

## Chapter und

# Undecidability

### und.1 Introduction

tur:und:int:  
sec It might seem obvious that not every function, even every arithmetical function, can be computable. There are just too many, whose behavior is too complicated. Functions defined from the decay of radioactive particles, for instance, or other chaotic or random behavior. Suppose we start counting 1-second intervals from a given time, and define the function  $f(n)$  as the number of particles in the universe that decay in the  $n$ -th 1-second interval after that initial moment. This seems like a candidate for a function we cannot ever hope to compute.

But it is one thing to not be able to imagine how one would compute such functions, and quite another to actually prove that they are uncomputable. In fact, even functions that seem hopelessly complicated may, in an abstract sense, be computable. For instance, suppose the universe is finite in time—some day, in the very distant future the universe will contract into a single point, as some cosmological theories predict. Then there is only a finite (but incredibly large) number of seconds from that initial moment for which  $f(n)$  is defined. And any function which is defined for only finitely many inputs is computable: we could list the outputs in one big table, or code it in one very big Turing machine state transition diagram.

We are often interested in special cases of functions whose values give the answers to yes/no questions. For instance, the question “is  $n$  a prime number?” is associated with the function

$$\text{isprime}(n) = \begin{cases} 1 & \text{if } n \text{ is prime} \\ 0 & \text{otherwise.} \end{cases}$$

We say that a yes/no question can be *effectively decided*, if the associated 1/0-valued function is effectively computable.

To prove mathematically that there are functions which cannot be effectively computed, or problems that cannot effectively decided, it is essential to fix a specific model of computation, and show about it that there are functions

it cannot compute or problems it cannot decide. We can show, for instance, that not every function can be computed by Turing machines, and not every problem can be decided by Turing machines. We can then appeal to the Church-Turing thesis to conclude that not only are Turing machines not powerful enough to compute every function, but no effective procedure can.

The key to proving such negative results is the fact that we can assign numbers to Turing machines themselves. The easiest way to do this is to enumerate them, perhaps by fixing a specific way to write down Turing machines and their programs, and then listing them in a systematic fashion. Once we see that this can be done, then the existence of Turing-uncomputable functions follows by simple cardinality considerations: the set of functions from  $\mathbb{N}$  to  $\mathbb{N}$  (in fact, even just from  $\mathbb{N}$  to  $\{0, 1\}$ ) are **non-enumerable**, but since we can enumerate all the Turing machines, the set of Turing-computable functions is only **denumerable**.

We can also define *specific* functions and problems which we can prove to be uncomputable and undecidable, respectively. One such problem is the so-called *Halting Problem*. Turing machines can be finitely described by listing their instructions. Such a description of a Turing machine, i.e., a Turing machine program, can of course be used as input to another Turing machine. So we can consider Turing machines that decide questions about other Turing machines. One particularly interesting question is this: “Does the given Turing machine eventually halt when started on input  $n$ ?” It would be nice if there were a Turing machine that could decide this question: think of it as a quality-control Turing machine which ensures that Turing machines don’t get caught in infinite loops and such. The interesting fact, which Turing proved, is that there cannot be such a Turing machine. There cannot be a single Turing machine which, when started on input consisting of a description of a Turing machine  $M$  and some number  $n$ , will always halt with either output 1 or 0 according to whether  $M$  machine would have halted when started on input  $n$  or not.

Once we have examples of specific undecidable problems we can use them to show that other problems are undecidable, too. For instance, one celebrated undecidable problem is the question, “Is the first-order **formula**  $\varphi$  valid?”. There is no Turing machine which, given as input a first-order **formula**  $\varphi$ , is guaranteed to halt with output 1 or 0 according to whether  $\varphi$  is valid or not. Historically, the question of finding a procedure to effectively solve this problem was called simply “the” decision problem; and so we say that the decision problem is unsolvable. Turing and Church proved this result independently at around the same time, so it is also called the Church-Turing Theorem.

## und.2 Enumerating Turing Machines

explanation

We can show that the set of all Turing-machines is **enumerable**. This follows from the fact that each Turing machine can be finitely described. The set of states and the tape vocabulary are finite sets. The transition function is a partial function from  $Q \times \Sigma$  to  $Q \times \Sigma \times \{L, R, N\}$ , and so likewise can be

tur:und:enu:  
sec

specified by listing its values for the finitely many argument pairs for which it is defined. Of course, strictly speaking, the states and vocabulary can be anything; but the *behavior* of the Turing machine is independent of which objects serve as states and vocabulary. So we may assume, for instance, that the states and vocabulary symbols are natural numbers, or that the states and vocabulary are all strings of letters and digits.

Suppose we fix a **denumerable** vocabulary for specifying Turing machines:  $\sigma_0 = \triangleright, \sigma_1 = 0, \sigma_2 = 1, \sigma_3, \dots, R, L, N, q_0, q_1, \dots$ . Then any Turing machine can be specified by some finite string of symbols from this alphabet (though not every finite string of symbols specifies a Turing machine). For instance, suppose we have a Turing machine  $M = \langle Q, \Sigma, q, \delta \rangle$  where

$$Q = \{q'_0, \dots, q'_n\} \subseteq \{q_0, q_1, \dots\} \text{ and}$$

$$\Sigma = \{\triangleright, \sigma'_1, \sigma'_2, \dots, \sigma'_m\} \subseteq \{\sigma_0, \sigma_1, \dots\}.$$

We could specify it by the string

$$q'_0 q'_1 \dots q'_n \triangleright \sigma'_1 \dots \sigma'_m \triangleright q \triangleright S(\sigma'_0, q'_0) \triangleright \dots \triangleright S(\sigma'_m, q'_n)$$

where  $S(\sigma'_i, q'_j)$  is the string  $\sigma'_i q'_j \delta(\sigma'_i, q'_j)$  if  $\delta(\sigma'_i, q'_j)$  is defined, and  $\sigma'_i q'_j$  otherwise.

