

Konrad Zuse's Legacy: The Architecture of the Z1 and Z3

RAÚL ROJAS

This paper provides a detailed description of the architecture of the Z1 and Z3 computing machines that Konrad Zuse designed in Berlin between 1936 and 1941. The necessary basic information was obtained from a careful evaluation of the patent application Zuse filed in 1941. Additional insight was gained from a software simulation of the machine's logic. The Z1 was built using purely mechanical components; the Z3 used electromechanical relays. However, both machines shared a common logical structure, and their programming model was the same. I argue that both the Z1 and the Z3 possessed features akin to those of modern computers: The memory and processor were separate units, and the processor could handle floating-point numbers and compute the four basic arithmetical operations as well as the square root of a number. The program was stored on punched tape and was read sequentially. In the last section of this paper, I put the architecture of the Z1 and Z3 into historical perspective by offering a comparison with computing machines built in other countries.

Introduction

Konrad Zuse is popularly recognized in Germany as the father of the computer, and his Z1, a programmable automaton built from 1936 to 1938, has been called the first computer in the world. Other nations reserve this honor for one of their own scientists, and there has been a long and often acrimonious debate on the issue of who is the true inventor of the computer. Sometimes the discussion is preempted by specifying in full detail the technological features of a specific machine. The Electronic Numerical Integrator and Computer (ENIAC), for example, has been called the first "large-scale general-purpose electronic computer in the world."² The ENIAC was built at the Moore School of Electrical Engineering of the University of Pennsylvania from May 1943 to 1945. It solved its first problem in December 1945 and was officially presented in February 1946. Another contender for the title of the first computer is the Mark I, built by Howard Aiken at Harvard University between 1939 and 1944. The Mark I was an electromechanical machine, not of the all-mechanical nature of previous computing devices and not built with the electronics available at the time.¹ The machine John Atanasoff built (later called the ABC) at Iowa State College from 1938 to 1942 used vacuum tubes but was restricted to the addition and subtraction of vectors and had a structure inappropriate for universal computation.³ In direct contrast to these three machines, the Z1 was more flexible and was designed to execute a long and modifiable sequence of instructions contained on a punched tape. Zuse's machines, the Z3 and the Z4, were not electronic and were of reduced size. Since the Z3 was completed and was successfully working prior to the Mark I, it has been called the first *program-*

mable calculating machine in the world. Of course the old debate will not be closed with this paper, but I want to show here just how advanced the machines Zuse built were when considered from the viewpoint of modern computer architecture and compared with other designs of that time.

The Berlin Polytechnic student Zuse started thinking about computing machines in the 1930s. He realized that he could construct an automaton capable of executing a sequence of arithmetical operations like those needed to compute mathematical tables. Coming from a civil engineering background, he had no formal training in electronics and was not acquainted with the technology used in conventional mechanical calculators. This nominal deficit worked to his advantage, however, because he had to rethink the whole problem of arithmetic computation and thus hit on new and original solutions.

Zuse decided to build his first experimental calculating machine exploiting two main ideas:

- the machine would work with binary numbers
- the computing and control unit would be separated from the storage.

Years before John von Neumann explained the advantages of a computer architecture in which the processor is separated from the memory, Zuse had already arrived at the same solution. However, it must be said that Charles Babbage had the same idea in the previous century when he designed his Analytical Engine. In 1936, Zuse completed the memory of the machine he had planned. (Zuse called it the *Speicherwerk* (storage mecha-

the Z3
omput-
worked
adopt
onds to

the Z3.
sor and
able of
proces-
metical
e expo-
n. The
ach in-
proces-
ynchro-
of each
The I/O
unit.

