

# How to Make Zuse's Z3 a Universal Computer

Raúl Rojas

September 5, 1997

## Abstract

The computing machine Z3, built by Konrad Zuse from 1938 to 1941, could only execute fixed sequences of floating-point arithmetical operations (addition, subtraction, multiplication, division and square root) coded in a punched tape. We show in this paper that a single program loop containing this type of instructions can simulate any Turing machine whose tape is of bounded size. This is achieved by simulating conditional branching and indirect addressing by purely arithmetical means. Zuse's Z3 is therefore, at least in principle, as universal as today's computers which have a bounded memory size. This result is achieved at the cost of blowing up the size of the program stored on punched tape.

## Universal Machines and Single Loops

*Nobody has ever built a universal computer.* The reason is that a universal computer consists, in theory, of a fixed processor and a memory of unbounded size. This is the case of Turing machines with their unbounded tapes. In the theory of general recursive functions there is also a small set of rules and some predefined functions, but there is no upper bound on the size of intermediate reduction terms. Modern computers are only *potentially* universal: They can perform any computation that a Turing machine with a *bounded* tape can perform. If more storage is required, more can be added without having to modify the processor (provided that the extra memory is still addressable).



It is the purpose of this paper to show that Konrad Zuse's Z3, a computing automaton built in Berlin from 1938 to 1941, can be programmed *in principle* as any other universal computer. This is a surprising result, since the Z3 can only compute sequences of arithmetical operations (addition, subtraction, multiplication and division) stored in a punched tape. There is no conditional branching. Since both ends of the punched tape can be glued together, the Z3 is a machine capable of executing a single loop of arithmetical operations using numbers stored in memory.

It is well known that any computer program containing conditional branches and the usual instructions of imperative languages, like for example FORTRAN, can be programmed using a single WHILE loop [1]. Also, all conditional branches can be eliminated from the loop [2]. I showed in [3] that if the Z3 is extended with indirect addressing it can simulate a Turing machine. We will adopt the techniques used in those papers in order to show that a Turing machine can be simulated by a single program loop of a machine capable of computing the four basic arithmetic operations.

Our computing model is the following: there exist memory locations which we will denote by lower case letters. We can only refer explicitly to memory addresses (there is no indirect addressing). Initially (for the sake of simplicity) we will restrict our programs for the Z3 to a language containing only statements of the form

$$a = b \text{ op } c,$$

where op represents one of the four basic arithmetic operations. Any statement  $a = b \text{ op } c$  can be "compiled" using the two registers of the Z3 and four assembler instructions (which load the two argument registers in the appropriate order):

```
LOAD b
LOAD c
op
STORE a
```

The store operation refers implicitly to the first register (accumulator) of the processor. All computations are performed with floating-point numbers. The mantissa has a precision of 16 bits for its fractional part. The Z3 uses normalized floating-point numbers (i.e. with a mantissa  $m$  such that  $1 \leq m < 2$ ). The special case of a zero mantissa



is handled with a special code (like in the IEEE standard). There is also a "halt" instruction in the Z3 (when a number is displayed at the console the machine stops). For more details on the architecture of the Z3 see [4] and [5].

## Simulating Branches

We show here how to simulate the operation of a CASE statement using a technique introduced in [2]. Define the state of the machine as the state of its memory. Assume that in a program  $P$  there are  $n$  consecutive sections of code  $P_1, \dots, P_n$  and that the variable  $z \in \{1, 2, \dots, n\}$  (a number stored in memory) is used to select which one of these sections should perform the computation we are interested in. The general strategy is to execute *all*  $n$  sections of code, one after the other, but we will allow only the  $z$ -th section to modify the memory contents. In order to implement this idea we transform each section of code  $P_j$  in equivalent code  $P'_j$  according to the following recipe: At the beginning of each section  $P_j$  a comparison is made and if  $z = j$  the auxiliary variable  $t$  is set to zero, otherwise it is set to one. The variable  $t$  can be interpreted as a flag for the "selected section" since it will be only zero in  $P_z$ . Now all original statements in the program  $P_1, \dots, P_n$  of the form  $a = b \text{ op } c$  are transformed to

$$\begin{aligned}u &= b \text{ op } c \\v &= a * t \\w &= 1 - t \\u &= w * u \\a &= v + u\end{aligned}$$

The expression computed by this piece of code is  $a = a * t + (1 - t)(b \text{ op } c)$ , that is the state of variable  $a$  will not be modified unless the computation is performed within the  $z$ -th code section. When all statements have been transformed in this way and the appropriate initialization of  $t$  has been introduced at the beginning of each code section, we can execute the whole transformed program  $P'_1, \dots, P'_n$  from beginning to end. Most of the computations are superfluous, since we execute all sections of code, but only  $P'_z$  modifies the contents of the memory, as we demand from our CASE statement.