**Theorem und.1.** *There are functions from  $\mathbb{N}$  to  $\mathbb{N}$  which are not Turing computable.*

*Proof.* We know that the set of finite strings of symbols from a **denumerable** alphabet is **enumerable**. This gives us that the set of descriptions of Turing machines, as a subset of the finite strings from the **enumerable** vocabulary  $\{q_0, q_1, \dots, \triangleright, \sigma_1, \sigma_2, \dots\}$ , is itself enumerable. Since every Turing computable function is computed by some (in fact, many) Turing machines, this means that the set of all Turing computable functions from  $\mathbb{N}$  to  $\mathbb{N}$  is also enumerable.

On the other hand, the set of all functions from  $\mathbb{N}$  to  $\mathbb{N}$  is not **enumerable**. This follows immediately from the fact that not even the set of all functions of one argument from  $\mathbb{N}$  to the set  $\{0, 1\}$  is **enumerable**. If all functions were computable by some Turing machine we could enumerate the set of all functions. So there are some functions that are not Turing-computable.  $\square$

### und.3 The Halting Problem

tms:und:hal:  
sec

Assume we have fixed some finite descriptions of Turing machines. Using these, we can enumerate Turing machines via their descriptions, say, ordered by the lexicographic ordering. Each Turing machine thus receives an *index*: its place in the enumeration  $M_1, M_2, M_3, \dots$  of Turing machine descriptions.

explanation

We know that there must be non-Turing-computable functions: the set of Turing machine descriptions—and hence the set of Turing machines—is enumerable, but the set of all functions from  $\mathbb{N}$  to  $\mathbb{N}$  is not. But we can find specific examples of non-computable function as well. One such function is the halting function.

**Definition und.2 (Halting function).** The *halting function*  $h$  is defined as

$$h(e, n) = \begin{cases} 0 & \text{if machine } M_e \text{ does not halt for input } n \\ 1 & \text{if machine } M_e \text{ halts for input } n \end{cases}$$

**Definition und.3 (Halting problem).** The *Halting Problem* is the problem of determining (for any  $e, n$ ) whether the Turing machine  $M_e$  halts for an input of  $n$  strokes.

explanation We show that  $h$  is not Turing-computable by showing that a related function,  $s$ , is not Turing-computable. This proof relies on the fact that anything that can be computed by a Turing machine can be computed using just two symbols: 0 and 1, and the fact that two Turing machines can be hooked together to create a single machine.

**Definition und.4.** The function  $s$  is defined as

$$s(e) = \begin{cases} 0 & \text{if machine } M_e \text{ does not halt for input } e \\ 1 & \text{if machine } M_e \text{ halts for input } e \end{cases}$$

**Lemma und.5.** *The function  $s$  is not Turing computable.*

*Proof.* We suppose, for contradiction, that the function  $s$  is Turing-computable. Then there would be a Turing machine  $S$  that computes  $s$ . We may assume, without loss of generality, that when  $S$  halts, it does so while scanning the first square. This machine can be “hooked up” to another machine  $J$ , which halts if it is started on a blank tape (i.e., if it reads 0 in the initial state while scanning the square to the right of the end-of-tape symbol), and otherwise wanders off to the right, never halting.  $S \frown J$ , the machine created by hooking  $S$  to  $J$ , is a Turing machine, so it is  $M_e$  for some  $e$  (i.e., it appears somewhere in the enumeration). Start  $M_e$  on an input of  $e$  1s. There are two possibilities: either  $M_e$  halts or it does not halt.

1. Suppose  $M_e$  halts for an input of  $e$  1s. Then  $s(e) = 1$ . So  $S$ , when started on  $e$ , halts with a single 1 as output on the tape. Then  $J$  starts with a 1 on the tape. In that case  $J$  does not halt. But  $M_e$  is the machine  $S \frown J$ , so it should do exactly what  $S$  followed by  $J$  would do. So  $M_e$  cannot halt for an input of  $e$  1's.
2. Now suppose  $M_e$  does not halt for an input of  $e$  1s. Then  $s(e) = 0$ , and  $S$ , when started on input  $e$ , halts with a blank tape.  $J$ , when started on a blank tape, immediately halts. Again,  $M_e$  does what  $S$  followed by  $J$  would do, so  $M_e$  must halt for an input of  $e$  1's.