equivalent to a significand of 15 bits. However, there is a problem with the number zero, which cannot be expressed using a normalized significand. The Z3 uses the convention that any significand with exponent -64 is to be considered equal to zero. Any number with exponent 63 is considered infinitely large. Operations involving zero and infinity are treated as exceptions, and special hardware monitors the numbers loaded in the processor in order to set the exception flags (see below). With this convention, the smallest number representable in the memory of the Z3 is $2^{-63} = 1.08 \times 10^{-19}$, and the largest is $1.999 \times 2^{62} = 9.2 \times 10^{18}$. The arguments for computations can be entered as decimal numbers on the keyboard of the Z3 (four digits). The exponent of the decimal representation is entered by pushing the appropriate button in a row of buttons labeled -8, -7, ..., 7, 8. The original Z3 could accept input only between 1×10^{-8} and $9,999 \times 10^8$. Zuse's reconstruction of the Z3 for the Deutsches Museum in Munich provides enough buttons for larger exponents. With this arrangement, the whole numerical capacity of the machine can be reflected on the acceptable input. The same can be said of the output. However, the Z3 does not print the numerical results the program produces. A single number is displayed on an array of lamps representing the digits from zero to nine. The largest number that can be displayed is 19,999. The smallest is 00001. The largest exponent that can be displayed is +8, the smallest -8.

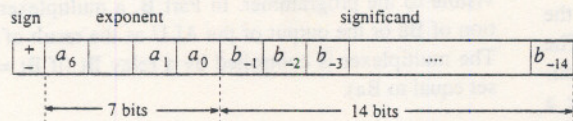


Fig. 2. The floating-point representation in memory.

Instruction Set

The program for the Z3 is stored on punched tape. One instruction is coded using eight bits for each row of the tape. The instruction set of the Z3 consists of the nine instructions shown in Table 1. There are three types of instructions: I/O, memory, and arithmetical operators. The opcode has a variable length of two or five bits. Memory operations encode the address of a word in the lower six bits, that is, the addressing space has a maximum size of 64 words, as mentioned above.

The instructions on the punched tape can be arranged in any order. The instructions Lu and Ld (read from keyboard and display result, respectively) halt the machine, so that the operator has enough time to input a number or write down a result. The machine is then restarted and continues processing the program.

The instruction most conspicuously absent from the instruction set of the Z3 is conditional branching. Loops can be implemented by the simple expedient of bringing together the two ends of the punched tape, but there is no way to implement conditional sequences of instructions. The Z3 is therefore not a universal computer in the sense of Turing.

Number of Cycles

The Z3 is a clocked machine. Each cycle is divided into five stages called I, II, III, IV, and V. The instruction in the punched tape is decoded in Stage I of a cycle. The two basic arithmetical

operations of the machine are addition and subtraction of exponents and significands. The operations can be executed in the first three stages of each cycle. Stages IV and V are used to prepare arguments for the next operation or to write back results.

TABLE 1
INSTRUCTION SET AND OPCODES OF THE Z3

Type	Instruction	Description	Opcode
I/O	Lu	read keyboard	01 110000
	Ld	display result	01 111000
memory	Pr z	load address z	11 $z_6 z_5 z_4 z_3 z_2 z_1$
	Ps z	store address z	10 $z_6 z_5 z_4 z_3 z_2 z_1$
arithmetic	Lm	multiplication	01 001000
	Li	division	01 010000
	Lw	square root	01 011000
	Ls ₁	addition	01 100000
	Ls ₂	subtraction	01 101000

The instructions implemented in the Z3 require the following number of cycles:

Multiplication:	16 cycles
Division:	18 cycles
Square root:	20 cycles
Addition:	3 cycles
Subtraction:	4 or 5 cycles, depending on the result
Read keyboard:	9 to 41 cycles, depending on the exponent
Display output:	9 to 41 cycles, depending on the exponent
Load from memory:	1 cycle
Store to memory:	0 or 1 cycle

According to Zuse, the time required for a multiplication was three seconds. Considering that a multiplication operation needs 16 cycles, one can estimate that the operating frequency of the Z3 was $16/3 \approx 5.33$ Hz. It is a curious coincidence that the gate-level simulation of the Z3 that my students implemented using a personal computer also required around three seconds for a multiplication.

The instruction most conspicuously absent from the instruction set of the Z3 is conditional branching.

