

# Formalising Computability Theory in Isabelle/HOL

Jian Xu, Xingyuan Zhang  
PLA University of Science and Technology Nanjing, China

Christian Urban  
King's College London, UK

**Abstract**—We present a formalised theory of computability in the theorem prover Isabelle/HOL. This theorem prover is based on classical logic which precludes *direct* reasoning about computability: every boolean predicate is either true or false because of the law of excluded middle. The only way to reason about computability in a classical theorem prover is to formalise a concrete model for computation. We formalise Turing machines and relate them to abacus machines and recursive functions. Our theory can be used to formalise other computability results: we give one example about the undecidability of Wang’s tiling problem, whose proof uses the notion of a universal Turing machine.

**Keywords**—Turing Machines, Computability, Isabelle/HOL, Wang tilings

## I. INTRODUCTION

We formalised in earlier work the correctness proofs for two algorithms in Isabelle/HOL—one about type-checking in LF [5] and another about deciding requests in access control [7]. The formalisations uncovered a gap in the informal correctness proof of the former and made us realise that important details were left out in the informal model for the latter. However, in both cases we were unable to formalise in Isabelle/HOL computability arguments about the algorithms. The reason is that both algorithms are formulated in terms of inductive predicates. Suppose  $P$  stands for one such predicate. Decidability of  $P$  usually amounts to showing whether  $P \vee \neg P$  holds. But this does *not* work in Isabelle/HOL, since it is a theorem prover based on classical logic where the law of excluded middle ensures that  $P \vee \neg P$  is always provable no matter whether  $P$  is constructed by computable means. The same problem would arise if we had formulated the algorithms as recursive functions, because internally in Isabelle/HOL, like in all HOL-based theorem provers, functions are represented as inductively defined predicates too.

The only satisfying way out of this problem in a theorem prover based on classical logic is to formalise a theory of computability. Norrish provided such a formalisation for the HOL4 theorem prover. He chose the  $\lambda$ -calculus as the starting point for his formalisation of computability theory, because of its “simplicity” [3, Page 297]. Part of his formalisation is a clever infrastructure for reducing  $\lambda$ -terms. He also established the computational equivalence between the  $\lambda$ -calculus and recursive functions. Nevertheless he concluded that it would be “appealing” to have formalisations for

more operational models of computations, such as Turing machines or register machines. One reason is that many proofs in the literature use them. He noted however that in the context of theorem provers [3, Page 310]:

*“If register machines are unappealing because of their general fiddliness, Turing machines are an even more daunting prospect.”*

In this paper we take on this daunting prospect and provide a formalisation of Turing machines, as well as abacus machines (a kind of register machines) and recursive functions. To see the difficulties involved with this work, one has to understand that interactive theorem provers, like Isabelle/HOL, are at their best when the data-structures at hand are “structurally” defined, like lists, natural numbers, regular expressions, etc. Such data-structures come with convenient reasoning infrastructures (for example induction principles, recursion combinators and so on). But this is *not* the case with Turing machines (and also not with register machines): underlying their definitions are sets of states together with transition functions, all of which are not structurally defined. This means we have to implement our own reasoning infrastructure in order to prove properties about them. This leads to annoyingly fiddly formalisations. We noticed first the difference between both, structural and non-structural, “worlds” when formalising the Myhill-Nerode theorem, where regular expressions fared much better than automata [6]. However, with Turing machines there seems to be no alternative if one wants to formalise the great many proofs from the literature that use them. We will analyse one example—undecidability of Wang’s tiling problem—in Section V. The standard proof of this property uses the notion of *universal Turing machines*.

We are not the first who formalised Turing machines in a theorem prover: we are aware of the preliminary work by Asperti and Ricciotti [1]. They describe a complete formalisation of Turing machines in the Matita theorem prover, including a universal Turing machine. They report that the informal proofs from which they started are *not* “sufficiently accurate to be directly useable as a guideline for formalization” [1, Page 2]. For our formalisation we followed mainly the proofs from the textbook [2] and found that the description there is quite detailed. Some details are left out however: for example, it is only shown how the universal Turing machine is constructed for Turing

machines computing unary functions. We had to figure out a way to generalize this result to  $n$ -ary functions. Similarly, when compiling recursive functions to abacus machines, the textbook again only shows how it can be done for 2- and 3-ary functions, but in the formalisation we need arbitrary functions. But the general ideas for how to do this are clear enough in [2]. However, one aspect that is completely left out from the informal description in [2], and similar ones we are aware of, is arguments why certain Turing machines are correct. We will introduce Hoare-style proof rules which help us with such correctness arguments of Turing machines.