This shows there cannot be a Turing machine  $S$ :  $s$  is not Turing computable.  $\square$

**Theorem und.6 (Unsolvability of the Halting Problem).** *The halting problem is unsolvable, i.e., the function  $h$  is not Turing computable.*

tms:und:hal:  
thm:halting-problem

*Proof.* Suppose  $h$  were Turing computable, say, by a Turing machine  $H$ . We could use  $H$  to build a Turing machine that computes  $s$ : First, make a copy of the input (separated by a blank). Then move back to the beginning, and run  $H$ . We can clearly make a machine that does the former, and if  $H$  existed, we would be able to “hook it up” to such a modified doubling machine to get a new machine which would determine if  $M_e$  halts on input  $e$ , i.e., computes  $s$ . But we’ve already shown that no such machine can exist. Hence,  $h$  is also not Turing computable.  $\square$

**Problem und.1.** The Three Halting (3-Halt) problem is the problem of giving a decision procedure to determine whether or not an arbitrarily chosen Turing Machine halts for an input of three strokes on an otherwise blank tape. Prove that the 3-Halt problem is unsolvable.

**Problem und.2.** Show that if the halting problem is solvable for Turing machine and input pairs  $M_e$  and  $n$  where  $e \neq n$ , then it is also solvable for the cases where  $e = n$ .

**Problem und.3.** We proved that the halting problem is unsolvable if the input is a number  $e$ , which identifies a Turing machine  $M_e$  via an enumeration of all Turing machines. What if we allow the description of Turing machines from [section und.2](#) directly as input? (This would require a larger alphabet of course.) Can there be a Turing machine which decides the halting problem but takes as input descriptions of Turing machines rather than indices? Explain why or why not.

## und.4 The Decision Problem

[tms:und:dec:](#)  
[sec](#)

We say that first-order logic is *decidable* iff there is an effective method for determining whether or not a given [sentence](#) is valid. As it turns out, there is no such method: the problem of deciding validity of first-order sentences is unsolvable.

In order to establish this important negative result, we prove that the decision problem cannot be solved by a Turing machine. That is, we show that there is no Turing machine which, whenever it is started on a tape that contains a first-order [sentence](#), eventually halts and outputs either 1 or 0 depending on whether the [sentence](#) is valid or not. By the Church-Turing thesis, every function which is computable is Turing computable. So if this “validity function” were effectively computable at all, it would be Turing computable. If it isn’t Turing computable, then, it also cannot be effectively computable.

Our strategy for proving that the decision problem is unsolvable is to reduce the halting problem to it. This means the following: We have proved that the function  $h(e, w)$  that halts with output 1 if the Turing-machine described by  $e$  halts on input  $w$  and outputs 0 otherwise, is not Turing-computable. We will show that if there were a Turing machine that decides validity of first-order sentences, then there is also Turing machine that computes  $h$ . Since  $h$  cannot

be computed by a Turing machine, there cannot be a Turing machine that decides validity either.

The first step in this strategy is to show that for every input  $w$  and a Turing machine  $M$ , we can effectively describe a sentence  $\tau(M, w)$  representing the instruction set of  $M$  and the input  $w$  and a sentence  $\alpha(M, w)$  expressing “ $M$  eventually halts” such that:

$$\models \tau(M, w) \rightarrow \alpha(M, w) \text{ iff } M \text{ halts for input } w.$$

The bulk of our proof will consist in describing these sentences  $\tau(M, w)$  and  $\alpha(M, w)$  and verifying that  $\tau(M, w) \rightarrow \alpha(M, w)$  is valid iff  $M$  halts on input  $w$ .

## und.5 Representing Turing Machines

explanation

In order to represent Turing machines and their behavior by a sentence of first-order logic, we have to define a suitable language. The language consists of two parts: predicate symbols for describing configurations of the machine, and expressions for numbering execution steps (“moments”) and positions on the tape.

tms:und:rep:  
sec

We introduce two kinds of predicate symbols, both of them 2-place: For each state  $q$ , a predicate symbol  $Q_q$ , and for each tape symbol  $\sigma$ , a predicate symbol  $S_\sigma$ . The former allow us to describe the state of  $M$  and the position of its tape head, the latter allow us to describe the contents of the tape.

In order to express the positions of the tape head and the number of steps executed, we need a way to express numbers. This is done using a constant symbol  $0$ , and a 1-place function  $t$ , the successor function. By convention it is written *after* its argument (and we leave out the parentheses). So  $0$  names the leftmost position on the tape as well as the time before the first execution step (the initial configuration),  $0'$  names the square to the right of the leftmost square, and the time after the first execution step, and so on. We also introduce a predicate symbol  $<$  to express both the ordering of tape positions (when it means “to the left of”) and execution steps (then it means “before”).

Once we have the language in place, we list the “axioms” of  $\tau(M, w)$ , i.e., the sentences which, taken together, describe the behavior of  $M$  when run on input  $w$ . There will be sentences which lay down conditions on  $0$ ,  $t$ , and  $<$ , sentences that describes the input configuration, and sentences that describe what the configuration of  $M$  is after it executes a particular instruction.

**Definition und.7.** Given a Turing machine  $M = \langle Q, \Sigma, q_0, \delta \rangle$ , the language  $\mathcal{L}_M$  consists of:

tms:und:rep:  
defn:tm-descr

1. A two-place predicate symbol  $Q_q(x, y)$  for every state  $q \in Q$ . Intuitively,  $Q_q(\bar{m}, \bar{n})$  expresses “after  $n$  steps,  $M$  is in state  $q$  scanning the  $m$ th square.”
2. A two-place predicate symbol  $S_\sigma(x, y)$  for every symbol  $\sigma \in \Sigma$ . Intuitively,  $S_\sigma(\bar{m}, \bar{n})$  expresses “after  $n$  steps, the  $m$ th square contains symbol  $\sigma$ .”