The number of cycles needed for the *read* and *display* instructions is variable, because it depends on the exponent of the arguments. Since the input has to be converted from decimal to binary representation, the number of multiplications needed with the factor 10 or 0.1 is dictated by the decimal exponent (see below).

Addition and subtraction require more than one cycle because, in the case of floating-point numbers, care has to be taken to set the size of the exponent of both arguments to the same value. This requires some extra comparisons and shifting.

A number can be stored in memory in *zero* cycles when the result of the last arithmetical operation can be redirected to the desired memory address. In this case, the cycle needed for the store instruction overlaps the last cycle of the arithmetical operation.

Bf and Ba and a shifter between Bf and Bb. The first shifter can displace the significand up to two positions to the right and one position to the left. This amounts to a division of Bf by four or a multiplication with the constant two. The second shifter can displace the significand in Af from one to 16 positions to the right and from one to 15 positions to the left. These shifts are needed for addition and subtraction of floating-point numbers. Multiplication and division with powers of two can therefore be performed when the operands for the next arithmetical operation are fetched and, in this sense, do not consume time.

The number of bits used in the registers is as follows:

Af	7 bits	Bf	17 bits
Aa	8	Ba	19
Ab	8	Bb	18
Ae	8	Be	18

As can be seen from this list, Ae uses one extra bit to handle the addition of the exponents of the arguments. Part B of the processor uses two extra bits for the significands (b_{-15} and b_{-16}) and makes explicit b_0 , which is not stored in memory. The extra bits at positions -15 and -16 are included to increase the precision of the computations. Therefore, the total number of bits needed to store the result of an arithmetical operation in Bf is 17 bits. Registers Ba and Bb require more extra bits (ba_2 , ba_1 , and bb_1) to handle intermediate results of some of the numerical algorithms. In particular, the square root algorithm can lead to partial computations in Ba requiring three bits to the left of the decimal point.

The basic primitive operation of the data path is the addition or subtraction of exponents or significands. When the relay As or Bs is set, the negation of the second argument (Ab or Bb) is fed into the ALU. Therefore, if the relay As is set to one, the ALU in Part A subtracts its arguments, otherwise they are added. The same is true for Part B and the relay Bs. The constant of one is needed to build the two's complement of a number.

Assume that two numbers with the same exponent are to be added. The first exponent is stored in Af, the second in Ab. Since they are equal, no operation has to be performed on this side of the machine. In Part B, the significand of the first number is stored in Bf and the significand of the second in Bb. The first step consists of loading Ba with Bf by setting the relay box Fa to one. The addition is performed next, the relay Bt is set to one, and so the result $Ba + Bb$ is assigned to Be. The relay box Ff is now set to one, and the result is stored in Bf. As one can see, the information can move between registers and so flow through the data path. The computer architect has to provide the correct sequence of activations of the relay boxes in order to get the desired operation. This is done in the Z3 using a technique very similar to microprogramming.

The Control Unit

Fig. 4 shows a more detailed diagram of the control unit and of the I/O panels. The circuit Pa decodes the opcode of the instruction read from the punched tape. If it is a memory instruction, circuit Pb sets the address bus to the value of the lower six bits of the opcode. The control unit determines the correct microsequencing of the instructions. There are special circuits for each of the operations in the instruction set.

Circuit Z represents the panel of buttons used to enter a decimal number in the machine. Only one button in each of the four columns can be activated. The exponent is set by pressing one of the buttons labeled -8 to 8 in circuit K. The output display is very similar to the input panel, but here lamps illuminate the appropriate decimal digits, the exponent of the number (circuit Q), as well as its sign. Note that there is a fifth digit for the output (which can be only one or zero).

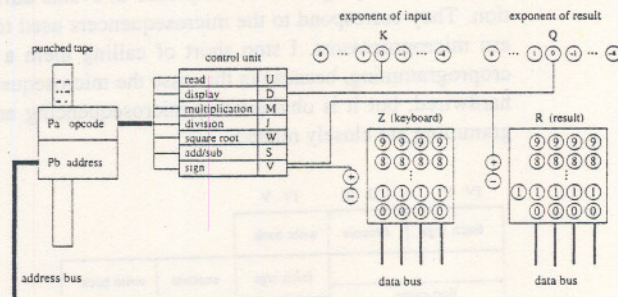


Fig. 4. The control unit and I/O panels.

