R. Kirchner and U. Kulisch

Fachbereich Informatik, Universität Kaiserlautern Fakultät für Mathematik, Universität Karlsruhe West Germany

<u>Abstract</u>: In electronic computers the elementary arithmetic operations are these days generally approximated by floating-point operations of highest accuracy. Vector processors and parallel computers often provide additional operations like "multiply and add", "accumulate" or "multiply and accumulate". Also these operations shall always deliver the correct answer whatever the data are. The user should not be oblighed to execute an error analysis for operations predefined by the manufacturer.

In the first part of this paper we discuss circuits which allow a fast and correct computation of sums and scalar products making use of a matrix shaped arrangement of adders and pipeline technology. In the second part a variant is discussed which permits a drastic reduction in the number of adders required. The methods discussed in this paper can also be used to build a fast arithmetic unit for micro computers in VLSI-technology.

1. <u>Introduction</u>

Modern computers of highest performance, the so-called vectorprocessors or supercomputers, are gaining considerably in importance in research and development. They serve for simulation of processes which cannot be measured at all or only with great effort, for solving large engineering design problems or for evaluation of large sets of measured data and for many other applications. It is commonly assumed that these computers open a new dimension for scientific computation. In sharp contrast to this is the fact that the arithmetic implemented on supercomputers differs only marginally from that of their slower predecessors, although results are much more sensitive to rounding errors, numerical instabilities, etc. due to the huge number of operations executed.

Research in numerical mathematics has shown that, with a more comprehensive and optimal vector arithmetic, reliable results can be more easily obtained when dealing with extensive and huge problems. Computers with this kind of arithmec have proved the significance of this development in many successful applications.

Until now, it has been assumed that an optimal vector arithmetic could not be implemented on supercomputers. The users, therefore, had to choose between either lengthy computation times and accurate results on general purpose computers or comparatively short computation times and possibly wrong results obtained on supercomputrs.

It was assumed, in particular, that correct computation of continued sums and scalar products, which are necessary for vector arithmetic, could not be implemented on supercomputers with pipeline processing. Well known circuits, which solve this problem, require several machine cycles for carrying out a single addition whereas a computer of highest performance with traditional arithmetic

carries out one addition in each cycle¹. This paper describes various circuits for the optimal computation of sums and scalar products at the speed of supercomputers. There is, in principle, no longer any reason to continue to accept inaccurate sums or scalar products by not using optimal vector arithmetic on vectorprocessors and supercomputers. The additional costs compared with the cost of the complete system are justified in any case. It takes the burden of an error analysis from the user.

The first electronic computers were developed in the middle of this century. Before then, highly sophisticated electromechanical computing devices were used. Several very interesting techniques provided the four basic operations of addition, subtraction, multiplication, and division. Many of these calculators were able to perform an additional operation which could be called "accumulating addition/subtraction" or continued summation. The machine was equipped with an input register of about 10 to 13 digits. Compared to that, the result register was much longer and had perhaps 30 digits. It was situated on a sled which could be shifted back and forth relatively to the input register. This allowed an accumulation of a large number of summands into different positions of the result register. There was no rounding executed after each addition. As long as no overflow occurred, this accumulating addition was error free. Addition was associative, the result being independent of the order in which the summands were added.

This accumulating addition without intermediate roundings was never implemented on electronic com-

¹By a cycle time or a machine cycle we understand the time which the system needs to deliver a summand or a product, in case of a scalar product computation, to the addition pipeline.

puters. Only recently, several /370 compatible systems have appeared which simulate this process on general purpose machines by accumulating into an area in main memory, which is kept in the cache memory for enhanced performance. [5], [6]. This allows the elimination of a large number of roundings and contributes essentially to the stability of the computational process. This paper desribes circuits for an implementation of the accumulating addition on very fast computers making use of pipelining and other techniques.

The first <u>electronic computers</u> executed their calculations in fixed-point arithmetic. Fixed-point addition and subtraction is error free. Even very long sums can be accumulated with only one final rounding in fixed-point arithmetic, if a carry counter is provided which gathers all intermediate positive or negative overflows or carries. At the very end of the summation a normalization and rounding is executed. Thus accumulation of fixed point numbers is associative again. The result is correct to one unit in the last figure and it is independent of the order in which the summands are added. Fixed-point arithmetic, however, imposed a scaling requirement. Problems needed to be preprocessed by the user so that they could be accommodated by the fixed-point number representation. With the increasing speed of computers, problems that could be solved became larger and larger. The necessary pre-processing soon became an enormous burden.

The introduction of floating-point representation in computation largely eliminated this burden. A scaling factor is appended to each number in floating-point representation. The arithmetic itself takes care of the scaling. Multiplication and division require an addition, respectively subtraction, of the exponents which may result in a large change in the value of the exponent. But multiplication and division are relatively stable operations in floating-point arithmetic. Addition and subtraction, in contrast, are troublesome in floating-point.

As an example let us consider the two floatingpoint vectors

$$\mathbf{x} = \begin{bmatrix} 10^{20} \\ 1223 \\ 10^{24} \\ 10^{18} \\ 3 \\ -10^{21} \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} 10^{30} \\ 2 \\ -10^{26} \\ 10^{22} \\ 2111 \\ 10^{19} \end{bmatrix}$$

A computation of the inner or scalar product of these two vectors gives

 $x.y = 10^{50} + 2,446 - 10^{50} + 10^{40} + 6,333 - 10^{40} = 8,779$

Most digital computers will return zero as the answer although the exponents of the data vary only within 5 % or less of the exponent range of large systems. This error occurs because the floating-point arithmetic in these computers is unable to cope with the large digit range required for this calculation.

Floating-point representation and arithmetic in computers was introduced in the middle of this

centry. Computers then were relatively slow, being able to execute only about 100 floating-point operations in a second. The fastest computers today are able to execute billions of floating-point operations in a second. This is a gigantic gain in speed by a factor of 10^7 over the electronic computers of the early fifties. Of course, the problems that can be dealt with, have become larger and larger. The question is whether floating-point representation and arithmetic which already fails in simple calculations, as illustrated above, are still adequate to be used in computers of such gigantic speed for huge problems.

We think that the set of floating-point operations should be extended by a fifth operation, the "accumulating addition/subtraction" without intermediate rounding, an operation which was already available on many electromechanical calculators. It is the purpose of this paper to show that this additional operation can be executed with extreme speed. We realize this operation by adding the floating-point summands into a fixed-point number over the full floating-point range. Thus "accumulating addition/subtraction" is error free. Even very long chains of additions/subtractions can be executed with only a single rounding at the very end of the summation. Such "Accumulating addition/ subtraction" is associative. The result is independent of the order in which the summands are added.