3. A constant symbol  $o$
4. A one-place function symbol  $'$
5. A two-place predicate symbol  $<$

For each number  $n$  there is a canonical term  $\bar{n}$ , the *numeral* for  $n$ , which represents it in  $\mathcal{L}_M$ .  $\bar{0}$  is  $o$ ,  $\bar{1}$  is  $o'$ ,  $\bar{2}$  is  $o''$ , and so on. More formally:

$$\begin{aligned}\bar{0} &= o \\ \overline{n+1} &= \bar{n}'\end{aligned}$$

The sentences describing the operation of the Turing machine  $M$  on input  $w = \sigma_{i_1} \dots \sigma_{i_k}$  are the following:

1. Axioms describing numbers:

- a) A sentence that says that the successor function is injective:

$$\forall x \forall y (x' = y' \rightarrow x = y)$$

- b) A sentence that says that every number is less than its successor:

$$\forall x x < x'$$

- c) A sentence that ensures that  $<$  is transitive:

$$\forall x \forall y \forall z ((x < y \wedge y < z) \rightarrow x < z)$$

- d) A sentence that connects  $<$  and  $=$ :

$$\forall x \forall y (x < y \rightarrow x \neq y)$$

2. Axioms describing the input configuration:

- a) After 0 steps—before the machine starts— $M$  is in the initial state  $q_0$ , scanning square 1:

$$Q_{q_0}(\bar{1}, \bar{0})$$

- b) The first  $k+1$  squares contain the symbols  $\triangleright, \sigma_{i_1}, \dots, \sigma_{i_k}$ :

$$S_{\triangleright}(\bar{0}, \bar{0}) \wedge S_{\sigma_{i_1}}(\bar{1}, \bar{0}) \wedge \dots \wedge S_{\sigma_{i_k}}(\bar{n}, \bar{0})$$

- c) Otherwise, the tape is empty:

$$\forall x (\bar{k} < x \rightarrow S_0(x, \bar{0}))$$

tms:und:rep:  
tm-rep-a

3. Axioms describing the transition from one configuration to the next:

For the following, let  $\varphi(x, y)$  be the conjunction of all **sentences** of the form

$$\forall z ((z < x \vee x < z) \wedge S_\sigma(z, y)) \rightarrow S_\sigma(z, y')$$

where  $\sigma \in \Sigma$ . We use  $\varphi(\bar{m}, \bar{n})$  to express “other than at square  $m$ , the tape after  $n + 1$  steps is the same as after  $n$  steps.”

- a) For every instruction  $\delta(q_i, \sigma) = \langle q_j, \sigma', R \rangle$ , the **sentence**:

tms:und:rep:  
rep-right

$$\begin{aligned} \forall x \forall y ((Q_{q_i}(x, y) \wedge S_\sigma(x, y)) \rightarrow \\ (Q_{q_j}(x', y') \wedge S_{\sigma'}(x, y') \wedge \varphi(x, y))) \end{aligned}$$

This says that if, after  $y$  steps, the machine is in state  $q_i$  scanning square  $x$  which contains symbol  $\sigma$ , then after  $y+1$  steps it is scanning square  $x+1$ , is in state  $q_j$ , square  $x$  now contains  $\sigma'$ , and every square other than  $x$  contains the same symbol as it did after  $y$  steps.

- b) For every instruction  $\delta(q_i, \sigma) = \langle q_j, \sigma', L \rangle$ , the **sentence**:

tms:und:rep:  
rep-left

$$\begin{aligned} \forall x \forall y ((Q_{q_i}(x', y) \wedge S_\sigma(x', y)) \rightarrow \\ (Q_{q_j}(x, y') \wedge S_{\sigma'}(x', y') \wedge \varphi(x, y))) \wedge \\ \forall y ((Q_{q_i}(0, y) \wedge S_\sigma(0, y)) \rightarrow \\ (Q_{q_j}(0, y') \wedge S_{\sigma'}(0, y') \wedge \varphi(0, y))) \end{aligned}$$

Take a moment to think about how this works: now we don't start with “if scanning square  $x \dots$ ” but: “if scanning square  $x+1 \dots$ ” A move to the left means that in the next step the machine is scanning square  $x$ . But the square that is written on is  $x + 1$ . We do it this way since we don't have subtraction or a predecessor function.

Note that numbers of the form  $x + 1$  are 1, 2,  $\dots$ , i.e., this doesn't cover the case where the machine is scanning square 0 and is supposed to move left (which of course it can't—it just stays put). That special case is covered by the second conjunction: it says that if, after  $y$  steps, the machine is scanning square 0 in state  $q_i$  and square 0 contains symbol  $\sigma$ , then after  $y + 1$  steps it's still scanning square 0, is now in state  $q_j$ , the symbol on square 0 is  $\sigma'$ , and the squares other than square 0 contain the same symbols they contained after  $y$  steps.