The main difference between our formalisation and the one by Asperti and Ricciotti is that their universal Turing machine uses a different alphabet than the machines it simulates. They write [1, Page 23]:

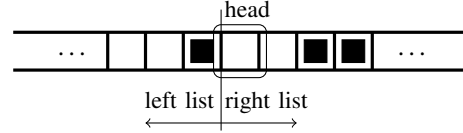
*“In particular, the fact that the universal machine operates with a different alphabet with respect to the machines it simulates is annoying.”*

In this paper we follow the approach by Boolos et al [2], which goes back to Post [4], where all Turing machines operate on tapes that contain only *blank* or *occupied* cells (represented by  $Bk$  and  $Oc$ , respectively, in our formalisation). Traditionally the content of a cell can be any character from a finite alphabet. Although computationally equivalent, the more restrictive notion of Turing machines in [2] makes the reasoning more uniform. In addition some proofs *about* Turing machines are simpler. The reason is that one often needs to encode Turing machines—consequently if the Turing machines are simpler, then the coding functions are simpler too. Unfortunately, the restrictiveness also makes it harder to design programs for these Turing machines. In order to construct a *universal Turing machine* we therefore do not follow [1], instead follow the proof in [2] by relating abacus machines to Turing machines and in turn recursive functions to abacus machines. The universal Turing machine can then be constructed as a recursive function.

## Contributions:

## II. TURING MACHINES

Turing machines can be thought of as having a read-write-unit, also referred to as *head*, “gliding” over a potentially infinite tape. Boolos et al [2] only consider tapes with cells being either blank or occupied, which we represent by a datatype having two constructors, namely  $Bk$  and  $Oc$ . One way to represent such tapes is to use a pair of lists, written  $(l, r)$ , where  $l$  stands for the tape on the left-hand side of the head and  $r$  for the tape on the right-hand side. We have the convention that the head, abbreviated  $hd$ , of the right-list is the cell on which the head of the Turing machine currently operates. This can be pictured as follows:



Note that by using lists each side of the tape is only finite. The potential infinity is achieved by adding an appropriate blank or occupied cell whenever the head goes over the “edge” of the tape. To make this formal we define five possible *actions* the Turing machine can perform:

$a$	::=	$W_{Bk}$	write blank ( $Bk$ )
		$W_{Oc}$	write occupied ( $Oc$ )
		$L$	move left
		$R$	move right
		$Nop$	do-nothing operation

We slightly deviate from the presentation in [2] by using the  $Nop$  operation; however its use will become important when we formalise halting computations and also universal Turing machines. Given a tape and an action, we can define the following tape updating function:

$$\begin{aligned}
 \text{update } (l, r) \ W_{Bk} &\stackrel{\text{def}}{=} (l, Bk::tl\ r) \\
 \text{update } (l, r) \ W_{Oc} &\stackrel{\text{def}}{=} (l, Oc::tl\ r) \\
 \text{update } (l, r) \ L &\stackrel{\text{def}}{=} \\
 &\quad \text{if } l = [] \text{ then } ([], Bk::r) \text{ else } (tl\ l, hd\ l::r) \\
 \text{update } (l, r) \ R &\stackrel{\text{def}}{=} \\
 &\quad \text{if } r = [] \text{ then } (Bk::l, []) \text{ else } (hd\ r::l, tl\ r) \\
 \text{update } (l, r) \ Nop &\stackrel{\text{def}}{=} (l, r)
 \end{aligned}$$

The first two clauses replace the head of the right-list with a new  $Bk$  or  $Oc$ , respectively. To see that these two clauses make sense in case where  $r$  is the empty list, one has to know that the tail function,  $tl$ , is defined in Isabelle/HOL such that  $tl\ [] \stackrel{\text{def}}{=} []$  holds. The third clause implements the move of the head one step to the left: we need to test if the left-list  $l$  is empty; if yes, then we just prepend a blank cell to the right-list; otherwise we have to remove the head from the left-list and prepend it to the right-list. Similarly in the fourth clause for a right move action. The  $Nop$  operation leaves the the tape unchanged (last clause).