Once a decimal number has been set, a data bus transmits the digits to register Ba, and a complex series of operations is started. The decimal input must be transformed into a binary number. This requires a chain of multiplications, which is longer according to the absolute magnitude of the exponent. If the exponent is zero, the whole transformation requires nine cycles, but if it is -8 , the operation requires $9 + 4 \times 8 = 41$ cycles.

There are a lot of details that the engineer designing the "microprogram" must keep in mind, otherwise short circuits can destroy the hardware.

Microcontrol of the Z3

The heart of the control unit is made up of its microsequencers. Before I describe the way they work, it is necessary to take a closer look at the chaining of arithmetical instructions in the Z3. Fig. 5 shows the main idea. Each cycle of the Z3 is divided into five stages. Stages IV and V are used to move information around in the machine. During Stages I, II, and III, an addition/subtraction is computed in Part A and another in Part B of the Z3. I call this the "execute" phase of an instruction. A typical instruction fetches its arguments, executes, and writes back the result. Zuse took great care to save execution time by overlapping the fetch stage of the next instruction with the write-back stage of the current one. One can think of an execution cycle as consisting of just two stages, as shown in Fig. 5, where the first two cycles of a series of instructions have been labeled. I have adopted this convention in the tabular diagrams of the numerical algorithms discussed later in this article.

The microsequencing is done by special control wheels. There is one for the multiplication algorithm, another to control division, and yet another for the square root instruction. The moving arm shown in Fig. 6 starts moving clockwise as soon as the control

to one if the number is stored in the register pair $\langle Af, Bf \rangle$. If the number is stored in the register pair $\langle Ab, Bb \rangle$, the relay Nn_2 is set to one. In this way, we always know if one or both of the arguments for an arithmetical operation are zero. Something similar is done for any exponent of value 63 (an infinite number, according to the convention). In this case, the relays Ni_1 or Ni_2 are set to one, according to the register pair in which the number is stored.

Operations involving "exceptional" numbers (zero or infinity) are performed as usual, but the result is overridden by the snooping circuit. Assume, for example, that a multiplication is computed and that the first argument is zero (Nn_1 is set to one). The computation proceeds as usual, but in each cycle the snooping circuit produces the result -64 at the output of the adder of Part A. It does not matter what operations are performed with the significands because the exponent of the result is set to -64 , and therefore the final result is zero. Division by an infinite number can be handled in a similar manner. The Z3 can detect undefined operations such as $0/0$, $\infty - \infty$, ∞/∞ , and $0 \times \infty$. In all these cases, the corresponding exception lamp lights on the output panel, and the machine is stopped. The Z3 always produces the correct result when one of the arguments is zero or ∞ and the other is a number within bounds. This was not the case for the Z1. Zuse thought of, but did not implement, exception handling in the Z1. The machine could not correctly perform some computations involving zero.¹⁶

An additional circuit looks at the exponent of the result at the output of the exponent's adder. If the exponent is greater than or equal to 63, overflow has occurred and the result must be set to ∞ . If the exponent is lower than -64 , underflow has occurred and the result must be set to zero. To do this, the appropriate relay (Nn_1 or Ni_1) is set to one.

Zuse managed to implement exception handling using just a few relays. This feature of the Z3 is one of the most elegant in the whole design. Many of the early microprocessors of the 1970s did not include exception handling and left it to the software. Zuse's approach is sounder, since it frees programmers from the tedium of checking the bounds of their numbers before each operation.

Addition and Subtraction

In order to add or subtract two floating-point numbers x and y , their representation must be reduced to the same exponent. After this has been done, only the significands have to be added or subtracted. If the exponents are different, the significand of the smaller number is shifted to the right as many places as necessary (and its exponent is incremented correspondingly to keep the number unchanged) until both exponents are equal. It can, of course, happen that the smaller number becomes zero after 17 shifts to the right.