- c) For every instruction  $\delta(q_i, \sigma) = \langle q_j, \sigma', N \rangle$ , the **sentence**:

tms:und:rep:  
rep-stay

$$\begin{aligned} \forall x \forall y ((Q_{q_i}(x, y) \wedge S_\sigma(x, y)) \rightarrow \\ (Q_{q_j}(x, y') \wedge S_{\sigma'}(x, y') \wedge \varphi(x, y))) \end{aligned}$$

Let  $\tau(M, w)$  be the conjunction of all the above **sentences** for Turing machine  $M$  and input  $w$



In order to express that  $M$  eventually halts, we have to find a **sentence** that says “after some number of steps, the transition function will be undefined.” Let  $X$  be the set of all pairs  $\langle q, \sigma \rangle$  such that  $\delta(q, \sigma)$  is undefined. Let  $\alpha(M, w)$  then be the **sentence**

$$\exists x \exists y \left( \bigvee_{\langle q, \sigma \rangle \in X} (Q_q(x, y) \wedge S_\sigma(x, y)) \right)$$

If we use a Turing machine with a designated halting state  $h$ , it is even easier: then the **sentence**  $\alpha(M, w)$

$$\exists x \exists y Q_h(x, y)$$

expresses that the machine eventually halts.

*tms:und:rep: prop:mlessk* **Proposition und.8.** *If  $m < k$ , then  $\tau(M, w) \models \bar{m} < \bar{k}$*

*Proof.* Exercise. □

**Problem und.4.** Prove **Proposition und.8.** (Hint: use induction on  $k - m$ ).

## und.6 Verifying the Representation

*tms:und:ver: sec* In order to verify that our representation works, we have to prove two things. *explanation* First, we have to show that if  $M$  halts on input  $w$ , then  $\tau(M, w) \rightarrow \alpha(M, w)$  is valid. Then, we have to show the converse, i.e., that if  $\tau(M, w) \rightarrow \alpha(M, w)$  is valid, then  $M$  does in fact eventually halt when run on input  $w$ .

The strategy for proving these is very different. For the first result, we have to show that a **sentence** of first-order logic (namely,  $\tau(M, w) \rightarrow \alpha(M, w)$ ) is valid. The easiest way to do this is to give a **derivation**. Our proof is supposed to work for all  $M$  and  $w$ , though, so there isn't really a single **sentence** for which we have to give a derivation, but infinitely many. So the best we can do is to prove by induction that, whatever  $M$  and  $w$  look like, and however many steps it takes  $M$  to halt on input  $w$ , there will be a **derivation** of  $\tau(M, w) \rightarrow \alpha(M, w)$ .

Naturally, our induction will proceed on the number of steps  $M$  takes before it reaches a halting configuration. In our inductive proof, we'll establish that for each step  $n$  of the run of  $M$  on input  $w$ ,  $\tau(M, w) \models \chi(M, w, n)$ , where  $\chi(M, w, n)$  correctly describes the configuration of  $M$  run on  $w$  after  $n$  steps. Now if  $M$  halts on input  $w$  after, say,  $n$  steps,  $\chi(M, w, n)$  will describe a halting configuration. We'll also show that  $\chi(M, w, n) \models \alpha(M, w)$ , whenever  $\chi(M, w, n)$  describes a halting configuration. So, if  $M$  halts on input  $w$ , then for some  $n$ ,  $M$  will be in a halting configuration after  $n$  steps. Hence,  $\tau(M, w) \models \chi(M, w, n)$  where  $\chi(M, w, n)$  describes a halting configuration, and since in that case  $\chi(M, w, n) \models \alpha(M, w)$ , we get that  $\tau(M, w) \models \alpha(M, w)$ , i.e., that  $\models \tau(M, w) \rightarrow \alpha(M, w)$ .

The strategy for the converse is very different. Here we assume that  $\models \tau(M, w) \rightarrow \alpha(M, w)$  and have to prove that  $M$  halts on input  $w$ . From the hypothesis we get that  $\tau(M, w) \models \alpha(M, w)$ , i.e.,  $\alpha(M, w)$  is true in every **structure**

in which  $\tau(M, w)$  is true. So we'll describe a **structure**  $\mathfrak{M}$  in which  $\tau(M, w)$  is true: its domain will be  $\mathbb{N}$ , and the interpretation of all the  $Q_q$  and  $S_\sigma$  will be given by the configurations of  $M$  during a run on input  $w$ . So, e.g.,  $\mathfrak{M} \models Q_q(\bar{m}, \bar{n})$  iff  $T$ , when run on input  $w$  for  $n$  steps, is in state  $q$  and scanning square  $m$ . Now since  $\tau(M, w) \models \alpha(M, w)$  by hypothesis, and since  $\mathfrak{M} \models \tau(M, w)$  by construction,  $\mathfrak{M} \models \alpha(M, w)$ . But  $\mathfrak{M} \models \alpha(M, w)$  iff there is some  $n \in |\mathfrak{M}| = \mathbb{N}$  so that  $M$ , run on input  $w$ , is in a halting configuration after  $n$  steps.

**Definition und.9.** Let  $\chi(M, w, n)$  be the **sentence**

$$Q_q(\bar{m}, \bar{n}) \wedge S_{\sigma_0}(\bar{0}, \bar{n}) \wedge \cdots \wedge S_{\sigma_k}(\bar{k}, \bar{n}) \wedge \forall x (\bar{k} < x \rightarrow S_0(x, \bar{n}))$$

where  $q$  is the state of  $M$  at time  $n$ ,  $M$  is scanning square  $m$  at time  $n$ , square  $i$  contains symbol  $\sigma_i$  at time  $n$  for  $0 \leq i \leq k$  and  $k$  is the right-most non-blank square of the tape at time 0, or the right-most square the tape head has visited after  $n$  steps, whichever is greater.