Note that our treatment of the tape is rather “unsymmetric”—we have the convention that the head of the right-list is where the head is currently positioned. Asperti and Ricciotti [1] also considered such a representation, but dismiss it as it complicates their definition for *tape equality*. The reason is that moving the head one step to the left and then back to the right might change the tape (in case of going over the “edge”). Therefore they distinguish four types of tapes: one where the tape is empty; another where the head is on the left edge, respectively right edge, and in the middle of the tape. The reading, writing and moving of the tape is then defined in terms of these four cases. In this way they can keep the tape in a “normalised” form, and

thus making a left-move followed by a right-move being the identity on tapes. Since we are not using the notion of tape equality, we can get away with the unsymmetric definition above, and by using the *update* function cover uniformly all cases including corner cases.

Next we need to define the *states* of a Turing machine. Given how little is usually said about how to represent them in informal presentations, it might be surprising that in a theorem prover we have to select carefully a representation. If we use the naive representation where a Turing machine consists of a finite set of states, then we will have difficulties composing two Turing machines: we would need to combine two finite sets of states, possibly renaming states apart whenever both machines share states.<sup>1</sup> This renaming can be quite cumbersome to reason about. Therefore we made the choice of representing a state by a natural number and the states of a Turing machine will always consist of the initial segment of natural numbers starting from 0 up to the number of states of the machine. In doing so we can compose two Turing machine by shifting the states of one by an appropriate amount to a higher segment and adjusting some “next states” in the other.

An *instruction*  $i$  of a Turing machine is a pair consisting of an action and a natural number (the next state). A *program*  $p$  of a Turing machine is then a list of such pairs. Using as an example the following Turing machine program, which consists of four instructions

$$dither \stackrel{def}{=} \underbrace{[(W_{Bk}, 1), (R, 2)]}_{1st\ state} \underbrace{[(L, 1), (L, 0)]}_{2nd\ state} \quad (1)$$

$\overbrace{W_{Bk}, 1}^{Bk\text{-case}}$      $\overbrace{R, 2}^{Oc\text{-case}}$   
 $\overbrace{L, 1, L, 0}$

the reader can see we have organised our Turing machine programs so that segments of two belong to a state. The first component of the segment determines what action should be taken and which next state should be transitioned to in case the head reads a *Bk*; similarly the second component determines what should be done in case of reading *Oc*. We have the convention that the first state is always the *starting state* of the Turing machine. The zeroth state is special in that it will be used as the “halting state”. There are no instructions for the 0-state, but it will always perform a *Nop*-operation and remain in the 0-state. Unlike Asperti and Ricciotti [1], we have chosen a very concrete representation for programs, because when constructing a universal Turing machine, we need to define a coding function for programs. This can be easily done for our programs-as-lists, but is more difficult for the functions used by Asperti and Ricciotti.

Given a program  $p$ , a state and the cell being read by the head, we need to fetch the corresponding instruction from

the program. For this we define the function *fetch*

$$\begin{aligned} fetch\ p\ 0\ _ &\stackrel{def}{=} (Nop, 0) \\ fetch\ p\ (Suc\ s)\ Bk &\stackrel{def}{=} \\ &\quad case\ nth\_of\ p\ (2 * s)\ of \\ &\quad \quad None \Rightarrow (Nop, 0) \mid \\ &\quad \quad Some\ i \Rightarrow i \\ fetch\ p\ (Suc\ s)\ Oc &\stackrel{def}{=} \\ &\quad case\ nth\_of\ p\ (2 * s + 1)\ of \\ &\quad \quad None \Rightarrow (Nop, 0) \mid \\ &\quad \quad Some\ i \Rightarrow i \end{aligned}$$