With the fifth operation "accumulating addition/subtraction", we combine the advantages of fixedpoint arithmetic - error free addition and subtraction even for very long sums - with the advantages of floating-point arithmetic - no scaling requirements.

2. The State of the Art

A normalized floating-point number z (in sign-magnitude representation) is a real number of the form

$$z = * m \cdot b^{e}$$

Here $* \in \{+,-\}$ denotes the sign (sign(z)), m the mantissa (mant(z)), b the base of the number system and e the exponent (exp(z)). b is an integer number with b > 1. The exponent is an integer and lies between two integers el \leq c2. In general, el \leq 0 and e2 > 0. m is the mantissa. It is of the form

$$m = \sum_{i=1}^{l} z[i] \cdot b^{-i}.$$

Here, the z[i] denote the digits of the mantissa; z [i] ϵ {0,1,...,b-1} for all i = 1(1)n and z[1] \neq 0. 1 is the length of the mantissa. It denotes the number of mantissa digits carried along. The set of normalized floating-point numbers does not contain the number 0. In order to obtain a unique definition of 0 one can additionally define: sign(0) = +, mant(0) = .000 ... 0 (1 zeros after the point) and exp(0) = el. This kind of floatingpoint system depends on four constants b,l,el and e2. We denote it with S = S(b,l,el,e2). Let

$$u = (u_{1}) = \begin{bmatrix} u_{1} \\ u_{2} \\ \vdots \\ \vdots \\ u_{n} \end{bmatrix} \qquad v = (v_{1}) = \begin{bmatrix} v_{1} \\ v_{2} \\ \vdots \\ v_{n} \end{bmatrix}$$

be two vectors, the components of which are normalized floating-point numbers, i.e. u_i , $v_i \in S$

forall i = i(1)n. The theory of computer arithmetic[1], [2], [3] demands that scalar products of two floating-point vectors u and v be computed with maximum accuracy by the computer for each relevant, finite n and different roundings. By doing so, millions of roundings can be eliminated in complicated calculations. This contributes essentially to the stability of the computational process and enlarges the reliability and accuracy of computed results. Furthermore, defect correction then becomes an effective mathematical instrument.

This requires, for example, the execution of the following formulae by the computer:

$$u \odot v = O\left(\sum_{i=1}^{n} u_{i} * v_{i}\right)$$
$$u \boxdot v = \Box\left(\sum_{i=1}^{n} u_{i} * v_{i}\right)$$
$$u \bigtriangledown v = \nabla\left(\sum_{i=1}^{n} u_{i} * v_{i}\right)$$
$$(1)$$
$$u \bigtriangleup v = \Delta\left(\sum_{i=1}^{n} u_{i} * v_{i}\right)$$

The multiplication- and addition-signs on the right side denote the correct multiplication and addition for real numbers. O, \Box , ∇ , Δ are rounding symbols. O denotes a rounding to the nearest floating-point number. \Box denotes the rounding towards zero, ∇ denotes the monotone downwardly directed rounding and Δ denotes the monotone upwardly directed rounding.

For an execution of formulae (I) first the products $u_i \neq v_i$ have to be correctly calculated by the computer. This leads to a mantissa of 21 digits and an exponent which lies in the range of 2e1-1 $\leq \leq 2e2$. So the computation of scalar products is reduced to the evaluation of sums of the following form:

$$(\sum_{i=1}^{n} w_i), n \in \mathbb{N}$$
 (II)

Here the w_i are floating-point numbers of double length $w_i cS(b, 21, 2e1-1, 2e2)$, for all i = 1(1)n. \diamond denotes a general rounding symbol, $\diamond \in \{O, \Box, \nabla, \Delta\}$. Measures have to be taken first to generate and represent the summands w_i correctly in the computer. In case of scalar products this can be done by very fast and well-known circuits.

For traditional general purpose computers there

are several ways to correctly compute (I) and (II) mentioned in the literature. It is the intention of this paper to describe circuits for high speed computation of (I) and (II) on vector computers by means of pipeline techniques. These circuits have to accept and process one summand from (I) resp. (II) during each machine cycle. To assist in the understanding of the following material, we first refer to one of the possibilities mentioned in [4]:

We consider a register of $L = k + 2e^2 + 2l + 2|e_1|$ digits of base b, which should be placed in the arithmetic unit (Figure 1).

k	2e2	21	2 e1
		Figure 1	

We divide this register into segments of length 1 (Fig. 2):

k			1		
		Figu	ire 2		

The summands in (I) and (II) are of length 21. They fit therefore, digitwise into a subrange of length 31 of this storage. This part of the register, which is determined by the exponent of the summand, is selected and loaded into an accumulator of length 31. The summand is loaded into a shiftregister of the same length, being correctly positioned according to the exponent, and then added into the accumulator (Figure 3).



The addition may produce a carry. In order to catch this carry, a few more digits than the three words of length 1 can be read from the long register into the accumulator, which is extended to the left accordingly. If not all of these digits are b-1, the carry is caught by these additional digits. Since it is possible that all these additional digits are b-1, a loop has to be provided which then adds the carry to the following digits of the long register. This loop may possibly have to be activated several times.

The addition of the summands of (I) resp. (II) into the long register, Fig. 1 resp. Fig. 2, may still produce a carry on the very far left of the register. In order to catch such carries the long register is extended on the left by a few more (k) digits of base b (Fig. 1). Then, any sum (I) or (II) of n summands can be added without loss of information into the long register of length L. b^{k} carries may occur and can be processed without loss of information.

Here we conclude our description of one possibility to solve the problems (I) and (II). See [4]. What we just described belongs to the state of the art.

9

3. Fast Computation of Sums and Scalar Products

The method described above is not suited for the computation of (I) resp. (II) on vector processors or supercomputers. The process of reading, shifting, carry handling, possibly by a loop, and writing back is certainly too slow to be executed in one cycle time of only a few nsecs of these computers. A solution of the problem by a very long adder is also very costly and probably too slow.

We therefore discuss here a variant of the possibilities mentioned above which makes processing of a summand of (I) resp. (II) possible within a very short cycle time. In comparison to general purpose computers, vector processors and supercomputers achieve their high speed of computation by means of pipeline technology whereby during each machine cycle a result is obtained. If scalar products and sums are to be computed with high speed on vector processors or supercomputers, one has to develop circuits which accept and process one summand (resp. a product) per machine cycle. This is only possible if the addition is done by means of pipeline technology. This paper describes various circuits which allow this.

