb{N}$  from their corresponding numerals in the language of arithmetic, as this usually will not lead to any confusion.

### 3. ON STABLE CONSISTENCY AND SOUNDNESS

**Definition 3.1.** Given an encoded map:

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

We say that  $b \in B$  is  **$M$ -stable** if there is an  $n_0$  with  $M(n_0) = (b, +)$  s.t. there is no  $m > n_0$  with  $M(m) = (b, -)$ . In the case above, we may also say that  $M$  **prints  $b$  stably**.

**Definition 3.2.** Given an encoded map

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

we define

$$M^s : \mathbb{N} \rightarrow B$$

to be the encoded map (which recall is by definition partial) enumerating, in order, all the  $M$ -stable  $b$ . 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 an abstract Turing machine

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

whose stabilization  $A^s$  enumerates every Diophantine (integer coefficients) polynomial with no integer roots, where  $Pol$  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.

---

<sup>3</sup>The key word is 'naturally' as any countable set can be encoded.

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,

$$Z : \mathbb{N} \rightarrow \text{Pol}, N : \mathbb{N} \rightarrow \mathbb{Z}$$

are total bijective computable maps. The encoding of  $\text{Pol}$  should be suitably natural so that in particular the map

$$E : \mathbb{Z} \times \text{Pol} \rightarrow \mathbb{Z}, \quad (n, p) \mapsto p(n)$$

is computable. In what follows, for each  $n \in \mathbb{N}$ ,  $L_n$  has the type of an ordered finite list of elements of  $\text{Pol} \times \{\pm\}$ , with order starting from 0.

- Initialize  $L_0 := \emptyset$ ,  $n := 0$ .
- Start. For each  $p \in \{Z(0), \dots, Z(n)\}$  check if one of  $\{N(0), \dots, N(n)\}$  is a solution of  $p$ , if no add  $(p, +)$  to  $L_n$ , if yes add  $(p, -)$ . Call the resulting list  $L_{n+1}$ . Explicitly,

$$L_{n+1} := L_n \cup \bigcup_{m=0}^n (Z(m), d^n(Z(m))),$$

where  $d^n(p) = +$  if none of  $\{N(0), \dots, N(n)\}$  are roots of  $p$ ,  $d^n(p) = -$  otherwise,

where  $\cup$  is set union operation of subsets of  $\text{Pol} \times \{\pm\}$ .

- Set  $n := n + 1$  go to Start and continue.

Set  $L := \cup_n L_n$ , which is an ordered infinite list. Define

$$A : \mathbb{N} \rightarrow \text{Pol} \times \{\pm\}$$

by  $A(m) := L(m)$ , with the latter the  $m$ 'th element of  $L$ . Since  $E$  is computable, it clearly follows that  $A$  is computable and its stabilization  $A^s$  enumerates Diophantine polynomials which have no integer roots.

**3.1. Decision maps.** By a **decision map** we mean an encoded map of the form:

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

**Definition 3.4.** For a  $D$  as above we say that  $b \in B$  is  **$D$ -decided** if there is an  $N$  s.t. for all  $n \geq N$   $D(b, n) = +$ .

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

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

which is defined as follows.

**Definition 3.5.** Define  $b \in B$  to be  **$(M, n)$ -stable** if there a  $m$ ,  $0 \leq m \leq n$ , with  $M(m) = (b, +)$  s.t. there is no  $k$  satisfying  $m < k \leq n$ , with  $M(k) = (b, -)$ .

$D_M$  is then defined by

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

**Lemma 3.6.** If  $M$  as above is (stably) computable then  $D_M$  is (stably) computable. Moreover,  $b$  is  $M$ -stable iff  $b$  is  $D_M$ -decided.

*Proof.* Suppose that  $T$  computes  $M$ . Define

$$\mathcal{B}^T = \{(b, n) \in B \times \mathbb{N} \mid b \text{ is } (T, n)\text{-stable}\},$$

then  $\mathcal{B}^T$  is obviously Turing decidable and so  $D_T$  is computable and  $D_M$  is computable. Now suppose that  $T$  stably computes  $M$  then again  $D_T$  is computable by a Turing machine  $Z$ , but clearly  $Z$  stably computes  $D_M$ . The last part of the lemma is immediate.  $\square$

In particular by the example above there is a Turing machine

$$D_A : Pol \times \mathbb{N} \rightarrow \{\pm\}$$

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, in this sense.