The signs of the two numbers are compared before deciding on the type of operation to be executed. If an addition has been requested and the signs are the same, the addition is performed. If the signs are different, a subtraction is executed. If a subtraction has been requested and the signs are different, an addition is executed. If the signs are the same, the subtraction is executed. A special circuit sets the sign of the final result according to the signs of the arguments and the sign of the partial result.

Addition and subtraction are controlled by a chain of relays (not by a control wheel), since the maximum number of cycles needed is low. Fig. 8 shows the synchronization required for the addition of two numbers. Initially, the arguments for the addition

are stored in the register pairs $\langle Af, Bf \rangle$ and $\langle Ab, Bb \rangle$. In the first cycle, the exponents are subtracted. In Cycle 2, the significand with the larger exponent is loaded into register Ba, and the significand with the smaller exponent is loaded into register Bb. The significand in register Bb is shifted as many places to the right as the absolute value of the difference of the exponents (exception handling takes care of the case in which the smaller number becomes zero after the shift). In Stages I, II, and III of Cycle 2, the significands are added, and finally the processor tests if the result is greater than two. If this is the case, the significand is shifted one position to the right and the exponent is incremented by one. Note that the test "if $(Be \geq 2)$ " in Part A of the arithmetical unit is done *after* Be has already been computed in Part B during Stages I, II, and III of Cycle 2.

Zuse managed to implement exception handling using just a few relays. This feature of the Z3 is one of the most elegant in the whole design.

In the case of a subtraction, four or five cycles are needed. Fig. 9 shows the synchronization required for a subtraction. The first two cycles are almost identical to the first two cycles of the addition algorithm, but now the significands are subtracted. Cycle 3 is executed only when the difference of the significands is negative. The effect of Cycle 3 is just to make the significand of the result positive. Cycle 4 is very important: The difference of two normalized significands can have many zeros in the first bit positions to the left. The result is normalized by shifting Be to the left as many places as necessary (this is done with the shifter between the relay box Fd and register Bb). The number of one-bit shifts is subtracted from the exponent in Part A of the processor. In Cycle 5, the result is stored in the register pair $\langle Af, Bf \rangle$.

cycle	stage	exponent	significand
0	I, II, III		
1	IV, V	$Aa := Af$	
		$Ae := Aa - Ab$	$Be := 0 + Bb$
2	IV, V	if $(Ae \geq 0)$ then $Ab := 0$, $Aa := Af$ else $Aa := 0$	if $(Ae \geq 0)$ then $Ba := Bf$, $Bb := Be$ (shifted) else $Ba := Be$, $Bb := Bf$ (shifted) (Be or Bf are shifted $ Ae $ places to the right)
		I, II, III if $(Be \geq 2)$ then $Ae := Aa + Ab + 1$ else $Ae := Aa + Ab$	$Be := Ba + Bb$
3	IV, V	$Af := Ae$	if $(Be \geq 2)$ then $Bf := Be/2$ else $Bf := Be$

Fig. 8. The three cycles needed for the addition algorithm. The arguments for the addition are stored in the register pairs $\langle Af, Bf \rangle$ and $\langle Ab, Bb \rangle$ before the operation is started.

Multiplication

The multiplication algorithm of the Z3 is like the one used for decimal multiplication by hand, that is, it is based on repeated additions of the multiplicator according to the individual digits of the multiplicand. At the beginning of the algorithm, the first argument is stored in the register pair $\langle Af, Bf \rangle$. The second argument is stored in the register pair $\langle Ab, Bb \rangle$. The temporary regis-

remainder is
position to
f register
ult bit is
s contin-

a+Bb
a
a+Bb
a
1+Bb
a
Ba+Bb
Ba
Bb

1. The *i*th
stored in
the opera-

sition to
ter con-
 $< r < 2$.
a one is
ne posi-

Z3. Fig.
) cycles

needed to compute the square root of a number. The argument for the operation is stored in the register pair <Af,Bf>. The register pair <Aa,Ba> is initialized to zeroes. The algorithm computes the square root of numbers with an even exponent. If the exponent is an odd number, the significand is shifted one place to the left, and the exponent is decremented by one. The final exponent (computed in Cycle 19) is half this initial exponent.

