

$\text{TR}$	$\supset$	$D$	$\supset$	$R$
$\text{VR}$	$\supset$	$VD$	$\supset$	$VR$
$\text{MR}$	$\supset$	$MD$	$\supset$	$MR$

$\text{PR}$	$\supset$	$IR$	$\supset$	$ID$	$\supset$	$IR$
$\text{PVR}$	$\supset$	$I\text{VR}$	$\supset$	$I\text{VD}$	$\supset$	$I\text{VR}$
$\text{PMR}$	$\supset$	$\text{IMR}$	$\supset$	$\text{IMD}$	$\supset$	$\text{IMR}$

$C$	$\supset$	$CD$	$\supset$	$CR$
$\text{VC}$	$\supset$	$VCD$	$\supset$	$VCR$
$MC$	$\supset$	$MCD$	$\supset$	$MCR$

$\text{PC}$	$\supset$	$IC$	$\supset$	$ICD$	$\supset$	$ICR$
$\text{PVC}$	$\supset$	$I\text{VC}$	$\supset$	$I\text{VCD}$	$\supset$	$I\text{VCR}$
$\text{PMC}$	$\supset$	$\text{IMC}$	$\supset$	$\text{IMCD}$	$\supset$	$\text{IMCR}$

Operat. in  $\text{TR}, \text{VR}, \text{MR}, C, VC, MC$

$\{M, *\}, \{\text{PM}, *\} : \wedge \quad x * y := \{x * y \mid x \in x \wedge y \in y\}$   
 $x, y \in \text{PM}$

$\Rightarrow$  Operat. in  $\text{PR}, \text{PVR}, \text{PMR}, \text{PC}, \text{PVC}, \text{PMC}$

Kulisch  
7 June 88

1. trad. Airthm.

$$\{R, \oplus, \ominus, \otimes, \oslash, \leq\}$$

$$\begin{array}{ccc} \mathbb{T}R & \longrightarrow & R \\ \downarrow & & \downarrow \\ \mathbb{C} & & CR \\ \downarrow & & \downarrow \\ M\mathbb{T}R & & MR \end{array}$$

$$\alpha = \alpha_1 + i\alpha_2, \beta = \beta_1 + i\beta_2 \in CR$$

$$\alpha \boxtimes \beta := (\alpha_1 \otimes \beta_1, \alpha_2 \otimes \beta_2, \alpha_1 \otimes \beta_2 + \alpha_2 \otimes \beta_1)$$

$$a = (a_{ij}), b = (b_{ij}) \in MR$$

$$a \boxtimes b := \left( \sum_{k=1}^n a_{ik} \otimes b_{kj} \right)$$

2. Semimorphism.  $\square : M \rightarrow N$

$$\begin{array}{ccc} \mathbb{T}R & \longrightarrow & R \\ \downarrow & & \downarrow \\ \mathbb{C} & \longrightarrow & CR \\ \downarrow & & \downarrow \\ M\mathbb{T}R & \longrightarrow & MR \end{array}$$

$$a \boxtimes b := \square(a * b) \quad \text{f.a. } a, b \in N, * \in \{+, -, \times, /\}$$

$$\square a = a \quad \text{f.a. } a \in N$$

$$a \leq b \Rightarrow \square a \leq \square b \quad \text{f.a. } a, b \in M$$

$$\square(-a) = -\square a \quad \text{f.a. } a \in M$$

$$a \leq \square a \quad \text{f.a. } a \in M (= I..)$$

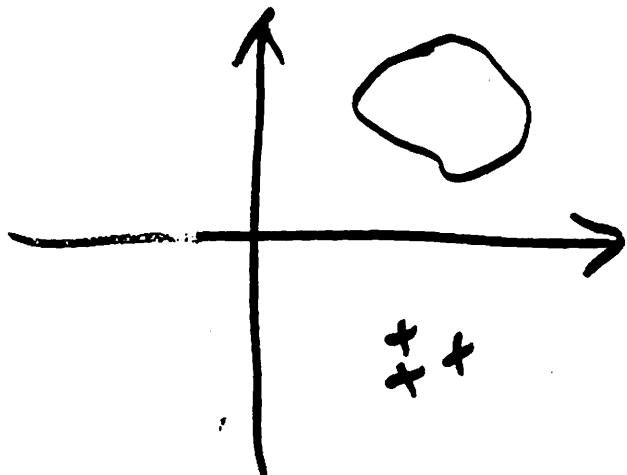
$$\begin{aligned} \alpha \boxtimes \beta &= \square(\alpha * \beta) = \square(\alpha_1 * \beta_1 - \alpha_2 * \beta_2, \alpha_1 * \beta_2 + \alpha_2 * \beta_1) \\ &= (\square(\alpha_1 * \beta_1 - \alpha_2 * \beta_2), \square(\alpha_1 * \beta_2 + \alpha_2 * \beta_1)) \end{aligned}$$

$$a \boxtimes b = \square(a * b) = \square\left(\sum_{k=1}^n a_{ik} b_{kj}\right) = \left(\square \sum_{k=1}^n a_{ik} b_{kj}\right)$$

$$a = \begin{pmatrix} a_1 \\ \vdots \\ a_n \end{pmatrix} \quad b = \begin{pmatrix} b_1 \\ \vdots \\ b_n \end{pmatrix}$$