**Definition 3.7.** *Given an encoded map*

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

*and a Turing machine*

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

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

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

**Definition 3.8.** *Given Turing machines*

$$T_1, T_2 : B \times \mathbb{N} \rightarrow \{\pm\},$$

*we say that they are **stably equivalent** and write  $T_1 \simeq_s T_2$  if  $T_1, T_2$  stably compute the same encoded map.*

**3.2. Preliminaries on arithmetic decision maps.** In what follows  $M$  denotes a total encoded map  $M : \mathbb{N} \rightarrow \mathcal{A} \times \{\pm\}$ .

**Lemma 3.9.** *Given  $M$ , there is an encoded map  $CM : \mathbb{N} \rightarrow \mathcal{A} \times \{\pm\}$  so that  $CM^s(\mathbb{N})$  is the deductive closure of the set  $M^s(\mathbb{N})$ , and so that if  $M$  is (stably) computable then so is  $CM$ .*

In this stable setting the formal argument requires some care.

*Proof.* Let

$$T : \mathbb{N} \rightarrow \mathcal{A},$$

be a computable encoded map, let  $\mathcal{P} = \mathcal{P}(T)$  be the set of proofs in the formal system  $T(\mathbb{N})$ , also abbreviated by just  $T$ , and let  $Fin(T)$  denote the set of ordered finite lists of sentences in the formal system  $T$ . These are meant to be naturally encoded. Explicitly,  $Fin(T)$  has the encoding so that

$$P_i : Fin(T) \rightarrow \mathcal{A},$$

$P_i(f)$  being the  $i$ 'th element of the list  $f$ , is computable.  $\mathcal{P}$  is encoded so that the map:

$$G : \mathcal{P} \rightarrow Fin(T),$$

$G(p)$  is the list of sentences of  $T$  appearing in the proof  $p$  is computable.

There is a total computable encoded map

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

with the property that  $pr_{\mathcal{A}} \circ PT$  enumerates the deductive closure of  $T(\mathbb{N})$ . And so that, for each  $n$ ,  $pr_{\mathcal{P}} \circ PT(n)$  is the proof of  $pr_{\mathcal{A}} \circ PT(n)$  in the formal system  $\{T(0), \dots, T(n)\}$ . Here,

$$pr_{\mathcal{A}} : \mathcal{A} \times \mathcal{P} \rightarrow \mathcal{A}, \quad pr_{\mathcal{P}} : \mathcal{A} \times \mathcal{P} \rightarrow \mathcal{P}$$

are the natural projections. Existence of  $PT$  is elementary and well known.

Also denote by

$$pr_{\mathcal{A}} : \mathcal{A} \times \{\pm\} \rightarrow \mathcal{A}, \quad pr_{\pm} : \mathcal{A} \times \{\pm\} \rightarrow \{\pm\}$$

the pair of projections, set

$$M' := pr_{\mathcal{A}} \circ M.$$

**Definition 3.10.**  $\alpha \in \mathcal{A}$  will be called  **$n$ -stable** (with respect to  $M$ ) if the following holds.

- There exists  $m \leq n$  s.t.  $PM'(m) = (\alpha, p)$  where  $p$  is a proof of  $\alpha$  in  $F = \{\alpha_1, \dots, \alpha_k\}$  with  $\alpha_1, \dots, \alpha_k \in M'(\mathbb{N})$ .
- All  $\alpha_i$ ,  $1 \leq i \leq k$  as above satisfy:

$$\exists d \in \{0, \dots, n\} : M(d) = (\alpha_i, +) \wedge \forall m : d \leq m \leq n \implies M(m) \neq (\alpha_i, -).$$

For each  $n \in \mathbb{N}$ ,  $L_n$  will have the type of an ordered finite list of elements in  $\mathcal{A} \times \{\pm\}$ .