cycle	stage	exponent	significand
0	I,II,III		
1	IV,V	Aa:=Af	Ba:=Bf
	I,II,III	Ae:=Aa-Ab	if (Ba-Bb ≥ 0) then Be:=Ba-Bb, bt:=1 else Be:=Ba, bt:=0
2	IV,V	Aa:=Ae Ab:=0	Bf:=0 if (bt=1) then Bf[0]:=1 Ba:=2×Be
	I,II,III	Ae:=Aa+Ab	if (Ba-Bb ≥ 0) then Be:=Ba-Bb, bt:=1 else Be:=Ba, bt:=0
3	IV,V	Aa:=Ae	if (bt=1) then Bf[-1]:=1 Ba:=2×Be
	I,II,III	Ae:=Aa+Ab	if (Ba-Bb ≥ 0) then Be:=Ba-Bb, bt:=1 else Be:=Ba, bt:=0
⋮	⋮	⋮	⋮
<i>i</i>	IV,V	Aa:=Ae	if (bt=1) then Bf[2-i]:=1 Ba:=2×Be
	I,II,III	Ae:=Aa+Ab	if (Ba-Bb ≥ 0) then Be:=Ba-Bb, bt:=1 else Be:=Ba, bt:=0
⋮	⋮	⋮	⋮
16	IV,V	Aa:=Ae	if (bt=1) then Bf[-14]:=1 Ba:=2×Be
	I,II,III	Ae:=Aa+Ab	if (Ba-Bb ≥ 0) then Be:=Ba-Bb, bt:=1 else Be:=Ba, bt:=0
17	IV,V	if (Bf[0] = 0) then Ab:=-1	Ba:=Bf, Bb:=0
	I,II,III	Ae:=Aa+Ab	Be:=Ba+Bb
18	IV,V	Af:=Ae	if (Bf[0]=0) then Bf:=2×Be else Bf:=Be

Fig. 11. The 18 cycles needed for the division algorithm. The *i*th bit of register Bf is denoted by Bf[*i*]. The dividend is stored in the register pair <Af,Bf> and the divisor in <Ab,Bb> before the operation is started.

The main idea of the classical algorithm is to reduce the square root operation to a division. If we want to compute the square root of *x*, we need a number *Q* such that $x/Q = Q$. The result *Q* is built sequentially by setting the *i*th bit to one and then testing whether the condition $x > Q^2$ still holds. If this is not the case, the *i*th bit must be set to zero.

Assume that we have already computed from bit zero to bit $-i+1$ of the final result. Denote by Q_{-i+1} the significand

$$Q_{-i+1} = Bf[0] \times 2^0 + Bf[-1] \times 2^{-1} + \dots + Bf[-i+1] \times 2^{-i+1}.$$

Bit $-i$ is then set to q_{-i} and it must hold that

$$x \geq Q_{-i}^2 = (Q_{-i+1} + q_{-i} 2^{-i})^2$$

This is true if

$$x - Q_{-i}^2 = (x - Q_{-i+1}^2) - 2^{-i} q_{-i} (2Q_{-i+1} + 2^{-i} q_{-i}) \geq 0$$

Define t_{-i} using the expression

$$2^{-i} t_{-i} = x - Q_{-i}^2 = (x - Q_{-i+1}^2) - 2^{-i} q_{-i} (2Q_{-i+1} + 2^{-i} q_{-i})$$