At first the most important issues and ideas of the circuitry are presented in the text referring to Figures 4 to 15. These Figures contain some more details which are not essential for a first understanding of the principles. These details are presented later in chapter 4 "Additional Remarks concerning the Figures".

The circuit described below consists of a shifter which is followed by a pipelined adder called summing matrix (Figure 4). The shifting device may be realized by standard technology and belongs to the state of the art.

The adder consists of registers of a total length of S \geq L. Here L denotes the length of the long

register as outlined above² (Figure 1). The register length S is divided into r identical parts which are arranged as rows one below the other (Figure 5). r denotes the number of rows. All rows are of the same length. Each of these rows is divided into $c \ge 1$ independent adders A (see Figure 6). Thus the whole summing device consists of $r \cdot$ c independent adders. Each of these adders A has a width of a digits. Between two of these independent adders, carry handling must be possible. Also between the last adder of a row and the first one of the next row a carry handling must be possible. The complete summing device which we call the summing matrix SM, has a width of $S = a \cdot c \cdot r di$ gits of base b. c denotes the number of columns of the summing matrix. It must be $S \ge L = k + 2e2 + e^{-k}$ 21 + 2 |e1 | (Figures 5, 6).

The summing matrix contains $c \cdot r$ independent adders A. Each of these adders must be able to add a digits of base b in parallel within one machine cycle, and to register a carry which possibly may occur. Since each row of the summing matrix consists of c identical adders, $h:= c \cdot a$ digits can be added in each row of the summing matrix. Each of the r rows of the summing matrix SM must be at least as long as the mantissa length of the summands which are to be added. Each digit of the summing matrix is characterized by a certain exponent corresponding to the digit's position. The upper right part of the summing matrix carries the least significant digit, the lower left part of the summing matrix carries the most significant digit of the full summing device (Figure 5, Figure 6).

Each summand resp. each product of (I) resp. (II) must now be added into the summing matrix at the proper position according to its exponent. The row selection is obtained by the more significant bits of the exponent (exp div h)³ and the selection of the columns is obtained by the less significant bits of the exponent (exp mod h)⁴. This complies roughly with the selection of the adding position in two steps by the process described in Fig. 3.

The incoming summands resp. products are now first shifted in the shifting unit (barrel shifter, cross bar switch) into the correct position according to their exponents. The shift is executed as a ringshift. This means that the part of the summand which hangs over the right end is reinserted at the left end of the shiftregister (Figure 6 upper part, summands 2 and 3, Figure 8). The summand is distributed onto the c independent parts of width a of the shiftregister. Each part receives an exponent identification according to a specific digit in it, e.g. the least significant one (Figures 5, 6 and 10). The individual adders A also carry an exponent identification. The shifted and expanded summand now drops into the top row of the summing matrix and thereafter proceeds row by row through the summing matrix, moving ahead one row in each machine cycle. The addition is executed as soon as the exponent identification of a transfer register in the summing matrix coincides with the exponent identification part of the summand.

A summand, which arrives at the summing unit, can remain connected after shifting to the correct position within the shifting unit. In this case, the addition is executed in only one row of the summing matrix. The shift procedure, however, can also cause an overhanging at the right end of the row. The overhanging part then is reinserted by a ringshift at the left end of the shifting unit (see Figures 6 and 8). In this case, the addition of both parts of the summand is then executed in neighbouring rows of the summing matrix. If the most significant part of the summand, which was situated at the right end of the shifter, is added in row y then the addition of the least significant part, which was situated at the left end of the shifter, is added in row y - 1. This means the next less significant row (see Figure 9).

It is, however, not at all necessary that each

²or a part of it. A reduction of the length S is discussed below.

 $[\]frac{4}{mod}$ denotes the remainder of integer division, i.e. 24 mod 10 = 4.

transfer unit carries a complete exponent identification. It is sufficient to identify the row by the exponent part exp \underline{diy} h of the summands in the shifter and to use it for selection of row y. The distinction whether the addition has to be executed in row y or in row y - 1 is made by a bit connected with each transfer register or by a suitable column signal which distinguishes the transfer registers of a row. (The principle is illustrated by the diagrams shown in Figures 11 and 12).

The addition may cause carries between the independent adders A. Carry registers between the independent adders absorb these carries. In the next machine cycle these carries are added into the next more significant adder A, possibly together with another summand. In this way, during each machine cycle one summand can be fed into the summing matrix, although the carry handling of on summand may take several machine cycles. The method displayed in the Figures shows one of diverse possibilities to handle the carries. There may be carry presencing or look-ahead or other techniques applied to speed up the carry processing within one row. In any way, the summing matrix allows the carry processing to be executed independently of the summations and in parallel with the processing that has to be done at all, e.g. adding further summands or reading out the result.

In principle, the summing matrix can only process positive summands. Negative summands or positive subtrahends are therefore marked and at the proper place not added but subtracted. Here negative carries instead of positive carries may occur.Similar to positive carries they have to be processed possibly over several machine cycles. In other words: The independent adders A must be able to carry out additions as well as subtractions and to process positive and negative carries in both cases (Figure 6, 12).

The design of the complete summing device containing the summing matrix SM described herewith can depend on the technology used. We have mentioned already that the width a of the individual adders A has to be chosen in such a way that an addition over the complete width can be executed within one machine cycle. Each row of the summing matrix must be at least as wide as the individual summands. The shorter the rows are, the faster the summands can be shifted into the right position. On the other hand, shortening the width of the rows of the summing matrix increases the number of rows and with it, the number of pipeline steps for the complete summation process.

After input of the last summand the rows can be read starting with the least significant row, provided the row in question does not require any carry handling. In this case the carries first have to be removed. The readout process can use the same data path by which the summands pass through the matrix. Thus the result rows follow the last summand on its way through the transfer registers. During the readout process additions and carry handling in the more significant rows may still be executed. Simultaneously with the readout process the rounding to the required floating-point format can be executed. The result can also be stored as an intermediate long variable for further processing. Several rounding possibilities can be carried out simultaneously as mentioned in [4]. During the readout process the computation of a new scalar product resp. a new sum can be started.

The width a of the independent adders A depends on the technology used and on the cycle time of the system. The width should be as large as possible. But on the other hand, it must permit the addition over the a digits in one machine cycle. (In the case of a scalar product, a machine cycle is the time in which the system delivers a product).

Depending on the technology there are several possibilities of transportation of the summands to one of the r rows of the summing matrix SN.