$$a \leq b$$

$a_i \in h_i$  f.a.i



$$\square \alpha \leq \kappa \leq \beta \quad R$$

$$\alpha \quad \kappa \quad \beta$$

$$\square / \alpha \leq \kappa \leq \beta$$

$$(R1) \quad \square \alpha \in \square \kappa \in \square \beta$$

$$(R1) \quad \alpha \leq \square \kappa \leq \beta$$

$$(R2) \quad \alpha \leq \kappa * y \leq \beta$$

$$(R3) \quad \square \alpha \in \square(\kappa * y) \leq \square \beta$$

$$(R4) \quad \alpha \leq \square(\kappa * y) \leq \beta$$

$$(R5) \quad \alpha \leq \kappa \boxplus y \leq \beta$$

## fixed-point arithmetic

$$\begin{array}{r} 0.2143769851 \\ 0.3214367534 \end{array}$$

$$\begin{array}{r} 0.7124342678 \\ 0.8123950690 \end{array}$$

$$\begin{array}{r} 0.1243467809 \\ 13.2467869097 \\ \hline 0.1324678691 \end{array}$$

addition / subtraction in fixed-point arithmetic  
is 'error free'

scaling requirement  $\rightarrow$  overscaling

$$0.0000021473$$

loss of accuracy

therefore: floating-point arithmetic  
exponent part takes care of the scaling automatically

multiplication and division rel. stable operations  
addition / subtraction are problematic

ideal computer:

multiplication, division in floating-point arithm.  
addition, subtraction in fixed-point arithm.

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

$$x \cdot y = 10^{50} + 2446 - 10^{50} + 10^{40} + 6333 - 10^{40} = 8779$$

$$x \square y = 0$$

$$x = 0.10005 \times 10^5$$

$$y = -0.99973 \times 10^4$$

$$-0.1000500 \times 10^5$$

$$-0.0999730 \times 10^5$$

$$0.0000770 \times 10^5$$

$$0.77000 \times 10^1$$

$$x = \square(x_1 \times x_2), \quad x_1 \times x_2 = 0.1000548241 \times 10^5$$

$$y = \square(y_1 \times y_2), \quad y_1 \times y_2 = 0.9997342213 \times 10^4$$

$$x_1 \times x_2 - y_1 \times y_2 = -0.10005482410 \times 10^5$$

$$-0.09997342213 \times 10^5$$

$$= 0.00008140197 \times 10^5$$

$$x_1 \times x_2 - y_1 \times y_2 = 0.8140197 \times 10^1$$

$$\square(x_1 \times x_2 - y_1 \times y_2) = 0.81402 \times 10^1$$

# scientific computing

standard probl. of num. analysis with verified results; lin. syst.; matrix inv.; eigenprob.; linear optimiz.; polyn. eval.; zeros; arithm. expr.; nonlin. equations; num. quadr.; ordin. diff. equations, ...

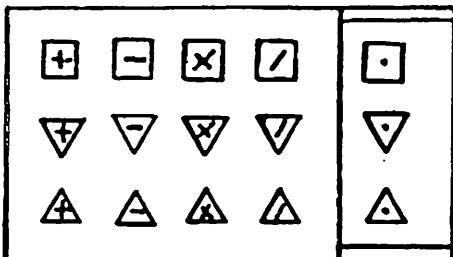
arithmetic in product spaces and corres.  
ponding interval spaces; language extens.

extended computer arithm.  $\square \triangleright \triangle$

elementary computer arithm.



implementation of semimorph. in  
 R, UR, MR, CR, VCR, MCR, IR, IUR, IMR, ICR, IVCR, IMCR:



$$a \square b = \square \sum_{i=1}^n a_i \cdot b_i$$

$$a \nabla b = \nabla \sum_{i=1}^n a_i \cdot b_i$$

$$a \Delta b = \Delta \sum_{i=1}^n a_i \cdot b_i$$

\* traditional numerical analysis  $* \in \{+, -, \times, /\}$

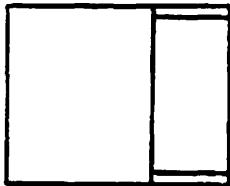
\nabla, \Delta traditional interval arithmetic



Z23 (1967) in software, ext. ALGOL

X8 (1968) in hardware, ext. ALGOL

IEEE-arithmetic-standard (1983)



high accuracy, operations in product sets

Z80 (1980)

PASCAL-SC

JBM-PC (1983)

MOTOROLA 68000 (1982)

hardware unit (1983, G. Schweizer)

JBM 1370, 4361, ACRITH (1983)

JBM 9377 (1986)

Siemens ARITHMOS (1986)

NAS, Hitachi, BASF, Nixdorf

GAMM-Resolution on Computer Arithmet

FORTRAN-SC

## Higher Order Computer Arithmetic

linear systems, matrix inversion, eigen-problems,  
linear programming problems, expression evalua-  
tion with maximum quality

1.  $A, B$  matrices,  $x, y, z$  vectors

$$x = x + A * y + B * z$$

$$x = \# * (x + A * y + B * z)$$

$$x = \# < ( ), x = \# > ( ), x = \# \# ( ), x = \# ( )$$

2.  $A_i, B_i, i=1(1)n$ , vectors or matrices

$$x = \# * (\text{sum}(A(i) * B(i), i=1, n))$$

computes  $x = \sum_{i=1}^n A_i \cdot B_i$  with max. quality

3. evaluation of expressions with max. quality

$$b = \# * (x + 4 * (3.0 e 8 * y / z))$$

$$b = \# < (((4 * x - 5) * x + 3) * x + 2.5 e 3)$$

$$c = \# \# (\text{sum}(a(i) * x^{**i}, i=1, n))$$

computes  $\sum_{i=1}^n a_i \cdot x^i$  to max. quality

4. computation of program parts with max. quality

accurate ( $x, y <, z >$ ) do

begin PROG

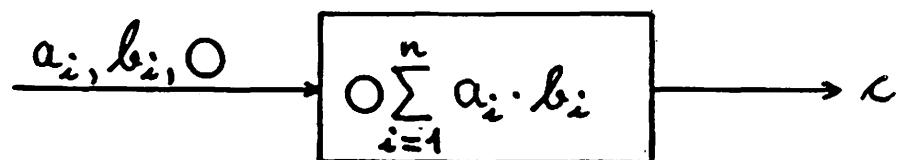
end

computes  $x, y, z$  to max. quality and rounds  
 $x$  to nearest,  $y$  downwards,  $z$  upwards