- (1) Initialize  $L_0 = \emptyset$  and  $n := 0$ .
- (2) Start. Set

$$U_n := \{pr_{\mathcal{A}}\sigma \mid \sigma \in L_n, pr_{\pm}\sigma = +, pr_{\mathcal{A}}\sigma \text{ is not } n\text{-stable}\}.$$

Set

$$L'_{n+1} := L_n \cup \bigcup_{\alpha \in U_n} \{(\alpha, -)\}$$

and set

$$L_{n+1} := L'_{n+1} \cup \{(pr_{\mathcal{A}} \circ PM'(n), +)\},$$

if  $pr_{\mathcal{A}} \circ PM'(n)$  is  $n$ -stable, otherwise set

$$L_{n+1} := L'_{n+1}.$$

- (3)  $n := n + 1$ , go to Start and continue.

The above recursion gives an ordered infinite list  $L := \cup_n L_n$ , for each  $m$  set  $CM(m)$  to be the  $m$ -th element of the list  $L$ .  $\square$

**Definition 3.11.** Let  $\mathcal{F}_0$  as in introduction denote the set of formulas  $\phi$  of arithmetic with one free variable so that  $\phi(n)$  is an RA-decidable sentence for each  $n$ . We say that  $M$  is **speculative** if the following holds. Let  $\phi \in \mathcal{F}_0$ , set

$$\alpha_\phi = \forall m : \phi(m),$$

then

$$\forall m : RA \vdash \phi(m) \implies M \vdash \alpha_\phi.$$

Note that of course the left hand side is not the same as  $M \vdash \alpha$ .

We may informally interpret this condition as saying that  $M$  initially prints  $\alpha$  as a hypothesis, and removes  $\alpha$  from its list (that is  $\alpha$  will not be in  $M^s$ ) only if for some  $m$ ,  $M \vdash \neg\phi(m)$ . If each  $\phi(m)$  is  $M^s$ -decidable, this is not in principle conflictory with stable consistency. Indeed in the example above we construct a stably sound Turing machine, with an analogue of this speculative property, deciding the halting problem. Moreover, we have the following crucial result.

**Theorem 3.12.** Given any encoded map

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

so that  $M^s(\mathbb{N}) \supset RA$  there is a speculative encoded map

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

with the properties: if  $M$  is (stably) computable so is  $M_{spec}$ ,  $M_{spec}^s(\mathbb{N}) \supset M^s(\mathbb{N})$ , if  $M$  is stably 1-consistent then  $M_{spec}$  is stably consistent.

Thus the speculative condition is not as forcing as it may sound in the context of stable consistency.

*Proof.*  $\mathcal{F}_0, \mathcal{A}$  are assumed to be encoded so that the maps

$$ev_m : \mathcal{F}_0 \rightarrow \mathcal{A}, \quad \phi \mapsto \phi(m)$$

are computable.

**Lemma 3.13.** There is a map  $F : \mathbb{N} \rightarrow \mathcal{F}_0 \times \{\pm\}$  with the property:

$$F^s(\mathbb{N}) = G := \{\phi \in \mathcal{F}_0 \mid \forall m : RA \vdash \phi(m)\}.$$

*Proof.* The construction is analogous to the construction in the Example 3.3 above. Fix any total, bijective, computable encoded map

$$Z : \mathbb{N} \rightarrow \mathcal{F}.$$

For a  $\phi \in \mathcal{F}_0$  we will say that it is *n-decided* if

$$\forall m \in \{0, \dots, n\} : RA \vdash \phi(m).$$

In what follows each  $L_n$  has the type of an ordered finite list of elements of  $\mathcal{F} \times \{\pm\}$ .

- Initialize  $L_0 := \emptyset$ ,  $n := 0$ .
- Start. For each  $\phi \in \{Z(0), \dots, Z(n)\}$  if it is *n-decided* then add  $(\phi, +)$  to the list  $L_n$ , otherwise add  $(\phi, -)$  to  $L_n$ . Explicitly, we set

$$L_{n+1} := L_n \cup \bigcup_{\phi \in \{Z(0), \dots, Z(n)\}} (\phi, d^n(\phi)),$$

where  $d^n(\phi) = +$  if  $\phi$  is *n-decided* and  $d^n(\phi) = -$  otherwise.