The method described above is based on the idea that each of the independent adders A is supplemented by a transfer register of the same width (plus tag-register for exponent identification and +/- control). During each machine cycle, each transfer register can pass on its contents to the transfer register in the corresponding position in the next row and receive a digit sequence from the transfer register in the corresponding position in the previous row. Attached to the transfer registers is the tag-register for exponent identification (Figure 5 and Figure 6). The contents of this register are always compared with the exponent identification of the corresponding adder. In case of coincidence, the addition resp. subtraction is activated (Figures 5, 6 and 12).

Alternatives to this procedure are also possible.

- 1. One of these alternatives could be to transfer the summand in one machine cycle directly into the appropriate row of transfer registers of the summing matrix as determined by the exponent. During the following machine cycle, the addition is executed. Simultaneously, a new summand can be transferred to the same, or another row, so that an addition in each machine cycle is carried out.
- 2. The procedure is similar to 1. The intermediate storage of the summands in transfer registers, however, is not necessary if it is possible to execute the transfer- and addition-process in one machine cycle. In this case, no transfer registers are necessary. The output of the result then also takes place directly.
- 3. The transfer of the summands to the target row can be carried out not only sequentially and directly but also with several intermediate steps, for example, by binary selection.

Each one of these alternatives also allows a direct and therefore faster readout of the result without dropping step by step through the transfer registers.

To each independent adder A of length a belongs a transfer register TR which is basically of the same length. The number of adders A resp. transfer registers TR in a row is chosen in such a way that

the mantissa length \overline{m} of the summands plus the length of the transfer registers t (=a) becomes less or equal to the length of the row ($\overline{m} + a \leq h$ = c • a). In this way, an overlapping of the less significant part of the mantissa with its most significant part in one transfer register is avoided. For typical floating-point formats this concition may result in long rows of the summing matrix or in short widths a of the adders resp. transfer registers. The former case causes lengthy shifts while the latter case causes more carries (Figure 6 upper part and Figure 8).

This disadvantage can be avoided by providing several (\geq 2) partial transfer registers for each adder of length a. Each partial transfer register TR of length t \leq a carries its own exponent identification. Finally, the length t of the transfer registers can be chosen independently of the length a of the adders A. Both only need to be integer divisors of the row length of the summing matrix $h = a \cdot c = t \cdot n$ (see Figures 13, 14 and 15).

Figures 6 and 13 show, in particular, that the summing matrix has a very systematic structure and that it can be realized by a few, very simple building blocks. It is suitable, therefore, for realization in various technologies.

Based on the same principle also, summands which consist of products of three and more factors can be added correctly.

If the summing matrix is to be realized in VLSI-technology it may happen that the complet summing matrix does not fit on a single chip. One should then try to develop components for the columns of the summing matrix since the number of connections (pins) between adjacent columns is much smaller than between neighbouring rows.

The following remarks and Figures 4 to 15 provide a more detailed description of the structure of the su ming matrix and its functioning.

4. Additional Remarks concerning the Figures

The following abbreviations are used in the Figures: A

- Adder
- Accumulator Register AC
- CY Carry
- Ε Tag-Register for Exponent Identification
- LSB Least Significant Bit
- · MSB Nost Significant Bit
- Summing Matrix SM
- SR Shifter
- TR **Transfer Register**

Figure 4 shows a structure diagram of the complete summing circuitry and illustrates the interaction of different parts of the whole circuitry, such as: separation of the summands into sign, exponent and mantissa, shifting unit, summing matrix, controller and rounding unit.

Figure 5: As mentioned in the text, we assume that $S \ge L$. Figure 5 shows the case S > L. There, for both the first and last rows part of the row is covered by transfer registers only. For the whole summing matrix this means that transfer registers exist for S digits but adders for L digits only. L is chosen such that it is a multiple of a.

The dotted lines through the independent adders A indicate that the transfer wires bypass the adders. Above the transfer registers, the tag-register for exponent identification is indicated by a box. This register is part of the transfer register.

Figure 6 shows a block diagram of the summing matrix. It is based on a special data format which uses 4 bits to describe one digit of base b.

Width of AC: a = 4 bytes = 32 bits Number of adders in one row c = 5Number of rows in SM r = 8k = 20 carry digits, 1 = 14 digits in the mantissa e1 = -64 and e2 = 64. Users of /370 compatible systems will recognize this data format as their double precision format. $L = 20 + 2 \cdot 64 + 2 \cdot 14 + 2 \cdot 64 = 304$ digits of 4 bits = 152 bytes. Width of the complete summing matrix $S = a \cdot c \cdot r = 4 \cdot 5 \cdot 8$ by tes = 160 by tes $\geq L =$ 152 bytes. In this example the width t of the transfer registers equals the width of the adders: t = a = 4bytes. The upper part of the Figure shows several positions of summands.

Figure 7 defines the exponent coordinates x and y of the digits in the summing matrix (x horizontal, y vertical). These coordinates are obtained according to the following formulae:

- eo denotes the reference point, the digit with the least exponent in the matrix (at the upper right end).
- e₁ denotes the least significant digit of the adder.
- e denotes the most significant digit of theadder. If the first and the last row of the complete matrix contains adders over the full width then $e_1 = e_0$ and $e_m = e_0 + r \cdot h - 1$.
- denotes the exponent of a digit to be added.
- e^{-e}o denotes the distance to the least significant end of the matrix.
- (e-e_o) <u>div</u> h is the row coordinate in which y =
- the digit with the exponent e is added. $(e-e_0) \mod h$ indicates the distance to the x ≃

least significant end of row y.

Figures 8 and 9 describe the task of the shift unit and its relation to the generation of the exponent identification which will be transferred into the summing matrix with the mantissa.

The task of the shift unit is:

- 1. adjust mantissa to the correct position for its addition, if necessary by a ring shift.
- fill the remaining positions of the transfer registers resp. the row with zeros.

Figure 8 shows the shifted mantissa in both possible cases.

Figure 9 describes the shift process. Two cases are to be distinguished:

- 1. $x = (e-e_0) \mod h \ge \overline{m}$: no overhanging,
 - the whole mantissa is added in one TOW.
- x < m : 2. overhanging. mantissa is added the in two

successive rows. Part m_1 remains within the width of the row. The overhanging part m_2 is reinserted at the left of the row. Both parts are furnished with a corresponding exponent identification. Part m_2 will be selected for addition in row y-1 whereas part m_1 will be added in row y.

The shifted and expanded mantissa row drops row by row through the matrix as a transfer row. Before that, each transfer section is characterized by its exponent which carries the information where the addition has to be executed.