This can be written as

$$2^{-i} t_{-i} = t_{-i+1} 2^{-i+1} - 2^{-i} q_{-i} (2Q_{-i+1} + 2^{-i} q_{-i})$$

where we have used the recursive definition

$$2^{-i+1} t_{-i+1} = (x - Q_{-i+1}^2).$$

Simplifying the last expression, we finally get:

$$t_{-i} = 2t_{-i+1} - q_{-i} (2Q_{-i+1} + 2^{-i} q_{-i})$$

If t_{-i} is positive for $q_{-i} = 1$, we set bit $-i$ of the final result to one, i.e., Bf[$-i$]=1. If t_{-i} is negative, we set Bf[$-i$]=0. The recursive computation is started with $t_0 = x$. Q_{-i+1} represents at each step the partial result contained in register Bf. Bit $-i$ is tentatively set to one, and the sign of t_{-i} is tested. The basic loop of the square root algorithm for bit $-i$ has the following form:

```

Ba := 2 × Be
Bb := 2 × Bf
Bb[-i] := 1
if (Ba - Bb ≥ 0) then Be := Ba - Bb, Bf[-i] := 1
else Be := Ba, Bf[-i] := 0

```

The square root algorithm is the jewel in the crown of the Z3.

All bits of register Bf are used for the computation of the square root. If the original number lies within bounds, the result is also within bounds.

Read and Display Instructions

The two most complex instructions of the Z3 are those related to the input and output of decimal numbers. A decimal number of four digits entered through the keyboard is first converted into a binary integer. This is done by reading each digit sequentially, transforming it into a binary number, and storing it in the bits Ba[-10], Ba[-11], Ba[-12], and Ba[-13] of register Ba. The number in register Ba is multiplied by 10, and the procedure is repeated for the other digits. After four iterations, the decimal input has been transformed to a binary number (the exponent of the binary representation is formed indirectly via shifts resulting from multiplication by 10). The difficult part is handling the exponent. If the exponent *e* is positive, the significand has to be multiplied *e* times with 10. If it is negative, it must be multiplied $|e|$ times with 0.1. Multiplying with 10 is relatively easy: The significand in Be can be shifted one bit to the left and then stored in Ba (i.e., Ba := 2 × Be). At the same time, Be can be shifted three places to the left and can be stored in Bb (i.e., Bb := 8 × Be). The addition of Ba and Bb then provides the desired result: the multiplication

times as
until the
between zero
that can be
number is
continues
number in
lications

Fig. 13.
er. Notice
transferred
which are

seful con-
which are
between
algorithm.
nt (A⁺)

zero, the
When it is
s relay is
a₁, b₁, and
needed.

of result
② ③ ④

(result)



nor
mal

The box labeled "zero, infinite" below Ae represents the circuits for exception handling. They snoop permanently on the data bus (results of operations and data from memory) and raise the corresponding exception flags when needed. The shifter below Be is used to displace the significand one bit to the right. This provides the normalization needed for the significand whenever $Be \geq 2$.

Fp and Fq are the relays that control the number and direction of one-bit shifts in the shifter below the relay boxes Fc and Fa. Fh, Fi, Fk, Fl, and Fm have the same function in relation to the other shifter. Using these five bits, the numbers between -16 and 15 can be represented, and this is also the range of the second shifter. When such a shift is performed, the number represented by the relays Fh to Fm is transferred through the relay box Bn to register Ab in order to modify the exponent of the result. If the number is shifted 10 positions to the left, then +10 is subtracted from the exponent of the result. Such drastic shifts are needed mostly after subtractions.

Look again at the diagram of the Z3. Everything makes sense now and looks as conventional as any modern small floating-point processor. It is indeed amazing how Zuse was able to find the adequate architecture right from the beginning. The Z3 processor employs just 600 relays; the memory needs three times as much. By having to optimize the design and by having to save hardware everywhere, Zuse was *forced* to think and rethink the logical structure of his machine. He was not allowed the luxury of the almost unlimited funding allocated by the U.S. military for the development of the ENIAC or by IBM for the Mark I. He was all alone. While this may have worked to his advantage from the conceptual side, it may also have worked to his disadvantage, considering the negligible impact that the Z1 and Z3 had on the emerging U.S. computer industry after World War II.¹³

The Invention of the Computer

The main defect of the Z3 was the absence of a conditional branch in the instruction set. When the program is stored on punched tape, a possible fix is to include multiple tapes and a mechanism to switch between them (as was done with the Harvard Mark I). Another possibility is having a "program counter," so that the tape can be advanced or rewound on demand.