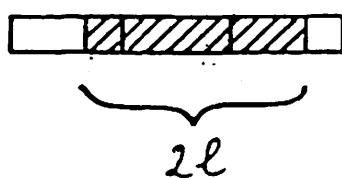
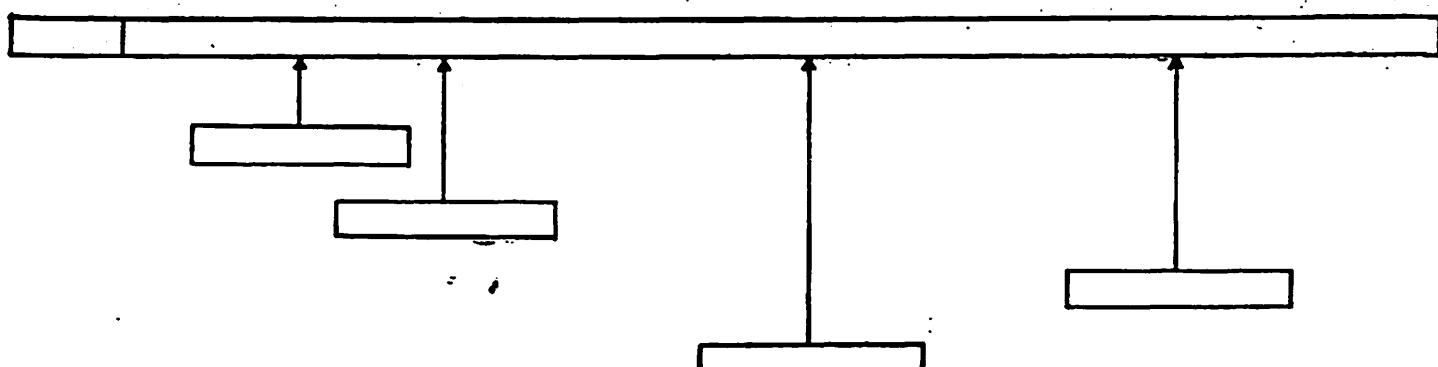
$$Z = m \cdot b^e, \quad m = 0.9784231, \quad b = 10, \quad e = 12$$

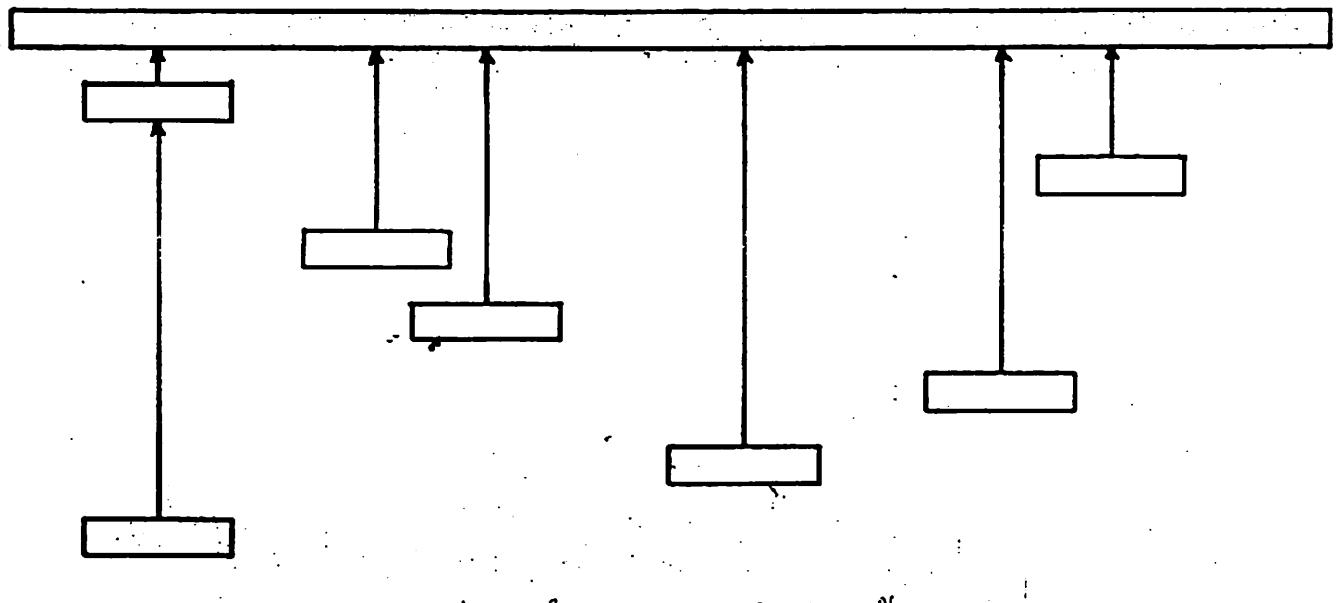
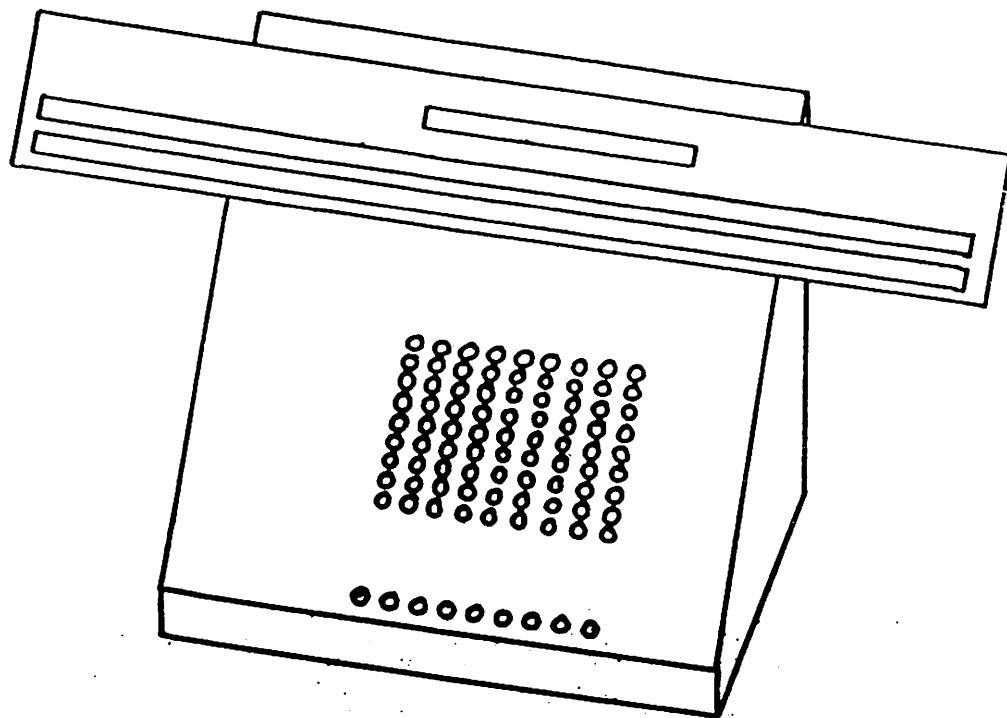
$$a = (a_i), \quad b = (b_i), \quad a_i, b_i \in R(b, l, e_1, e_2)$$