Figure 10 shows the exponent identification of the sections of the transfer rows. Each row of the transfer matrix consists of n transfer sections of length t. Figure 10 defines the exponent identification t_e (transfer exponent) of these transfer sections of the matrix. If e_t denotes the exponent

of the e.g. least significant digit of a transfer section then this transfer section can be characterized by the exponent identification t_e

with $t_e = (e_t - e_0) \underline{div} t$.

Before a summand enters the matrix, each transfer section of the summand receives an exponent identification. During the passage through the matrix, this exponent identification is then compared with $t_{\rm e}.$ Equality triggers the addition. The lower part

of Figure 10 shows how transfer sections of the summand get their exponent identification.

A mantissa with the exponent e (= exponent of its most significant digit), receives the exponent identification (e - e₀) <u>div</u> t = e_m in the most significant transfer section, and exponent identification e_m - 1, e_m - 2, etc. in the less significant transfer sections.

Figure 10 shows in the lower part the two typical cases. (Addition of the complete summand in one row resp. in two consecutive rows).

<u>Figure 11</u> explains the simplified adder selection by row identification y_i . This row identification

is transfered through the matrix with the transferrow. The addition is triggered off as soon as the row identification and the row index coincide. The row selection switch RS generates two selection signals which activate the adders of the row in question (see Figure 12, too). An activating signal is sent via the wire "z-selection" if the row identification equals the row index. An activating signal is sent via the wire "z-1-selection" if y - 1 equals the row index.

Then the transfer sections only carry the information (z-1,z)-summation.

Since the transfer rows may only contain positive values the information addition or subtraction is additionally transferred.

Thus the controller contains transfer registers with specific information for each row which leads about to the structure shown in Figure 11.

Figure 12 shows a block diagram for an adder cell. For simplicity the case t = a is selected. The cell contains centrally an "adder/subtractor" and a "partial accumulator section". The right upper corner shows the corresponding transfer register with wires from the next less significant row and to the next more significant row. Additionally, the transfer register contains a tag register for "z/z-1" identification which indentifies through which selection wire the cell can be activated. The "adder/subtractor" receives the operands from the "partial accumulator section" and in case of selection from the transfer register. Zero is added if no selection takes place. In addition, the carry (positive or negative) arriving from the right is processed during each addi-tion/subtraction and, if necessary, a carry is passed on to the next adder cell on the left. This carry is temporarily stored in an auxiliary register. Figure 15 further shows a control wire which selects the operation (addition/subtraction) as well as a control wire for the read out process (at the bottom of the figure). All control wires traverse the whole row.

Figure 13 is very similar to Figure 6. It shows one row of the summing matrix, but with t \langle a. The Figure is based on the same data format as Figure 6, i.e.: one digit of basis b is described by 4 bits, k = 20 carry digits, l = 14 digits in the mantissa, e1 = -64 and e2 = 64. Furthermore: Width of AC: a = 4 bytes = 32 bits. Number of adders in one row c = 4. Number of rows in SM r = 10. $L = 20 + 2 \cdot 64 + 2 \cdot 14 + 2 \cdot 64 = 304$ digits per 4 bits = 152 bytes. Width of the complete summing matrix $S = a \cdot c \cdot r = 4 \cdot 4 \cdot 10$ by tes = 160 by tes $\geq L =$ 152 bytes. In this example the width of the transfer registers is smaller than the width a of the adders: t $=\frac{a}{5}=2$ bytes. This permits a smaller row width of only c = 4adders. The upper part of the Figure shows the position of a summand of $\overline{m} = 2 \cdot 1 = 14$ bytes at a critical position.

Figure 14 shows another case where the width of the adders differs from that of the transfer registers ($t \neq a$). In the Figure the transfer registers are shown without exponent identification. Dotted lines again indicate transfer wires which bypass the adder in question.

Figure 15 shows a section of a row of the summing matrix with $t \neq a$. Here the case 3t = 2a has been selected. It shows how digits of the same transfer register are distributed and added into neighbouring adders.

5. <u>Summation with only one Row of Adders</u>

We now discuss a further variant of the above circuitry for which adders exist only for one row of the summing matrix. The complete structure of this variant is similar to the one before (Figure 16). I.e. the complete circuitry consists of an input adjusting unit, the summing unit with the actual accumulator and a device for carry handling, result row filtering and rounding.

The complete fixed-point word, over which summation takes place, is divided into rows and columns, as before. The transfer width and the adder width, however, must now be identical. The width can be chosen according to the criteria as outlined above. The columns of the matrix shaped summing unit are now completely disconnected, i.e. no transmission of carries takes place between the individual columns of the matrix during the process of summation. The carries occurring during the summation are collected in carry counters and processed at the end of the summation process.

Figure 17 shows the circuit of a "column" of the matrix shaped summing unit. The full "long accumulator" is spread over the various columns of the summing unit. The part allotted to one column is called "accu-memory", see (1) in Figure 17.⁵

To each cell of the accu-memory belongs a carry counter. The collection of carry counters of a column is called "carry-memory", see (2) in Figure 17. In these cells of the carry-memory <u>all</u> carries emerging from the adder/subtractor are collected and incorporated in the result at the very end of the summing process. The individual cells of the carry-memory must be so wide that they can take a carry (positive or negative) from each summand. For a vector length of 128 one needs, for example, 7 bits plus a sign bit resp. an 8 bit number in twos'-complement.

In Figure 17, for example, the column width is 32 bits and the width of the individual carry-memory cells is 16 bits. This allows a correct computation of sums with less than or equal to 32 K summands. The exponent identification (in Figure 17) has a width of e bit; consequently the column has

 2° cells resp. the memory matrix 2° rows.

During the normal summation process the following happens:

- 1. The mantissa section MANT, sign sg. and exponent identification EPI reach the input register RI, (3).
- 2. In the next cycle
 - the memory is addressed through EPI and the accu-part as well as the carry part are transferred to the corresponding section of the register <u>b</u>efore the <u>sum-mation RBS</u>, (4);
 - the mantissa section, sg, and EPI are also transferred to the corresponding section of RBS, (5).

3. In the next cycle

- addition resp. subtraction according to sg is executed in the adder/subtracter (6). The result is transferred to the corresponding section of the register after the gummation RAS, (7). According to the carry, the carrypart is adjusted in (8) by +1, -1 or not at all and it is also transferred to RAS. (9);

- EPI is also transferred to RAS, (10).
- 4. In the next cycle
 - EPI of RAS addresses the memory, and the accu-part together with the carrypart are written back into the memory.