Sometimes the dividing line between calculating machines and universal computers is drawn by differentiating between machines with externally or internally stored programs. I have argued elsewhere¹⁰ that this is not a valid criterion. An external program can work as an interpreter of numerical data. The external program becomes a fixed part of the processor, and the data become the program, much in the same way as a universal Turing machine works as an interpreter. I have argued that what is needed for universal computation is a minimal instruction set and indirect addressing.¹¹ Indirect addressing can be simulated by writing self-modifying programs, so that the instruction set becomes the defining criterion. A machine with enough addressable memory and an accumulator and that is capable of executing the instructions CLR (clear), INC (increment), LOAD, STORE, and BZ (branch if zero) is a universal computer. In this sense, the Z1 was *not* a fully fledged computer, but neither were any of the other early machines. The ABC was a special-purpose machine for solving sets of linear equations by Gaussian elimination; the Harvard Mark I lacked conditional

branching, although it featured loops; the ENIAC was not even programmable through software—the building blocks had to be hardwired in dataflow fashion. Conditional branching was available in the ENIAC only in a limited way, and self-modifying programs were, of course, out of the question.

TABLE 2
COMPARISON OF ARCHITECTURAL FEATURES

Machine	memory and CPU separated?	conditional branching?	soft or hard programming	self-modifying programs?	indirect addressing?
Zuse's Z1	✓	x	soft	x	x
Atanasoff's	✓	x	hard	x	x
H-Mark I	x	x	soft	x	x
ENIAC	x	partially	hard	x	x
M-Mark I	✓	✓	soft	✓	x

TABLE 3
SOME ADDITIONAL ARCHITECTURAL FEATURES

Machine	internal coding	fixed-point or floating-point?	bit-sequential arithmetic?	architecture	technology
Zuse's Z1	binary	floating	no	sequential	mechanical
Atanasoff's	binary	fixed-point	yes	vectorized	electronic
H-Mark I	decimal	fixed-point	no	parallel	electromechanical
ENIAC	decimal	fixed-point	no	dataflow	electronic
M-Mark I	binary	fixed-point	yes	sequential	electronic

Tables 2 and 3 show the most relevant information about the early computing machines mentioned earlier. As should be clear from the tables, none of the early computing machines fulfills all the necessary requirements for a universal computer. I also include the Mark I machine built in Manchester from 1946 to 1948, because as far as I know this was the first machine to fit my definition of a universal computer. The Mark I was built under the direction of F.C. Williams and T. Kilburn. This machine stored its program in random-access digital memory implemented with CRT tubes. All necessary instruction primitives were available (in modified form), and although it lacked indirect addressing, self-modifying programs could be written. The first program ran in June 1948 and calculated the highest proper factor of a large number.⁷ In September 1948, Turing was appointed as a reader in mathematics in Manchester and wrote some programs for the first universal computer in the world. His vision of universal computation published in 1936, the same year in which the storage unit of the Z1 was completed, had at last become a reality. Tables 2 and 3 are emphatic: The invention of the computer was a collective achievement encompassing two continents and 12 years.

Acknowledgments

Deciphering the sketchy documentation available was possible only with the collaboration of several of my students at the Universities of Halle and Berlin. I thank Alexander Thurm and Axel Bauer, who implemented a gate-level simulation of the Z3 processor. We became aware of synchronization problems when the simulation refused to run. I also thank Franz Konieczny, Reimund Spitzer, and Roland Schultes, who wrote part of a stand-alone simulation of the processor in C. We started working on the Z3 with the help of Konrad Zuse, who gladly answered our questions. It was amazing to see how, after almost 60 years, the whole design of the Z3 was still in his head. Unfortunately, Zuse died in December 1995 before this description of his work was ready. This paper is dedicated to his memory.