$$c = O \sum_{i=1}^n a_i \cdot b_i, \quad a_i \cdot b_i \in R(b, 2l, 2e_1, 2e_2) \\ O \in \{\square, \nabla, \Delta\}$$



$$L = k + 2e_2 + e_2 + 2l + 21e_11$$





simpler      applied as often as possible  
faster  
fully accurate

generalization of this operation for  
electronic computers has the same  
properties

simpler

# GESELLSCHAFT FÜR ANGEWANDTE MATHEMATIK UND MECHANIK ( GAMM )

## Resolution on Computer Arithmetic

The elementary floating-point operations +, -, \*, / in electronic computers are currently required to be of highest machine accuracy: For any choice of operands, the computed result must coincide with the rounded exact result of the operation, rounded according to the rounding mode in use (if no overflow occurs). For reference, see the IEEE Arithmetic Standards 754 (binary floating-point arithmetic) and 854 (general floating-point arithmetic).

In recent years there has been a significant shift of numerical computation from general-purpose computers towards vector and parallel computers - so-called supercomputers. Along with the 4 elementary operations +, -, \*, /, these computers usually offer compound operations as additional elementary operations. This leads to an increase of several orders of magnitude in computing power. Some of these elementary compound operations are:

- multiply and add:  $a * b + c$
  - multiply and subtract:  $a * b - c$
  - accumulate: computes the sum of the components of a vector
  - multiply and accumulate: computes the inner (or scalar) product of two vectors
- and others.

GAMM requires that all elementary compound operations be implemented by the manufacturer in such a way that guaranteed bounds are delivered for the deviation of the floating-point result from the exact result. It is desirable and usually achievable that for all possible data the computed result of such a compound floating-point operation agrees with the result that would be obtained if the exact result were computed and then rounded by the rounding in use (if no overflow occurs). In this case no explicit error bounds need be delivered. The user should not be obliged to perform an error analysis every time an elementary compound operation, predefined by the manufacturer, is employed.

All elementary compound operations should also be provided with directed roundings, a feature needed both for fast computation of reliable and narrow bounds in numerical algorithms and for verification of the correctness of computed results. It must be ensured that the final floating-point result can differ from the exact result only in the direction defined by the rounding in use. This is already required of the elementary floating-point operations by the arithmetic standards mentioned above.

(single), double, (extended, quadruple)  
real  $a, b, c, d, a_i, b_i$

$a+b$

$a \oplus b$

$a-b$

$a \ominus b$

$a * b$

$a \otimes b$

$a/b$

$a \oslash b$

$a+c*d$

multiply and add

$a \oplus c \otimes d$

$a-c*d$

multiply and subtr.

$a \ominus c \otimes d$

$a_1 + a_2 + \dots + a_n$  accumulate

$a_1 \oplus a_2 \oplus \dots \oplus a_n$

$a_1 * b_1 + a_2 * b_2 + \dots + a_n * b_n$

$a_1 \otimes b_1 \oplus a_2 \otimes b_2 \oplus \dots \oplus a_n \otimes b_n$

multiply and acc.

$\square(a+b)$

$\square(a-b)$

$\square(a*b)$

$\square(a/b)$

$\square(a*b+c*d)$

$\square(a*b-c*d)$

$\square(a_1+a_2+\dots+a_n)$

$\square(a_1*b_1+a_2*b_2+\dots+a_n*b_n)$

$\nabla(a+b)$

$\nabla(a-b)$

$\nabla(a*b)$

$\nabla(a/b)$

$\nabla(a*b+c*d)$

$\nabla(a*b-c*d)$

$\nabla(a_1+a_2+\dots+a_n)$

$\nabla(a_1*b_1+a_2*b_2+\dots+a_n*b_n)$

$\triangle(a+b)$

$\triangle(a-b)$

$\triangle(a*b)$

$\triangle(a/b)$

$\triangle(a*b+c*d)$

$\triangle(a*b-c*d)$

$\triangle(a_1+a_2+\dots+a_n)$

$\triangle(a_1*b_1+a_2*b_2+\dots+a_n*b_n)$

JEEE-St.