**Lemma und.10.** *If  $M$  run on input  $w$  is in a halting configuration after  $n$  steps, then  $\chi(M, w, n) \models \alpha(M, w)$ .*

tms:und:ver:  
lem:halt-config-implies-halt

*Proof.* Suppose that  $M$  halts for input  $w$  after  $n$  steps. There is some state  $q$ , square  $m$ , and symbol  $\sigma$  such that:

1. After  $n$  steps,  $M$  is in state  $q$  scanning square  $m$  on which  $\sigma$  appears.
2. The transition function  $\delta(q, \sigma)$  is undefined.

$\chi(M, w, n)$  is the description of this configuration and will include the clauses  $Q_q(\bar{m}, \bar{n})$  and  $S_\sigma(\bar{m}, \bar{n})$ . These clauses together imply  $\alpha(M, w)$ :

$$\exists x \exists y \left( \bigvee_{\langle q, \sigma \rangle \in X} (Q_q(x, y) \wedge S_\sigma(x, y)) \right)$$

since  $Q_{q'}(\bar{m}, \bar{n}) \wedge S_{\sigma'}(\bar{m}, \bar{n}) \models \bigvee_{\langle q, \sigma \rangle \in X} (Q_q(\bar{m}, \bar{n}) \wedge S_\sigma(\bar{m}, \bar{n}))$ , as  $\langle q', \sigma' \rangle \in X$ .  $\square$

explanation

So if  $M$  halts for input  $w$ , then there is some  $n$  such that  $\chi(M, w, n) \models \alpha(M, w)$ . We will now show that for any time  $n$ ,  $\tau(M, w) \models \chi(M, w, n)$ .

**Lemma und.11.** *For each  $n$ , if  $M$  has not halted after  $n$  steps,  $\tau(M, w) \models \chi(M, w, n)$ .*

tms:und:ver:  
lem:config

*Proof.* Induction basis: If  $n = 0$ , then the conjuncts of  $\chi(M, w, 0)$  are also conjuncts of  $\tau(M, w)$ , so entailed by it.

Inductive hypothesis: If  $M$  has not halted before the  $n$ th step, then  $\tau(M, w) \models \chi(M, w, n)$ . We have to show that (unless  $\chi(M, w, n)$  describes a halting configuration),  $\tau(M, w) \models \chi(M, w, n + 1)$ .

Suppose  $n > 0$  and after  $n$  steps,  $M$  started on  $w$  is in state  $q$  scanning square  $m$ . Since  $M$  does not halt after  $n$  steps, there must be an instruction of one of the following three forms in the program of  $M$ :

tms:und:ver:  
right  
tms:und:ver:  
left  
tms:und:ver:  
stay

1.  $\delta(q, \sigma) = \langle q', \sigma', R \rangle$
2.  $\delta(q, \sigma) = \langle q', \sigma', L \rangle$
3.  $\delta(q, \sigma) = \langle q', \sigma', N \rangle$

We will consider each of these three cases in turn.

1. Suppose there is an instruction of the form (1). By [Definition und.7\(3a\)](#), this means that

$$\forall x \forall y ((Q_q(x, y) \wedge S_\sigma(x, y)) \rightarrow (Q_{q'}(x', y') \wedge S_{\sigma'}(x, y') \wedge \varphi(x, y)))$$

is a conjunct of  $\tau(M, w)$ . This entails the following [sentence](#) (universal instantiation,  $\bar{m}$  for  $x$  and  $\bar{n}$  for  $y$ ):

$$(Q_q(\bar{m}, \bar{n}) \wedge S_\sigma(\bar{m}, \bar{n})) \rightarrow (Q_{q'}(\bar{m}', \bar{n}') \wedge S_{\sigma'}(\bar{m}, \bar{n}') \wedge \varphi(\bar{m}, \bar{n})).$$

By induction hypothesis,  $\tau(M, w) \models \chi(M, w, n)$ , i.e.,

$$Q_q(\bar{m}, \bar{n}) \wedge S_{\sigma_0}(\bar{0}, \bar{n}) \wedge \cdots \wedge S_{\sigma_k}(\bar{k}, \bar{n}) \wedge \forall x (\bar{k} < x \rightarrow S_0(x, \bar{n}))$$

Since after  $n$  steps, tape square  $m$  contains  $\sigma$ , the corresponding conjunct is  $S_\sigma(\bar{m}, \bar{n})$ , so this entails:

$$Q_q(\bar{m}, \bar{n}) \wedge S_\sigma(\bar{m}, \bar{n})$$

We now get

$$\begin{aligned} & Q_{q'}(\bar{m}', \bar{n}') \wedge S_{\sigma'}(\bar{m}, \bar{n}') \wedge \\ & S_{\sigma_0}(\bar{0}, \bar{n}') \wedge \cdots \wedge S_{\sigma_k}(\bar{k}, \bar{n}') \wedge \\ & \forall x (\bar{k} < x \rightarrow S_0(x, \bar{n}')) \end{aligned}$$

as follows: The first line comes directly from the consequent of the preceding conditional, by modus ponens. Each conjunct in the middle line—which excludes  $S_{\sigma_m}(\bar{m}, \bar{n}')$ —follows from the corresponding conjunct in  $\chi(M, w, n)$  together with  $\varphi(\bar{m}, \bar{n})$ .