In this definition the function *nth\_of* returns the  $n$ th element from a list, provided it exists (*Some*-case), or if it does not, it returns the default action *Nop* and the default state 0 (*None*-case). In doing so we slightly deviate from the description in [2]: if their Turing machines transition to a non-existing state, then the computation is halted. We will transition in such cases to the 0-state. However, with introducing the notion of *well-formed* Turing machine programs we will later exclude such cases and make the 0-state the only “halting state”. A program  $p$  is said to be well-formed if it satisfies the following three properties:

$$\begin{aligned} twf\ p &\stackrel{def}{=} 2 \leq length\ p \\ &\quad \wedge\ iseven\ (length\ p) \\ &\quad \wedge\ \forall (a, s) \in p. s \leq length\ p\ div\ 2 \end{aligned}$$

The first says that  $p$  must have at least an instruction for the starting state; the second that  $p$  has a *Bk* and *Oc* instruction for every state, and the third that every next-state is one of the states mentioned in the program or being the 0-state.

A *configuration*  $c$  of a Turing machine is a state together with a tape. This is written as  $(s, (l, r))$ . If we have a configuration and a program, we can calculate what the next configuration is by fetching the appropriate action and next state from the program, and by updating the state and tape accordingly. This single step of execution is defined as the function *step*

$$\begin{aligned} step\ (s, (l, r))\ p &\stackrel{def}{=} \\ &\quad let\ (a, s) = fetch\ p\ s\ (read\ r) \\ &\quad in\ (s', update\ (l, r)\ a) \end{aligned}$$

where *read*  $r$  returns the head of the list  $r$ , or if  $r$  is empty it returns *Bk*. It is impossible in Isabelle/HOL to lift the *step*-function realising a general evaluation function for Turing machines. The reason is that functions in HOL-based provers need to be terminating, and clearly there are Turing machine programs that are not. We can however define an evaluation function so that it performs exactly  $n$  steps:

$$\begin{aligned} nsteps\ c\ p\ 0 &\stackrel{def}{=} c \\ nsteps\ c\ p\ (Suc\ n) &\stackrel{def}{=} nsteps\ (step\ c\ p)\ p\ n \end{aligned}$$

<sup>1</sup>The usual disjoint union operation in Isabelle/HOL cannot be used as it does not preserve types.

Recall our definition of *fetch* with the default value for the  $0$ -state. In case a Turing program takes in [2] less than  $n$  steps before it halts, then in our setting the *nsteps*-evaluation does not actually halt, but rather transitions to the  $0$ -state and remains there performing *Nop*-actions until  $n$  is reached.

Given some input tape  $(l_i, r_i)$ , we can define when a program  $p$  generates a specific output tape  $(l_o, r_o)$

$$\text{runs } p(l_i, r_i)(l_o, r_o) \stackrel{\text{def}}{=} \exists n. \text{nsteps}(l, (l_i, r_i)) p n = (0, (l_o, r_o))$$

where  $l$  stands for the starting state and  $0$  for our final state. A program  $p$  with input tape  $(l_i, r_i)$  halts iff

$$\text{halts } p(l_i, r_i) \stackrel{\text{def}}{=} \exists l_o r_o. \text{runs } p(l_i, r_i)(l_o, r_o)$$

Later on we need to consider specific Turing machines that start with a tape in standard form and halt the computation in standard form. To define a tape in standard form, it is useful to have an operation  $\ulcorner \_ \urcorner$  that translates lists of natural numbers into tapes.

$$\begin{aligned} \ulcorner \_ \urcorner &\stackrel{\text{def}}{=} \_ \\ \ulcorner [n] \urcorner &\stackrel{\text{def}}{=} Oc^n + l \\ \ulcorner n::ns \urcorner &\stackrel{\text{def}}{=} Oc^n + l @ [Bk] @ \ulcorner ns \urcorner \end{aligned}$$

By this we mean

$$\text{stdhalt } p n \stackrel{\text{def}}{=} \exists k l m. \text{run } p(\ulcorner \_ \urcorner, Oc^n)(Bk^k, Oc^l @ Bk^m)$$

This means the Turing machine starts with a tape containing  $n$  *Ocs* and the head pointing to the first one; the Turing machine halts with a tape consisting of some *Bks*, followed by a “cluster” of *Ocs* and after that by some *Bks*. The head in the output is pointing again at the first *Oc*. The intuitive meaning of this definition is to start the Turing machine with a tape corresponding to a value  $n$  and producing a new tape corresponding to the value  $l$  (the number of *Ocs* clustered on the output tape).