(1)

(2)

(3)

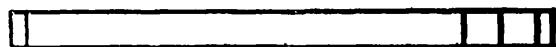
(4)

(5)

(6)

(7)

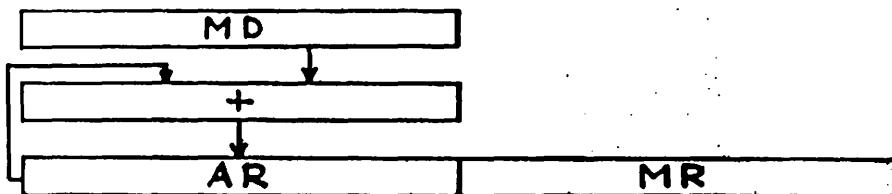
(8)



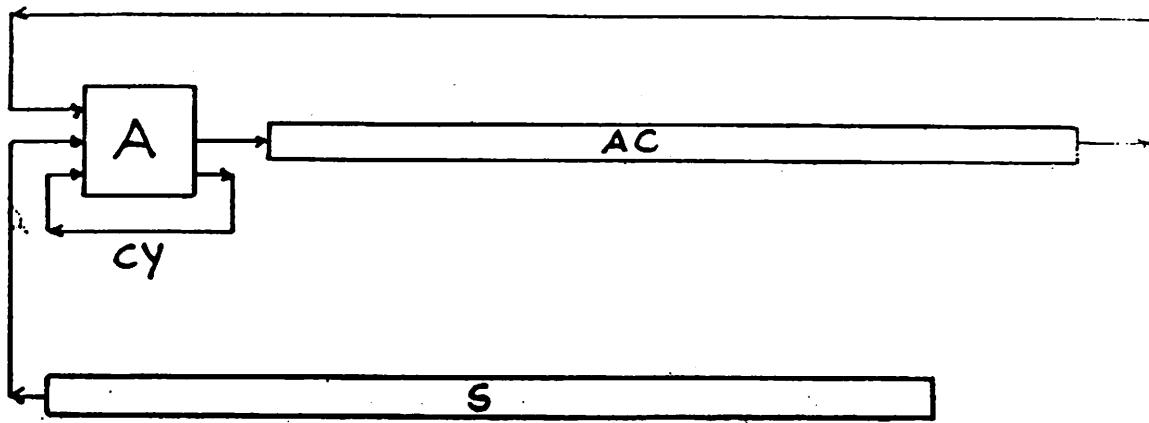
(1) ... (4) JEEE - Pt.

(1) ... (4)

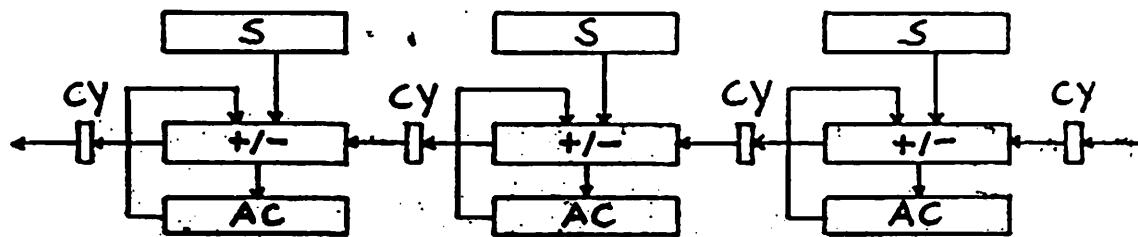
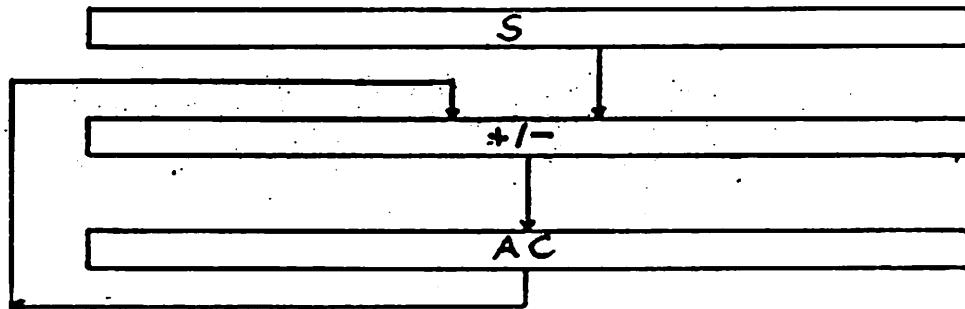
(1) ... (6) & double real