Since in each machine cycle a mantissa section is supplied, these phases must be pipelined. This means, in particular, that all phases need to be active simultaneously. It <u>must</u> be possible therefore, to read from the memory and to write into the same or another row of the memory during each machine cycle. This, however, is usual for register memories.

If in two consecutive cycles the same accu- and carry-memory cell is addressed, the previously described procedure may lead to a wrong result, since in the second cycle the result of the just started summing procedure should be read, which does not yet exist. We have a typical pipeline conflict. These difficulties can be overcome by duplicating the accu-carry-memory several times which, however, is very costly.

Therefore, we suggest an easier alternative. We suppose that during consecutive cycles mantissa sections with the same exponent identification arrive. We distinguish the following two cases:

- a) directly one after the other;
- b) with another exponent identification in between and both arbitrarily often and mixed.

We first deal with case a).

1) The registers EPI of RI and EPI of RBS contain the same exponent identification. The two are compared in (11) and in case of coincidence the read process from the memory to RBS is blocked off in part (13) of the selection unit (12). Instead, the result of the addition of the first of the two consecutive summands is directly transferred to RBS via (14) so that the second summand can immediately be added. Events of the summand can immediately be added.

Furthermore, (15) causes a dummy exponent to be read into EPI of RAS. So, if in the same cycle a further third value with the same exponent identification is transferred to RI the case EPI/RI = EPI/RBS = EPI/RAS is avoided. This case would cause a conflict in the selection unit (12).

Thus, consecutive summands with the same exponent identification can be added without memory involvement. The intermediate values may be written into the memory or discarded (storage blockade on). Only the last value must be written into the memory via RAS.

We now deal with case b).

b) Three values EPI_1 . EPI_2 . EPI_3 with $\text{EPI}_1 = \text{EPI}_3 \neq \text{EPI}_2$. In this case EPI/RI and EPI/RAS-contain the same exponent identification. The two registers are compared in (16). In the following cycle the contents of RAS is directly transferred to RBS through part (17) of the selection unit (12). The read process from the memory is again suppressed in (13). The intermediate value may be written into the memory. It can also be suppressed.

In this way, any consecutive mantissa sections can

⁵The numbers enclosed in round parentheses in the text indicate in the corresponding Figure that part of the circuitry which is marked with the same number.

be added and the carries collected in the carry counters.

We now consider the process of reading the result. The central read control produces continuous addresses so that the accu-memory is read from the least significant to the most significant row. This sequence is a must because of the necessary carry handling. The addresses reach the memory through the multiplexer (18).

Wires (19), (20) for transfer of the carries lead from column to column. The carry-parts of a column are fed to the next more significant column. There they are taken into the mantissa section of RBS. To get there the multiplexer (21) is switched over. The carry, which is stored in the twos'-complement for convenience, first has to be changed into sign-magnitude-representation and, if necessary, expanded in length (22). In the next cycle, the carry is added and together with a possible lbit-carry (positive or negative) transferred to the unit for preparing the result after temporary storage in RAS. The above mentioned carry can there be stored either in a part of the RAS-carry register or in a 2bit auxiliary register (23).

During the process of reading it is advisable to delete the particular storage cell immediately by a circuitry part which is not shown. This can, for example, be done by writing zeros into it. If various scalar products resp. sums are to be accumulated, the process of reading is not started until the computation of the full sum is finished. The summands are continuously accumulated into the accu-carry-memory.

From the most significant column the carry part of the memory is transferred into an auxiliary carry register. (24) in Figure 16. From there, this carry is transferred with a delay of one cycle via wire (20) to the least significant column to have it available for the read out process of the more significant row.

The final carry treatment (25) contains a singleresp. multi-stage pipeline where the still remaining carries are included in the result. At the end of this part of the circuitry the ready rows of the result appear, the least significant ones first.

In another part of the circuitry (26), which is shown in Figure 18, the two rows with the significant digits must be found. The most significant digit of the more significant register (28) contains the result sign; smallest digit (preferably zero) means positive, largest digit (dual 1, decimal 9, hexadecimal F) means negative. t is advisable to initialize both registers with zero. The circuitry for filtering the rows with significant information now checks in each row presented to the circuit whether there is at least one digit not equal to the sign digit already stored in the higher significant register (28). If this is the case or if there is no sign digit (e.g. 1..8 in a decimal system) at position (28) then the transfer is enabled for the actual and the next clock cycle to fill both registers with two new consecutive rows. If, however, the transfer was already enabled in the previous cycle, then it must be reenabled for one cycle only. The control circuit (29) may therefore be described by the following state table with entries "next state/transfer enable".

State	output sign O	of check 1	
1	1/0	2/1	
2	1/1	3/1	
3	1/0	3/1	

The transfer into the registers ends if only rows with sign digits follow. Finally, in both registers those rows appear, which contain the mantissa of the floating-point result. One obtains the exponent from the position as well as from the initial address resp. from the number of cycles necessary for reading. Furthermore, the information required for the rounding is easily obtained during output. It serves for a possible adaptation of the result.

The circuitry shown in Figure 17 may be varied to reduce the number of input/output lines, e.g. by transferring the carry count (19) through the MANT inputs. The Figure is intended just to show principles, and not tricky details.

6. <u>Systems with large Exponent Range and further</u> <u>Remarks</u>

Many computers have a very modest exponent range. This is for instance the case for the system /370 architecture. If in the decimal system, for instance, l = 17, el = -75 and e2 = 75 the full length L = k + 2e2 + 2l + 2 |e| of the registers (see Figure 1 and Figure 2) can more or less easily be provided. Then sums and scalar products of the form (I) and (II) can be correctly computed for all possible combinations of the data by the technique discussed in this paper without ever getting an overflow or an interrupt.

However, there are also computers on the market with a very large exponent range of several hundred or thousand. In such a case it may be costly to provide the full register lengths of L = k + 2e2 + 2l + 2 |e1| for the techniques discussed in this paper. It is most useful then to reduce the register lengths to the single exponent range and instead of L to choose $L^{H} = k + e2 + 2l + |e1|$ or even a smaller range e' $\leq e \leq e''$ with el $\langle e'$ and

 $e'' < e^2$ and correspondingly L' = k + e'' + 21 + |e'|.

Traditionally, sums and scalar products are computed in the single exponent range $el \leq e \leq e2$. If |e1| and e2 are relatively large most scalar products will be correctly computable within this range or even in $e' \leq e \leq e''$. Whenever, in this case, the exponent of a summand in a sum or scalar product computation exceeds this range $e' \leq e \leq e''$

an overflow has to be signalled which may cause an interrupt.