If  $m < k$ ,  $\tau(M, w) \vdash \bar{m} < \bar{k}$  ([Proposition und.8](#)) and by transitivity of  $<$ , we have  $\forall x (\bar{k} < x \rightarrow \bar{m} < x)$ . If  $m = k$ , then  $\forall x (\bar{k} < x \rightarrow \bar{m} < x)$  by logic alone. The last line then follows from the corresponding conjunct in  $\chi(M, w, n)$ ,  $\forall x (\bar{k} < x \rightarrow \bar{m} < x)$ , and  $\varphi(\bar{m}, \bar{n})$ . If  $m < k$ , this already is  $\chi(M, w, n + 1)$ .

Now suppose  $m = k$ . In that case, after  $n + 1$  steps, the tape head has also visited square  $k + 1$ , which now is the right-most square visited. So

$\chi(M, w, n + 1)$  has a new conjunct,  $S_0(\bar{k}', \bar{n}')$ , and the last conjunct is  $\forall x (\bar{k}' < x \rightarrow S_0(x, \bar{n}'))$ . We have to verify that these two **sentences** are also implied.

We already have  $\forall x (\bar{k} < x \rightarrow S_0(x, \bar{n}))$ . In particular, this gives us  $\bar{k} < \bar{k}' \rightarrow S_0(\bar{k}', \bar{n})$ . From the axiom  $\forall x x < x'$  we get  $\bar{k} < \bar{k}'$ . By modus ponens,  $S_0(\bar{k}', \bar{n})$  follows.

Also, since  $\tau(M, w) \vdash \bar{k} < \bar{k}'$ , the axiom for transitivity of  $<$  gives us  $\forall x (\bar{k}' < x \rightarrow S_0(x, \bar{n}'))$ . (We leave the verification of this as an exercise.)

2. Suppose there is an instruction of the form (2). Then, by **Definition und.7(3b)**,

$$\begin{aligned} &\forall x \forall y ((Q_q(x', y) \wedge S_\sigma(x', y)) \rightarrow \\ &\quad (Q_{q'}(x, y') \wedge S_{\sigma'}(x', y') \wedge \varphi(x, y))) \wedge \\ &\forall y ((Q_{q_i}(o, y) \wedge S_\sigma(o, y)) \rightarrow \\ &\quad (Q_{q_j}(o, y') \wedge S_{\sigma'}(o, y') \wedge \varphi(o, y))) \end{aligned}$$

is a conjunct of  $\tau(M, w)$ . If  $m > 0$ , then let  $l = m - 1$  (i.e.,  $m = l + 1$ ). The first conjunct of the above **sentence** entails the following:

$$\begin{aligned} &(Q_q(\bar{l}', \bar{n}) \wedge S_\sigma(\bar{l}', \bar{n})) \rightarrow \\ &\quad (Q_{q'}(\bar{l}, \bar{n}') \wedge S_{\sigma'}(\bar{l}', \bar{n}') \wedge \varphi(\bar{l}, \bar{n})) \end{aligned}$$

Otherwise, let  $l = m = 0$  and consider the following **sentence** entailed by the second conjunct:

$$\begin{aligned} &((Q_{q_i}(o, \bar{n}) \wedge S_\sigma(o, \bar{n})) \rightarrow \\ &\quad (Q_{q_j}(o, \bar{n}') \wedge S_{\sigma'}(o, \bar{n}') \wedge \varphi(o, \bar{n}))) \end{aligned}$$

Either sentence implies

$$\begin{aligned} &Q_{q'}(\bar{l}, \bar{n}') \wedge S_{\sigma'}(\bar{m}, \bar{n}') \wedge \\ &\quad S_{\sigma_0}(\bar{0}, \bar{n}') \wedge \cdots \wedge S_{\sigma_k}(\bar{k}, \bar{n}') \wedge \\ &\quad \forall x (\bar{k} < x \rightarrow S_0(x, \bar{n}')) \end{aligned}$$

as before. (Note that in the first case,  $\bar{l}' \equiv \overline{l+1} \equiv \bar{m}$  and in the second case  $\bar{l} \equiv o$ .) But this just is  $\chi(M, w, n + 1)$ .

3. Case (3) is left as an exercise.

We have shown that for any  $n$ ,  $\tau(M, w) \vDash \chi(M, w, n)$ . □

**Problem und.5.** Complete case (3) of the proof of **Lemma und.11**.

**Problem und.6.** Give a **derivation** of  $S_{\sigma_i}(\bar{i}, \bar{n}')$  from  $S_{\sigma_i}(\bar{i}, \bar{n})$  and  $\varphi(m, n)$  (assuming  $i \neq m$ , i.e., either  $i < m$  or  $m < i$ ).