Before we can prove the undecidability of the halting problem for Turing machines, we have to define how to compose Turing machines. Given our setup, this is relatively straightforward, if slightly fiddly.

$$\text{abacus.tshift } p n \stackrel{\text{def}}{=} \text{map } (\lambda(a, s). (a, \text{if } s = 0 \text{ then } 0 \text{ else } s + n)) p$$

(\* HERE \*)

shift and change of a  $p$   
composition of two  $ps$   
assertion holds for all tapes  
Hoare rule for composition

For showing the undecidability of the halting problem, we need to consider two specific Turing machines. copying TM and dithering TM

correctness of the copying TM  
measure for the copying TM, which we however omit.

halting problem

### III. ABACUS MACHINES

Boolos et al [2] use abacus machines as a stepping stone for making it less laborious to write programs for Turing machines. Abacus machines operate over an unlimited number of registers  $R_0, R_1, \dots$  each being able to hold an arbitrary large natural number. We use natural numbers to refer to registers, but also to refer to *opcodes* of abacus machines. Opcodes are given by the datatype

$$\begin{aligned} o ::= & \text{Inc } R && \text{increment register } R \text{ by one} \\ & | \text{Dec } R o && \text{if content of } R \text{ is non-zero,} \\ & && \text{then decrement it by one} \\ & && \text{otherwise jump to opcode } o \\ & | \text{Goto } o && \text{jump to opcode } o \end{aligned}$$

A *program* of an abacus machine is a list of such opcodes. For example the program clearing the register  $R$  (setting it to 0) can be defined as follows:

$$\text{clear } R o \stackrel{\text{def}}{=} [\text{Dec } R o, \text{Goto } 0]$$

The second opcode *Goto* 0 in this program means we jump back to the first opcode, namely *Dec*  $R o$ . The *memory*  $m$  of an abacus machine holding the values of the registers is represented as a list of natural numbers. We have a lookup function for this memory, written *lookup*  $m R$ , which looks up the content of register  $R$ ; if  $R$  is not in this list, then we return 0. Similarly we have a setting function, written *set*  $m R n$ , which sets the value of  $R$  to  $n$ , and if  $R$  was not yet in  $m$  it pads it appropriately with 0s.

Abacus machine halts when it jumps out of range.

### IV. RECURSIVE FUNCTIONS

#### V. WANG TILES

Used in texture mappings - graphics

#### VI. RELATED WORK

The most closely related work is by Norrish [3], and Asperti and Ricciotti [1]. Norrish bases his approach on lambda-terms. For this he introduced a clever rewriting technology based on combinators and de-Brujin indices for rewriting modulo  $\beta$ -equivalence (to keep it manageable)

#### REFERENCES

- [1] A. Asperti and W. Ricciotti. Formalizing Turing Machines. In *Proc. of the 19th International Workshop on Logic, Language, Information and Computation (WoLLIC)*, volume 7456 of *LNCS*, pages 1–25, 2012.
- [2] G. Boolos, J. P. Burgess, and R. C. Jeffrey. *Computability and Logic (5th ed.)*. Cambridge University Press, 2007.
- [3] M. Norrish. Mechanised Computability Theory. In *Proc. of the 2nd Conference on Interactive Theorem Proving (ITP)*, volume 6898 of *LNCS*, pages 297–311, 2011.

- [4] E. Post. Finite Combinatory Processes-Formulation 1. *Journal of Symbolic Logic*, 1(3):103–105, 1936.
- [5] C. Urban, J. Cheney, and S. Berghofer. Mechanizing the Metatheory of LF. *ACM Transactions on Computational Logic*, 12:15:1–15:42, 2011.
- [6] C. Wu, X. Zhang, and C. Urban. A Formalisation of the Myhill-Nerode Theorem based on Regular Expressions (Proof Pearl). In *Proc. of the 2nd Conference on Interactive Theorem Proving*, volume 6898 of *LNCS*, pages 341–356, 2011.
- [7] C. Wu, X. Zhang, and C. Urban. ??? Submitted, 2012.