Here  $\cup$  is set union operation of subsets of  $\mathcal{F} \times \{\pm\}$ .

- Set  $n := n + 1$  go to Start and continue.

Set  $L := \cup_n L_n$ , which is an ordered infinite list. Define

$$F : \mathbb{N} \rightarrow \mathcal{F}_0 \times \{\pm\}$$

by

$$F(m) := L(m),$$

with the latter the  $m$ 'th element of  $L$ . □

Define:

$$M_{spec}(n) := \begin{cases} M(k) & \text{if } n = 2k + 1 \\ (\alpha_{pr_{\mathcal{F}}F(k)}, pr_{\pm}F(k)) & \text{if } n = 2k, \end{cases}$$

where  $pr_{\mathcal{F}} : \mathcal{F} \times \{\pm\} \rightarrow \mathcal{F}$ , and  $pr_{\pm} : \mathcal{F} \times \{\pm\} \rightarrow \{\pm\}$  are the natural projections. Then  $M_{spec}$  is speculative and (stably) computable if  $M$  is. Set

$$S := M_{spec}^s(\mathbb{N}).$$

Then  $S$  is consistent unless for some  $\phi \in G$

$$M \vdash \neg \forall m : \phi(m),$$

that is

$$M \vdash \exists m : \neg \phi(m).$$

So that by 1-consistency of  $M^s$

$$\exists m : M \vdash \neg \phi(m).$$

But  $\phi$  is in  $G$ , and  $M^s \supset RA$  so that  $M^s$  is inconsistent, a contradiction, so  $S$  is consistent. □

#### 4. SEMANTIC INCOMPLETENESS FOR STABLY CONSISTENT TURING MACHINES

We start with the simpler semantic case, but we will also use most of the concepts here in the syntactic version of the main theorem.

Let  $\mathcal{D}$  denote the set of encoded maps of the form:

$$D : \mathcal{T} \times \mathbb{N} \rightarrow \{\pm\}.$$

And set

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

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 2.2. Note that  $\text{image } i|_{\mathcal{D}^t}$  is definable, meaning that

$$e_{\mathcal{T}}(\text{image } i|_{\mathcal{D}^t}) \subset \text{Strings} \simeq \mathbb{N}$$

is a set definable by a first order formula in arithmetic. So that in particular the sentence:

$$T \in \mathcal{D}^t$$

is logically equivalent to a first order sentence in arithmetic.

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  $T$  is not  $i^{-1}(T)$ -decided. Note that we can also say  $\neg(T \text{ is } T\text{-decided})$ , as this can be clearer in formulas.

**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:

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

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 informal interpretation of the above is that each  $D \in \mathcal{D}$  is understood as an operation with the properties:

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

In what follows for  $T \in \mathcal{T}$ ,  $D \in \mathcal{D}$ ,  $\Theta_{D,T}$  is shorthand for the sentence:

$$(T \in \mathcal{D}^t) \wedge (\Theta_{D,T}),$$

with  $T \in \mathcal{D}^t$  interpreted as above.

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

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

*Proof.* If

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

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

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

*Proof.* Suppose that  $D \in \mathcal{D}$  stably soundly decides  $\mathcal{P}$  then by the above lemma we obtain:

$$(4.4) \quad \forall T \in \mathcal{D}^t : \Theta_{D,T} \implies \neg(T \text{ is } D\text{-decided}).$$

But it is immediate:

$$(4.5) \quad \forall T \in \mathcal{D}^t : \Theta_{D,T} \implies (\neg(T \text{ is } D\text{-decided}) \implies \neg(T \text{ is } T\text{-decided})).$$

So combining (4.4), (4.5) above we obtain

$$\forall T \in \mathcal{D}^t : \Theta_{D,T} \implies \neg(T \text{ is } T\text{-decided}).$$

But  $D$  stably soundly decides  $\mathcal{P}$  so we conclude:

$$\forall T \in \mathcal{D}^t : \Theta_{D,T} \implies (T \text{ is } D\text{-decided}).$$

But this is a contradiction to (4.4) unless

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

which is what we wanted to prove. □

We can strengthen the result as follows.