In such a case the exponent range could be extended to a larger size on the negative or the positive side or even on both sides. We may very well assume that the necessity for such an extension of the exponent range occurs rather rarely. The supplementary register extensions, which are necessary for the techniques discussed in this paper, could then, for instance, be arranged in the main memory of the system and the summation within the extended register part may then be executed in software. Such procedure would slow down the computation of scalar products in rather rare cases. But it still always will deliver the correct answer.

£

We further discuss a few slightly different methods how to execute accumulating addition/subtraction and the scalar product summation on processors with large exponent range.

On a more sophisticated processor the exponent range covered by the summing matrix could even be made adjustable to gain most out of this special hardware. This could be done by an automatic process of three stages:

- 1. A special vector instruction analyzes the two vectors and computes the exponent range that covers most of the summands or products of the vector components. This step may be discarded if the best range is already known.
- 2. The summing matrix gets properly adjusted to the range found in 1. and in a vector instruction the fitting part of the summand or products is accumulated into the summing matrix. If a summand or product does not fit into it it can be dealt by one of the two alternatives:
 - a) Interrupt the accumulation and add that summand or product by software to the not covered extended parts of the accumulator which resides in main memory.
 - b) Do not interrupt the accumulation, but discard this summand or product and mark this element in a vector flag register. Later the marked elements are added by software to the extended parts of the accumulator. This second way avoids interrupting and restarting the pipeline and will thus lead to higher performance than a).
- 3. In a final step the content of the summing matrix part of the accumulator is properly inserted between the extended parts to get the complete result in form of a correspondingly long variable in main memory.

Another cure of the overflow situation e € [e', e"] may be the following: Summands with an exponent e, which is less than e', are not added, but gathered on a "negative heap". Similarly summands with an exponent, which is greater than e", are gathered on a "positive heap". The negative and the positive heap may consist of a bit string or a vector flag register where each summand or vector component is represented by a bit. This bit is set zero if the summand was already added. It is set 1 if the component belongs to the corresponding heap. After a first summation pass over all summands the computed sum is stored. Then the positive and/or negative heap is shifted into te middle of the exponent range e' $\leq e \leq e''$ by an exponent transformation and then added by the same procedure. After possibly several such steps the stored parts of the sum are put together and the final sum is computed. In many cases it will be possible to obtain the final result without summing up the negative heap.

Another possibility to obtain the correct result with a reduced register length L' = k + e' + 2l +e" is the following: The process of summation starts as usual. As soon as the exponent e of a : ummand exceeds the range [e', e"] an exponent part is built up which interprets the digit sequence of L' as a very long mantissa of a normalized floating-point number. The normalization, in general, will require a shift. Then a "positive heap" is no longer necessary. And in most cases it will be possible to obtain the correct rounded result without summing up a possibly still necessary "negative heap". The method computes all accumulating sums or scalar products correctly without considering the negative heaps as long as less than e" - e' digits cancel. The negative heap can only influence the k least significant digits of L'.

The reduction of the full accumulator length L to a smaller size L' \langle L may cause exponent under- or overflows in special summation processes. This always makes some event handling routine necessary. Whatever this is, this procedure represents a trade off between hardware expenditure and runtime.

A rather primitive event handling would consist in a traditional summation of the positive and negative heap. In this case a message should be delivered to the user that the result is probably not precise.

In the context of programming languages the accumulator of length L' = $k + e^{"} + 21 + e^{'}$ represents a new data type which could be called <u>precise</u>. As long as no exponent under- or overflow occurs (e' $\leq e \leq e^{"}$) <u>addition of variables of type real</u>, of products of such variables as well as of scalar products of real vectors into a variable of this type can precisely be executed and it is error free. Accumulation of real variables, products or scalar products into a variable of type precise is associative. The result is independent of the order in which the summands are added.

Vectorprocessors belong to the fastest computers which are presently available. Their main field of application is scientific computation. It should be natural that vectorprocessors compute vector operations correctly. The vector operations consist basically of the componentwise addition and subtraction, the componentwise multiplication and the scalar product. The implementation of highly accurate vector addition/subtraction and componentwise multiplication belongs to the state of the art. The computation of accurate scalar products has been dealt with in this paper.

Due to their high speed of computation, vectorprocessors must, however, also be able to support an automatic error analysis resp. verification of the computed result. In order to achieve this it is necessary that all operations, mentioned above, such as componentwise addition/subtraction, componentwise multiplication and scalar products can optionally be called with several roundings, in particular with the monotone downwardly directed rounding, the monotone upwardly directed rounding and the rounding to the least including interval. We do not discuss the implementation of these roundings here. It belongs to the state of the art. For further information we refer to the literature.

Finally, we remark that the methods and procedures outlined in this paper are also suitable to add up sums of products correctly which consist of more than two factors, for example

$$\sum_{i=1}^{n} a_i * b_i * c_i$$

Application to Multiple Precision Arithmetic 7.

We show in this chapter that the essential parts of multiple precision arithmetic can easily be executed with high speed if a fast scalar product unit is available.

We consider

- 1. Double Precision Arithmetic⁶ 1.1 Sum and Difference

It is clear that sums of two or n double precision summands a + b or a + b + c ... + z can be accumulated. The same holds for sums of vectors or matrices.

1.2 Product

If a product a • b of two double precision factors a and b has to be computed, each factor can be represented as a sum of two single precision numbers $a = a_1 + a_2$ and $b = b_1 + b_2$, where a_1 and b_1 represent the first (higher significant) 1 digits and a_2 and b_2 represent the last (lower signifi-

cant) 1 digits of a and b. The multiplication then requires the execution of a scalar product:

$$a \cdot b = (a_1 + a_2) (b_1 + b_2) = a_1b_1 + a_1b_2 + a_2b_1 + a_2b_2 ,$$
 (1)

where each summand is of double precision. These can be added by the techniques developed in this paper.

Similarly, products of more than two factors can be computed. As in (1) products of two double precision numbers are expressed by a scalar product of single precision numbers. On the right hand side of (1) each summand is a double precision number which can be expressed by a sum of two single precision numbers. In the case of a product of four double precision numbers this leads to the following formulas, which are self-explanatory.

$$\begin{array}{c}
8 \\
\Sigma a^{i} \cdot \Sigma c^{i} = \Sigma \\
i=1 \\
\text{with } a \cdot b = \sum_{i=1}^{8} a^{i} \text{ and } c \cdot d = \sum_{i=1}^{8} c^{i}.
\end{array}$$

Thus a b c d can be computed as the sum of 64 products of two single precision numbers each. The case of products of two or more double precision matrices is a little more difficult. But it can, in principle, be treated similarily. If a

product of two double precision matrices has to be computed the two matrices are first represented as ⁶High speed scientific computation is usually done in the long data format. Double precision here means the double mantissa length of that format.