We must only show now that it is in fact possible to perform the computation

if ( $z = i$ ) then  $t = 0$  else  $t = 1$ .

We do this using the following piece of code (where  $e$  represents a positive number much smaller than 1, for example 1/1000).

$$\begin{aligned}d &= z - i \\d &= d * d \\f &= d - e \\g &= d/f\end{aligned}$$

The variable  $g$  is equal to zero only if  $z = i$ . When  $z \neq i$  the variable  $g$  is equal to 1 plus a small fractional part. We want to get rid of these extra bits after the integer part of the mantissa (those to the right of the binary point). The way to do this is to compute

$$t = (2^{16} + g) - 2^{16},$$

respecting the implicit order. Remember that the Z3 represents numbers internally using 16 bits for the fractional part of the mantissa. In the addition above the smaller number (the variable  $g$ ) is shifted 16 places to the right, whereby the extra bits we want to eliminate are lost. The subtraction "restores" the integer part of  $g$  to its initial state. In this and other programs all necessary constants ( $2^{16}$  for example), can be precomputed and stored before we start the CASE statement. We assume that the processor is able to handle big enough floating-point numbers so that no overflows arise.

## Simulating a Turing Machine

A Turing machine (TM) is defined by a table of state transitions: given the current state  $Q$  and the tape symbol at the current position  $pos$  of the read/write head, we read from the table and find the new state  $Q'$ , the symbol to be written  $o$  and the direction  $dir$  of motion of the read/write head (+1 or -1). The new position of the head is given by  $pos = pos + dir$ . Before simulating a Turing machine, the memory of the Z3 is prepared. All necessary tables are loaded at specific addresses and the initial contents of the TM tape too. All necessary auxiliary constants are also loaded.



It is clear that any Turing machine can be simulated using the following master loop:

- read tape symbol,
- look-up new state, output symbol, and direction of movement
- modify tape,
- update state and position of read/write head.

The simulation can be done using table look-up. For example, reading the tape symbol amounts to the operation

$$s = \text{memory}(\text{tape}, \text{pos}),$$

where “tape” is the initial address of the simulated tape and “pos” the current position. Only basic arithmetic is needed to compute the position of entries in a table: A table starting at address  $T$  can be accessed at position  $k$  by computing  $T + k$  and using the result as an address. Thus, the only thing we need for the simulation of the Turing machine using the Z3, and which is still lacking, is indirect addressing: We want to use the results of arithmetic operations as addresses.

Assume that we want to implement the indirect addressing operation

$$a = (x)$$

where  $x$  is the result of an arithmetic operation with integers and  $n_a < x < n_b$ . The integer constants  $n_a$  and  $n_b$  are the limits of the memory segment that we want to address indirectly.

We can implement the above operation using a CASE statement with one section for each integer between  $n_a$  and  $n_b$ . In each section  $i$  of the CASE statement we load address  $a$  with the contents of address  $i$ . Assume that  $n_a = 10$  and  $n_b = 20$ . The code, before transforming it to work as a CASE statement, would be:

```
P10:  LOAD 10
      STORE a
P11:  LOAD 11
      STORE a
...
P20:  LOAD 20
      STORE a
```



Now we apply a transformation similar to the one discussed above using  $x$  as the CASE variable. The transformed program will select the contents of address  $x$  and will store it in address  $a$ , since only section  $P_x$  will modify the contents of  $a$ .