**Definition 4.6.** 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 4.7.** 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,

$$\mathcal{R}(D) \implies (\exists T' \in \mathcal{D}^t : (T' \simeq_s T) \wedge (T' \text{ is } D\text{-decided})).$$

Combining with Lemma 4.2 we get:

$$\mathcal{R}(D) \implies \exists T' \in \mathcal{D}^t : (T' \simeq_s T) \wedge \neg\Theta_{D,T'},$$

so that if  $\mathcal{R}(D)$  then we obtain a contradiction since:

$$(T' \simeq_s T) \wedge \neg\Theta_{D,T'} \implies \neg\Theta_{D,T}.$$

□

*Proof of Theorem 1.2.* Suppose that we have such an  $M$ , so  $M$  is stably computable and is stably sound. Let  $CM$  be as in the Lemma 3.9 so that in particular it is stably computable. As already discussed for each  $T \in \mathcal{T}$  the sentence

$$T \in \mathcal{D}^t$$

is logically equivalent to a first order sentence in arithmetic, likewise the sentence:

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

is logically equivalent to a first order sentence in arithmetic we call  $s(T)$ , and the translation is computable. Indeed this kind of translation already appears in the original work of Turing [1].

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

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

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

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

for  $D_{CM}$  defined as in Section 3. Then  $\tilde{D}_M$  is stably sound. By Lemma 3.6  $D_{CM}$  is stably computable and so  $\tilde{D}_M$  is stably computed by some abstract Turing machine  $(Z, Z_e)$ . So by Lemma 4.2

$$\neg(Z_e \text{ is } \tilde{D}_M\text{-decided}),$$

so that  $M \not\models s(Z_e)$  by construction of  $\tilde{D}_M$ . On the other hand  $\neg(Z_e \text{ is } \tilde{D}_M\text{-decided})$ , since  $Z_e \in \mathcal{D}^t$  by construction, entails

$$Z_e \in \mathcal{D}^t \wedge Z_e \text{ is not } Z_e\text{-decided}$$

that is  $s(Z_e)$ . Set  $\alpha(M) := s(Z_e)$  and we are done. □



## 5. SYNTACTIC INCOMPLETENESS FOR STABLY CONSISTENT TURING MACHINES

For each  $T \in \mathcal{T}$  let  $s(T) \in \mathcal{A}$  be the sentence logically equivalent to:

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

where this is interpreted as in the previous section.

Let  $M : \mathbb{N} \rightarrow \mathcal{A} \times \{\pm\}$ , be such that  $M$  is stably stably computable and  $M^s \supset RA$ . Let  $N = M_{spec}$  be as Theorem 3.12 and let  $CN$  be as in the Lemma 3.9. Denote by

$$\tilde{D} : \mathcal{T} \times \mathbb{N} \rightarrow \{\pm\}$$

the decision map defined by:

$$\tilde{D}(T, n) := D_{CN}(s(T), n),$$

where the latter is as in the previous section. Then by Lemma 3.9 and by Lemma 3.6  $\tilde{D}$  is stably computed by some Turing machine  $(Z, Z_e)$ .

**Proposition 5.1.** *For  $M, Z_e$  as above, set  $s(Z) := s(Z_e)$ , then:*

$$M \text{ is stably 1-consistent} \implies M \not\models s(Z).$$

$$M \text{ is stably 2-consistent:} \implies M \not\models \neg s(Z).$$

Moreover, if  $M$  is stably computed by  $(G, G_e)$  then the sentence:

$$G \text{ is stably 1-consistent} \implies s(Z)$$

is a theorem of  $PA$  under standard interpretation of all terms, (we will in part formalize this in the proof).

*Proof.* The first half of the following argument is meant to be a proof of a theorem in  $ZF$ , as this where all terms naturally fit, and we wish to avoid meta-logic. This is also the case elsewhere in the paper implicitly. Arithmetic is interpreted in set theory the standard way, using the standard set  $\mathbb{N}$  of natural numbers. In other words  $RA, PA$ , are understood as theorems about the standard  $\mathbb{N}$ , and we may write  $RA \subset PA \subset ZF$ . In particular  $s(Z)$  as a sentence of arithmetic is also a sentence of set theory under our interpretation.