If the usual long format is already called double precision our double precision corresponds to

quadruple or extended precision.

sums of two single precision matrices. Multiplication of these sums then leads to a sum of products of single precision matrices:

$$a \cdot b = (a_1 + a_2) (b_1 + b_2) =$$

 $a_1b_1 + a_1b_2 + a_2b_1 + a_2b_2$ (2)

Each component of the products on the right hand side of (2) is computed as a scalar product. Thus each component of the product matrix a • b consists of a sum of scalar products which itself is a scalar product.

In case of matrix products, which consist of more than two double precision matrix factors, one has to take into account that the components of (2) may already be pretty long. They may consist of 10 or 20 consecutive digit sequences of single precision lengths. These sums of single precision matrices then have to be multiplied with other such sums, which leads to a sum of matrix products. Each component of this sum can be computed as a scalar product of single precision numbers.

2.

Arithmetic of triple precision is a special case of quadruple precision arithmetic.

Quadruple Precision Arithmetic 3.

3.1 Sum and Difference

Each summand of quadruple precision can be represented as a sum of two double precision summands. Thus sums of two or more quadruple precision summands can be added as expressed by the following formulas:

$$a + b = a_1 + a_2 + b_1 + b_2$$

a + b + c + ... + z = $a_1 + a_2 + b_1 + b_2 + c_1 + c_2 + \dots + z_1 + z_2$.

Sums of quadruple precision vectors or matrices can be treated correspondingly.

3.2 Products

Each quadruple precision number can be represented as a sum of four single precision numbers $a = a_1 + a_2$

 $a_2 + a_3 + a_4$. Multiplication of such sums requires the execution of a scalar product:

Similarily, products of more than two quadruple precision factors can be computed. We indicate this process by the following formulas, which are self-explanatory.

$$\begin{array}{l} \mathbf{a} \cdot \mathbf{b} \cdot \mathbf{c} \cdot \mathbf{d} = \\ (\mathbf{a} \cdot \mathbf{b}) \ (\mathbf{c} \cdot \mathbf{d}) = \begin{pmatrix} 4 & 4 \\ \Sigma & \Sigma & \mathbf{a}_i & \mathbf{b}_j \end{pmatrix} \begin{pmatrix} 5 & \Sigma & \mathbf{c}_i \mathbf{d}_j \end{pmatrix} = \\ i = 1 \quad j = 1 \end{array}$$

$$\begin{array}{c} 32 \\ = (\sum_{i=1}^{32} a^{i}) (\sum_{j=1}^{32} c^{j}) = \sum_{i=1}^{32} \sum_{j=1}^{32} a^{i}c^{j} . \quad (4) \\ \end{array}$$

There the 16 double precision summands $a_i b_j$ and $c_i d_j$ of the two factors of (4) are each represented as sums of two single precision-numbers. This leads to the product of the two sums over 32 single precision numbers a^i resp. c^j in the next line.

If a product of two quadruple precision matrices is to be computed each factor is represented by a sum of four single precision floating-point matrices as in (3).

Multiplication of these sums leads to a sum of matrix products. Each component of these matrix products is computed as a scalar product. The sum of these scalar products is again a scalar product.

It was the intention of this section to demonstrate that with a fast accumulating addition/subtraction or scalar product unit a big step towards multiple precision arithmetic, even for product spaces, can be done.

- 8. Literature
- U. Kulisch: Grundlagen des Numerischen Rechnens - Mathematische Begründung der Rechnerarithmetik, Bibliographisches Institut, Mannheim 1976
- [2] U. Kulisch and W.L. Miranker: Computer Arithmetic in Teory and Practice, Academic Press 1981
- [3] U. Kulisch and W.L. Miranker: The Arithmetic of the Digital Computer: A New Approach, SIAM-Review, March 1986, pp. 1-40
- [4] IBM System /370 RPQ. High Accuracy Arithmetic, Publication Number SA 22-7093-0
- [5] High Accuracy Arithmetic, Subroutine Library, General Information Manual, IBM Program Number 5664-185
- [6] High Accuracy Arithmetic, Subroutine Library, Program Description and User's Guide, IBM Program Number 5664-185, Publication Number GC 33-6163
- [7] T. Teufel: Ein optimaler Gleitkommaprozessor. Dissertation, Universität Karlsruhe, 1984
- [8] G. Bohlender and T. Teufel: BAP-SC: A Decimal Floating-Point Processor for Optimal Arithmetic, to appear in: Computer Arithmetic, Scientific Computing and Programming Languages (E. Kaucher, U. Kulisch, Ch. Ullrich, Eds), B.G. Teubner, 1987
- [9] Arithmos Benutzerhandbuch, SIEMENS AG., Bestell-Nr.: U 2900-J-Z 87-1

For a supplementary bibliography see the literature listed in [3].







each row contains c adders of a digits and n transfer registers of t digits, i.e. h=c+a=n+t (here: n=c and t=a) Figure 5: Structure of the summing matrix





267



digit with exponent e 00...0 00.....0 00.....0.001 row with a transfer registers of length t Case_2: x < # 100 00....0 00.....0 00....0re ۰, ۰, row with a transfer registers of length t Figure 8: Task of the shift unit Case 1: # 2 8 8 226 addition of all digits of the mantissa in row y. Case 2: x < 8 يتكر ã••₁• •₂ part m_1 of the mantissa is added in row y, part m_2 in row y-1. •2___ y = (e-e_o) <u>div</u> h ment of the most significant digit of the montissa e: expo B: length of mantissa, number of digits of the mantissa h: number of digits of a row of the summing matrix Figure 9: Description of the shift process Exponent identification t_{e} of the transfer sections of the matrix: ۰ n+1 0 2n-1 1 (x-1)-n-1 a(5-3) r-2 {r-1]0 g·n -1 z-1 . A centifies transfer row with mastisse exponent 0 following exponent identifications in the corresp transfer sections: $\theta_{01} \theta_{0-1}, \dots$. e gets the 00.....0 00 [en en-1 9₈-2

<u>Lase 1:</u> = 2 =



This is independent of the location in the transfer row. For instance:





Figure 7: Exponent coordinates of the digits in the summing matrix







Figure 16: Structure of the summing unit with only one row of adders