## Seriendarbeiter



## Paralleladdierer



aufgebaut aus Teilstücken

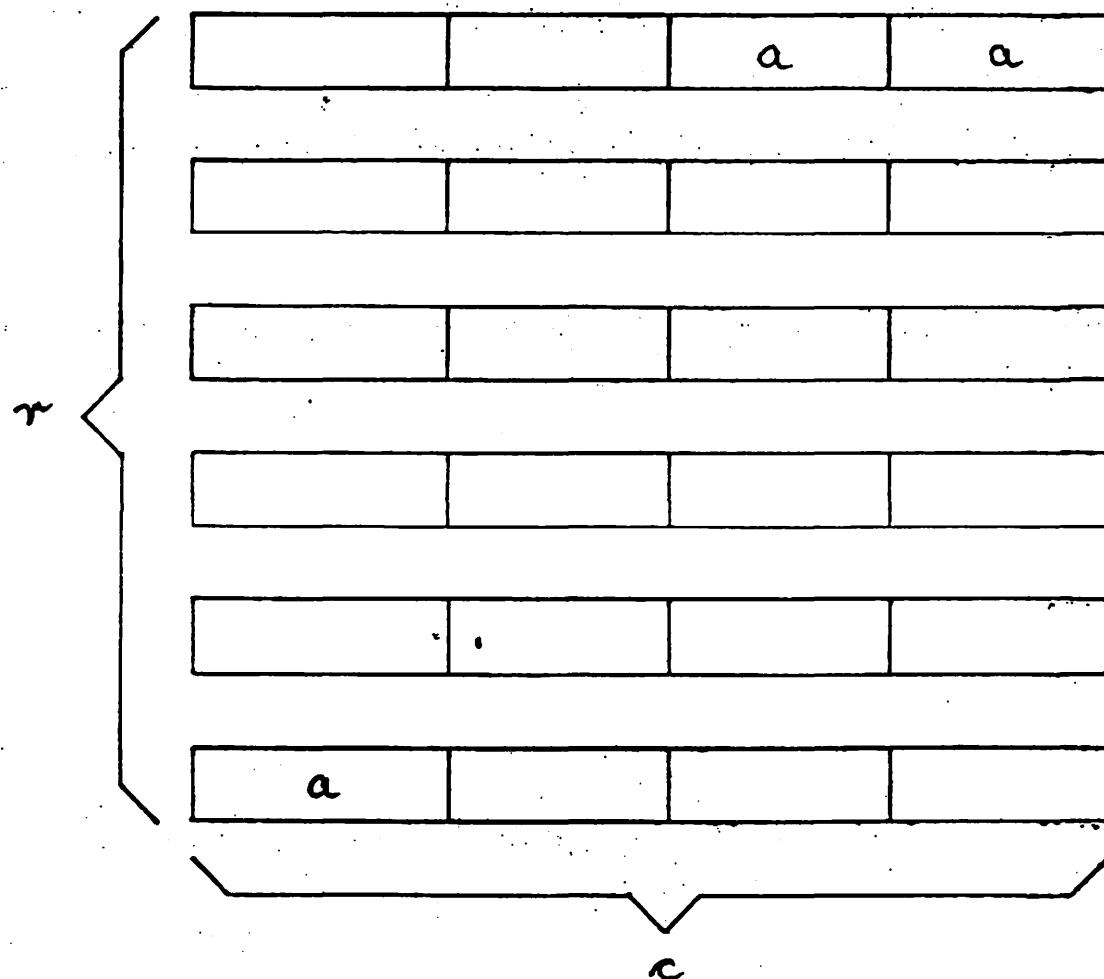
$k$	$2e2$	$2l$	$21e11$
-----	-------	------	---------

$$L = k + 2e2 + 2l + 21e11$$

$k$				$l$			
-----	--	--	--	-----	--	--	--



$$S \geq L$$



$$S = r \cdot c \cdot a \text{ Ziffern der Basis } b$$

a : Anzahl der Ziffern eines Teiladditiviers

r : Anzahl der Zeilen (Reihen, rows)

c : Anzahl der Spalten (columns)