Let  $M, N$  be as above. So  $N$  is speculative and stably consistent by Theorem 3.12. Suppose by contradiction that  $M$  is stably 1-consistent and  $M \vdash s(Z)$ , then  $N \vdash s(Z)$  and so by the construction,

$$Z_e \text{ is } Z_e\text{-decided}$$

or more explicitly we deduce the sentence  $\eta_Z$ :

$$\exists m \forall n \geq m : Z(Z_e, m) = +.$$

In other words:

$$M \vdash s(Z) \implies \eta_Z$$

is a theorem of  $ZF$ . If we translate  $\eta_Z$  to the arithmetic sentence  $s(\eta_Z)$  then

$$RA \vdash s(\eta_Z) \implies \neg s(Z),$$

by inference.

Now the sentence  $s(\eta_Z)$  is logically equivalent to a sentence in arithmetic of the form:

$$\exists m \forall n : \gamma(m, n),$$

where  $\gamma(m, n)$  is  $RA$ -decidable. So  $ZF$  proves:

$$s(\eta_Z) \implies \exists m \forall n : RA \vdash \gamma(m, n),$$

or

$$s(\eta_Z) \implies \exists m : N \vdash \forall n : \gamma(m, n),$$

since  $N$  is speculative and  $N^s \supset RA$ . So that

$$s(\eta_Z) \implies N \vdash s(\eta_Z),$$

and so

$$s(\eta_Z) \implies N \vdash \neg s(Z).$$

And so

$$M \vdash s(Z) \implies N^s \text{ is not consistent,}$$

which is a contradiction and hence  $ZF$  proves:

$$M \text{ is stably 1-consistent} \implies M \not\vdash s(Z) \quad (\text{moreover } N \not\vdash s(Z)).$$

Now suppose

$$M \text{ is stably 2-consistent} \wedge M \vdash \neg s(Z),$$

so either  $M \vdash s(\eta_Z)$  or  $M \vdash s(\rho_Z)$ , where  $\rho_Z$  is the sentence

$$Z_e \notin \mathcal{D}^t$$

and  $s(\rho_Z)$  the corresponding sentence in arithmetic which is logically equivalent to a sentence of the form:

$$\exists n : \alpha(n),$$

with  $\alpha$   $RA$ -decidable. To elaborate,  $Z_e \notin \mathcal{D}^t$  means that either  $Z_e$  permits some input it shouldn't, does not permit some input it should, or halts with value not in  $\{\pm\}_e$ . Because the set of permissible inputs of  $Z_e$  is computable (decidable) by assumption, the above can be summarized by the condition that for some input the Turing machine  $Z_e$  halts with a particular kind of total state. Which can be expressed by the predicate  $\exists n : \alpha(n)$  by standard elementary arguments. (Consider the zig-zag computation sequence.)

So in case of the latter by 1-consistency of  $M^s$  we would also have

$$\exists n : RA \vdash \alpha(n),$$

So

$$RA \vdash s(\rho_Z).$$

So syntactically

$$ZF \vdash (M \vdash s(\eta_Z) \implies s(\rho_Z)),$$

but also  $ZF \vdash \neg s(\rho_Z)$  by construction of  $Z$ . So that

$$ZF \vdash M \not\vdash s(\eta_Z).$$

Or if we do not care to avoid meta-logic then we may also argue semantically.  $RA$  is sound and  $\neg s(\rho_Z)$  is satisfied by construction of  $Z$  so that  $s(\rho_Z)$  is not satisfied and so  $RA \not\vdash s(\rho_Z)$ , and so  $M \not\vdash s(\eta_Z)$ .