**Problem und.7.** Give a **derivation** of  $\forall x (\bar{k}' < x \rightarrow S_0(x, \bar{n}'))$  from  $\forall x (\bar{k} < x \rightarrow S_0(x, \bar{n}'))$ ,  $\forall x x < x'$ , and  $\forall x \forall y \forall z ((x < y \wedge y < z) \rightarrow x < z)$ .

tms:und:ver:  
lem:valid-if-halt

**Lemma und.12.** *If  $M$  halts on input  $w$ , then  $\tau(M, w) \rightarrow \alpha(M, w)$  is valid.*

*Proof.* By **Lemma und.11**, we know that, for any time  $n$ , the description  $\chi(M, w, n)$  of the configuration of  $M$  at time  $n$  is entailed by  $\tau(M, w)$ . Suppose  $M$  halts after  $k$  steps. It will be scanning square  $m$ , say. Then  $\chi(M, w, k)$  describes a halting configuration of  $M$ , i.e., it contains as conjuncts both  $Q_q(\bar{m}, \bar{k})$  and  $S_\sigma(\bar{m}, \bar{k})$  with  $\delta(q, \sigma)$  undefined. By **Lemma und.10** Thus,  $\chi(M, w, k) \models \alpha(M, w)$ . But since  $(M, w) \models \chi(M, w, k)$ , we have  $\tau(M, w) \models \alpha(M, w)$  and therefore  $\tau(M, w) \rightarrow \alpha(M, w)$  is valid.  $\square$

To complete the verification of our claim, we also have to establish the reverse direction: if  $\tau(M, w) \rightarrow \alpha(M, w)$  is valid, then  $M$  does in fact halt when started on input  $m$ . explanation

tms:und:ver:  
lem:halt-if-valid

**Lemma und.13.** *If  $\tau(M, w) \rightarrow \alpha(M, w)$ , then  $M$  halts on input  $w$ .*

*Proof.* Consider the  $\mathcal{L}_M$ -**structure**  $\mathfrak{M}$  with domain  $\mathbb{N}$  which interprets 0 as 0, ' as the successor function, and  $<$  as the less-than relation, and the predicates  $Q_q$  and  $S_\sigma$  as follows:

$$\begin{aligned} Q_q^{\mathfrak{M}} &= \{ \langle m, n \rangle : \begin{array}{l} \text{started on } w, \text{ after } n \text{ steps,} \\ M \text{ is in state } q \text{ scanning square } m \end{array} \} \\ S_\sigma^{\mathfrak{M}} &= \{ \langle m, n \rangle : \begin{array}{l} \text{started on } w, \text{ after } n \text{ steps,} \\ \text{square } m \text{ of } M \text{ contains symbol } \sigma \end{array} \} \end{aligned}$$

In other words, we construct the **structure**  $\mathfrak{M}$  so that it describes what  $M$  started on input  $w$  actually does, step by step. Clearly,  $\mathfrak{M} \models \tau(M, w)$ . If  $\tau(M, w) \rightarrow \alpha(M, w)$ , then also  $\mathfrak{M} \models \alpha(M, w)$ , i.e.,

$$\mathfrak{M} \models \exists x \exists y \left( \bigvee_{\langle q, \sigma \rangle \in X} (Q_q(x, y) \wedge S_\sigma(x, y)) \right).$$

As  $|\mathfrak{M}| = \mathbb{N}$ , there must be  $m, n \in \mathbb{N}$  so that  $\mathfrak{M} \models Q_q(\bar{m}, \bar{n}) \wedge S_\sigma(\bar{m}, \bar{n})$  for some  $q$  and  $\sigma$  such that  $\delta(q, \sigma)$  is undefined. By the definition of  $\mathfrak{M}$ , this means that  $M$  started on input  $w$  after  $n$  steps is in state  $q$  and reading symbol  $\sigma$ , and the transition function is undefined, i.e.,  $M$  has halted.  $\square$

## und.7 The Decision Problem is Unsolvable

tms:und:uns:  
sec  
tms:und:uns:  
thm:decision-prob

**Theorem und.14.** *The decision problem is unsolvable.*

*Proof.* Suppose the decision problem were solvable, i.e., suppose there were a Turing machine  $D$  of the following sort. Whenever  $D$  is started on a tape that contains a **sentence**  $\psi$  of first-order logic as input,  $D$  eventually halts, and outputs 1 iff  $\psi$  is valid and 0 otherwise. Then we could solve the halting problem as follows. We construct a Turing machine  $E$  that, given as input the number  $e$  of Turing machine  $M_e$  and input  $w$ , computes the corresponding **sentence**  $\tau(M_e, w) \rightarrow \alpha(M_e, w)$  and halts, scanning the leftmost square on the tape. The machine  $E \frown D$  would then, given input  $e$  and  $w$ , first compute  $\tau(M_e, w) \rightarrow \alpha(M_e, w)$  and then run the decision problem machine  $D$  on that input.  $D$  halts with output 1 iff  $\tau(M_e, w) \rightarrow \alpha(M_e, w)$  is valid and outputs 0 otherwise. By **Lemma und.13** and **Lemma und.12**,  $\tau(M_e, w) \rightarrow \alpha(M_e, w)$  is valid iff  $M_e$  halts on input  $w$ . Thus,  $E \frown D$ , given input  $e$  and  $w$  halts with output 1 iff  $M_e$  halts on input  $w$  and halts with output 0 otherwise. In other words,  $E \frown D$  would solve the halting problem. But we know, by **Theorem und.6**, that no such Turing machine can exist.  $\square$

## Photo Credits

# Bibliography