Note that the whole CASE statement contains one load statement for each consecutive memory address. We read all memory addresses between  $n_a$  and  $n_b$ , but we only keep in  $a$  the one we are interested in, namely address  $x$ . In the extreme case, when the indirect addressing operation refers to the whole memory, we would need to read all addresses in order to implement a single indirect addressing. But since the number of indirect memory references during one simulation cycle of the Turing machine is constant, the size of the program that we need is also constant (for a given memory size).

Using an analogous approach we can store a number to the address represented by an arithmetical result  $x$  (indirect addressing in STORE operations). We can, for example, update the simulated tape of the TM using this approach.

It is clear that we have been helped here by the fact that the program is stored in a punched tape which is external to the memory. The punched tape is allowed to be as large as necessary to read the sections of memory that we need to address indirectly (state tables and tape of the TM). Given a maximal memory size, the size of the punched tape needed is bounded.

This proves that Zuse's Z3 can, in principle, do any computation that any other computer with a bounded memory can perform.

## The Halting Problem

The attentive reader will have noticed that the master loop of the simulation never stops. Algorithms, however, must stop after a finite number of steps. Fortunately, there is an additional feature of the Z3 which provides the solution for this problem.

Whenever an undefined operation is performed, the Z3 stops and a lamp is set on the console. This is the case, for example, for the operation  $0/0$ . Thus we define state  $Q_0 = 0$  of the simulation as the "halting state" (for all other states  $Q_i$  is a positive integer) and the computation  $0/Q$  is performed at the beginning of the master loop ( $Q$  is the current state). If the simulation reaches state  $Q_0$  the machine



stops.

If Zuse had not thought of trapping undefined operations, we would have been unable to stop the master loop. In that case, a possible way out could be considering those cycles in which nothing is altered as the "halting state" of the machine, but the operator would have some problems identifying this situation.

## Conclusions

The main result of this paper is intriguing because it looks so artificial. From the theoretical point of view it is interesting to see that limited precision floating-point arithmetic embedded in a WHILE loop can compute anything computers can compute. It could be argued that whenever we expand the memory (to accommodate more tape positions for a Turing machine) the program in the punched tape has to be expanded also (to cover the new memory addresses). If we think of the punched tape as part of the processor (when simulating a Universal Turing Machine), then we are extending the processor when we enlarge the program in the punched tape. This is undesirable. However, in real computers, there is also a limit for the size of the memory we can manage (given by the addressable space, i.e. the number of bits in the address registers). If we expand the memory we need more addressing bits and the processor has to be expanded (going for example from 16-bit to 32-bit registers).

The result of this paper seems counterintuitive, until we realize that operations like multiplication and division are iterative computations in which branching decisions are taken by the hardware. The conditional branchings we need are *embedded* in these arithmetical operations and the whole purpose of the transformations used, is to lift the branches up from the hardware, in which they are buried, to the software level so that we can control the program flow. The whole magic of the transformation is making the hardware branchings visible to the programmer.

A possible criticism is that the approach discussed in this paper produces a terrible slowdown of the computation. From a purely theoretical point of view this is irrelevant. From a practical point of view it is clearly the case that nobody would program the Z3 in this way, in the same way that nobody solves real-world problems using Turing



machines.

We can therefore say that, from an *abstract theoretical perspective*, the computing model of the Z3 is equivalent to the computing model of today's computers. From a practical perspective, and in the way the Z3 was really programmed, it was not equivalent to modern computers.

I am of course aware that these conclusions are curious enough to re ignite the whole discussion about the invention of the computer.

## References

- [1] D. Harel, "On Folk Theorems", *Communications of the ACM*, Vol. 23, N. 7, 1980, pp. 379–389.
- [2] O. Ibarra, S. Moran, L.E. Rosier, "On the Control Power of Integer Division", *Theoretical Computer Science*, Vol. 24, 1983, pp. 35–52.
- [3] R. Rojas, "Conditional Branching is not Necessary for Universal Computation in von Neumann Computers", *Journal of Universal Computer Science*, Vol. 2, N. 11, 1996, pp. 756–767.
- [4] R. Rojas, "Konrad Zuse's Legacy: the Architecture of the Z1 and Z3", *Annals of the History of Computing*, Vol. 19, N. 2, 1997, pp. 5–16.
- [5] R. Rojas, *Die Rechenmaschinen von Konrad Zuse*, Springer-Verlag, Berlin, 1997.