Thus we must have  $M \vdash s(\eta_Z)$ , or

$$M \vdash \exists m \forall n : \gamma(m, n),$$

and by 2-consistency

$$\neg(\forall m : M \vdash \neg \phi(m)),$$

where

$$\phi(m) = \forall n : \gamma(m, n).$$

So we deduce

$$\exists m : M \vdash \phi(m).$$

Now, since  $\gamma$  is  $RA$ -decidable, since  $M^s \supset RA$ , and since  $M^s$  is 1-consistent:

$$\exists m : M \vdash \phi(m) \implies \exists m \forall n : RA \vdash \gamma(m, n).$$

In other words since  $ZF \supset RA$ ,  $ZF$  proves:

$$(M \text{ is stably 2-consistent} \wedge M \vdash \neg s(Z)) \implies s(\eta_Z).$$

And  $ZF$  proves:

$$s(\eta_Z) \implies N \vdash s(Z),$$

by construction of  $Z$ . So we obtain:

$$(M \text{ is stably 2-consistent} \wedge M \vdash \neg s(Z)) \implies N \vdash s(Z),$$

so that we obtain that  $N^s$  is inconsistent, which is again a contradiction and so

$$M \text{ is stably 2-consistent} \implies M \not\models \neg s(Z).$$

Now for the last part of the theorem. Suppose that  $M$  is stably computed by a Turing machine  $(G, G_e)$ , let  $(H, H_e)$  likewise be a Turing machine stably computing  $N$ . Then all subsequent terms and predicates involving  $M, N$  can be interpreted in Peano arithmetic in terms of  $G_e, H_e$  (via Gödel coding or in our language the encoding map  $e_{\mathcal{T}}$ ). These interpretations are standard. Let  $\mathcal{M}_{\mathcal{A}}$  denote the set of encoded maps  $M : \mathbb{N} \rightarrow \mathcal{A} \times \{\pm\}$ , and set

$$\mathcal{T}_{\mathcal{A}} = \{T \in \mathcal{T} \mid \exists M \in \mathcal{M}_{\mathcal{A}} : T \text{ computes } M_e\}.$$

For  $T \in \mathcal{T}$  the sentence  $T \in \mathcal{T}_{\mathcal{A}}$  is logically equivalent to an arithmetic sentence we call  $\zeta(T)$ . For  $T_e \in \mathcal{T}_{\mathcal{A}}$  we say that it is stably 1-consistent if the corresponding  $T : \mathbb{N} \rightarrow \mathcal{A} \times \{\pm\}$  is stably 1-consistent. This is then logically equivalent to an arithmetic sentence we denote by  $1 - \text{con}^s(T_e)$ . Finally  $N \not\models s(Z)$  is likewise logically equivalent to an arithmetic sentence we call  $\omega$ . By the first part of the above argument,

$$\zeta(G_e) \wedge 1 - \text{con}^s(G_e) \implies \omega,$$

which can now be interpreted in  $PA$ , is a theorem of  $ZF$ . But it is also readily seen to be a theorem of  $PA$ . Clearly,

$$PA \vdash \zeta(G_e) \implies \neg s(\rho_{Z_e}).$$

And

$$PA \vdash \neg s(\rho_{Z_e}) \wedge \omega \iff s(Z),$$

by inference, and defining properties of  $(Z, Z_e)$ . So that

$$PA \vdash \zeta(G_e) \wedge 1 - \text{con}^s(G_e) \implies s(Z).$$

□

*Proof of Theorem 1.4.* Suppose that we have a stably computable

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

s.t.  $M^s \supset RA$ . For a Turing machine  $G$  stably computing  $M$ , there is a constructively from  $G$  determined Turing machine  $Z$ , which as a decision map

$$Z : \mathcal{T} \times \mathbb{N} \rightarrow \{\pm\}$$

is defined by:

$$Z(T, n) := D_{CN}(s(T), n),$$

as in the proposition above. Set  $\alpha(M) := s(Z)$ . Then the theorem follows by the proposition 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. KREMNIER 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).
- [17] S. SALEHI AND P. SERAJI, *Gödel-rosser’s incompleteness theorems for non-recursively enumerable theories*, Journal of Logic and Computation, 27-5 (2017).
- [18] R. I. SOARE, *Turing computability. Theory and applications*, Berlin: Springer, 2016.

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