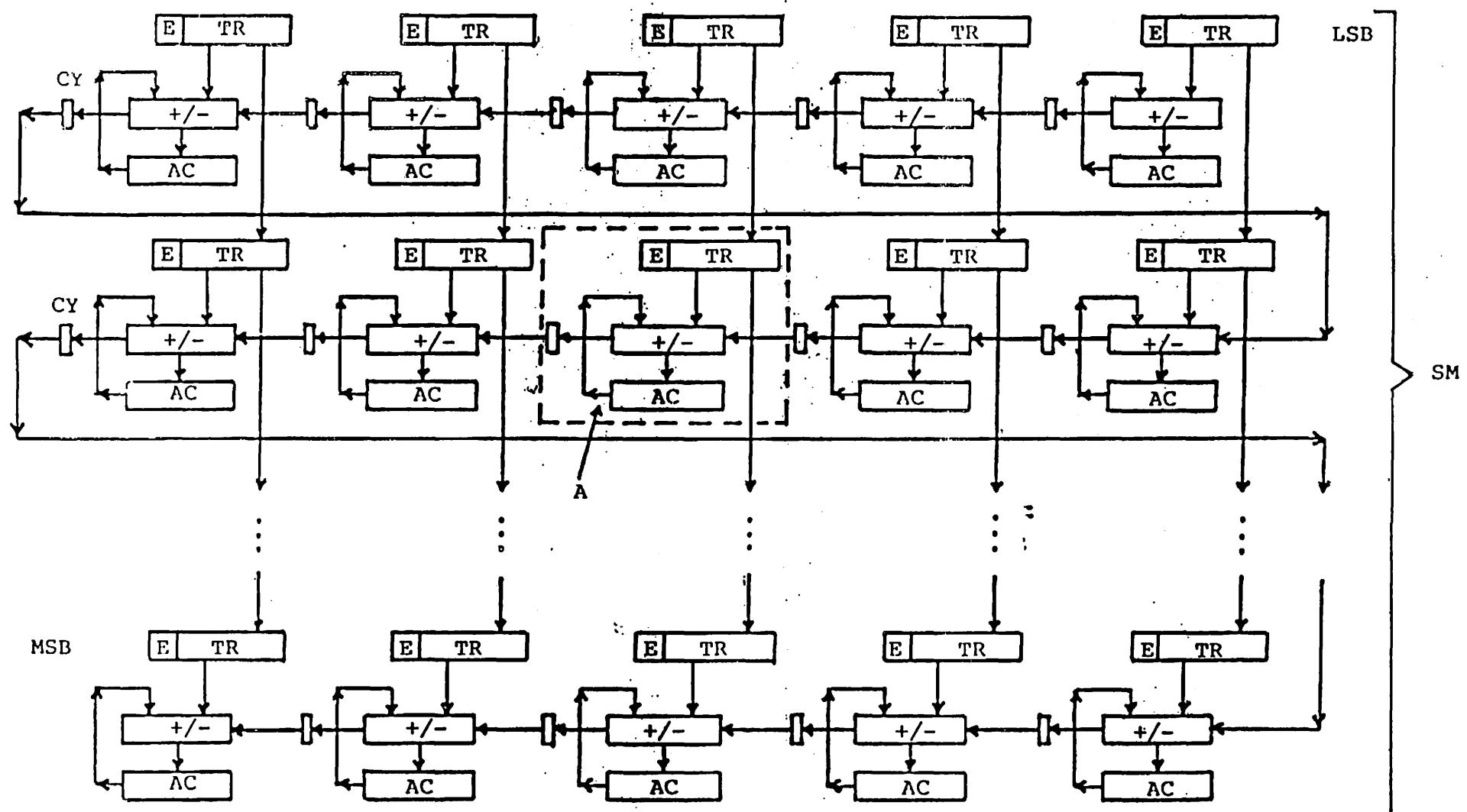
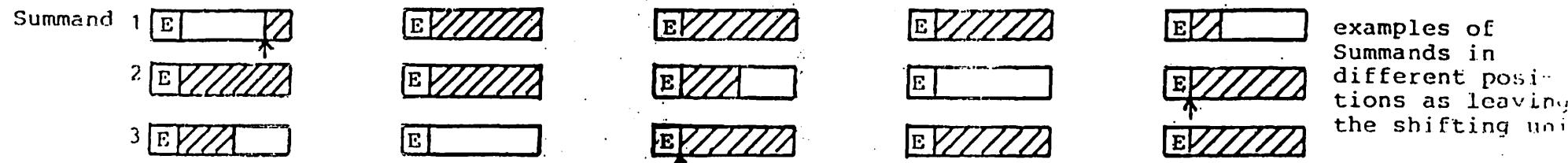
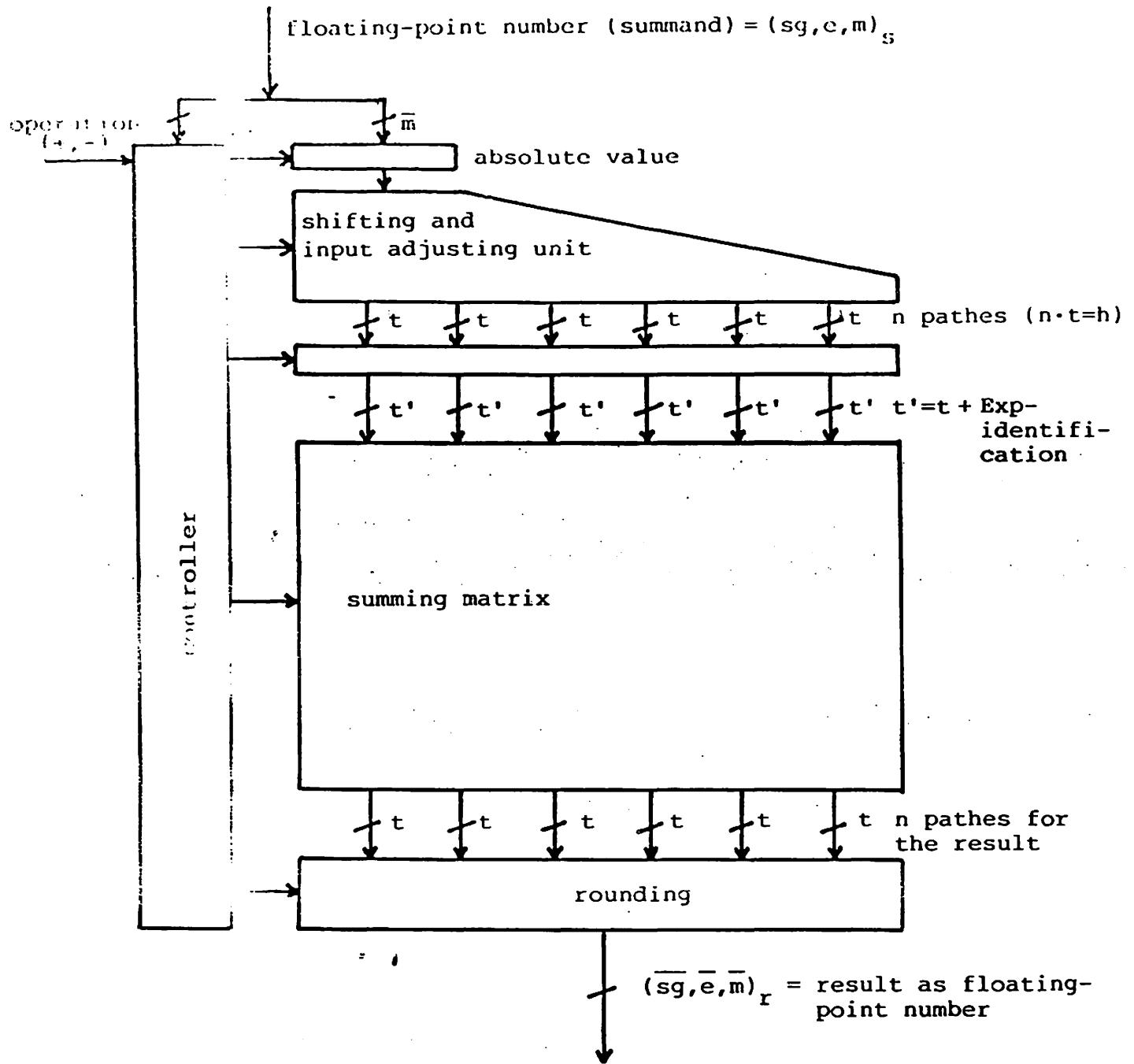
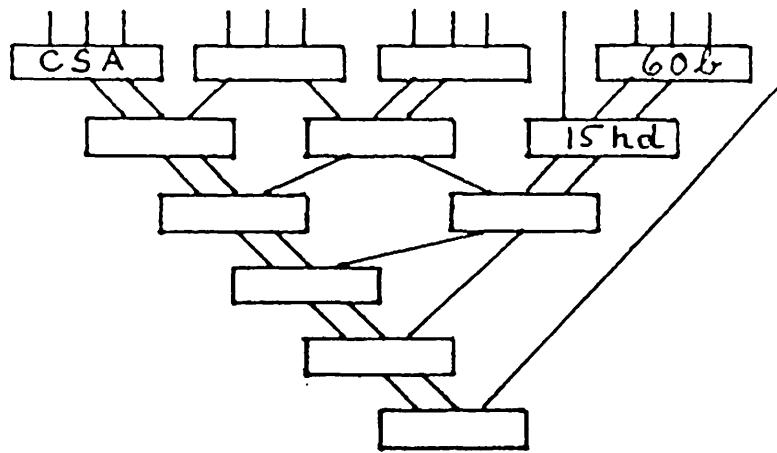


Figure 6: Summing matrix SM consisting of  $h=c \cdot r$  independent adders A  
 E: tag-register for exponent identification, TR: transfer register,  
 AC: accumulator register, CY: carry, t: most significant digit of summand

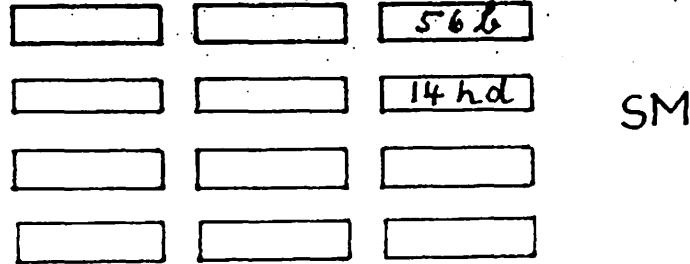


**Figure 4:** Structure of the whole circuitry

value denotes the number of figures of 'value'



W.Tr.:  $15 \times 12 = \underline{180 \text{ hd}}$        $720 \text{ b}$



SM:  $64 + 28 + 64 = \underline{156 \text{ hd}} < 42 \times 4 = 168 \text{ hd}$

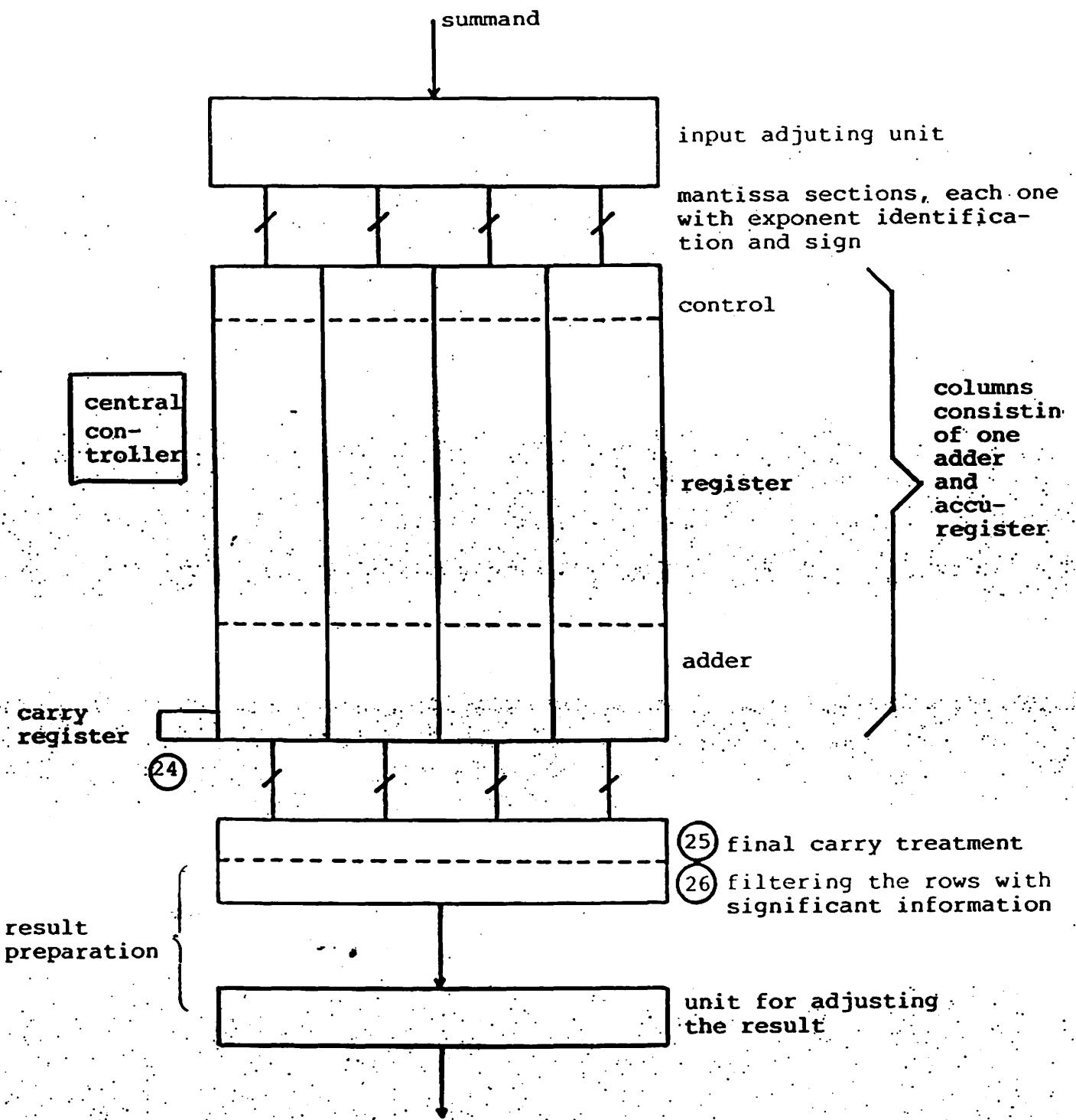


Figure 16: Structure of the summing unit with only one row of adders

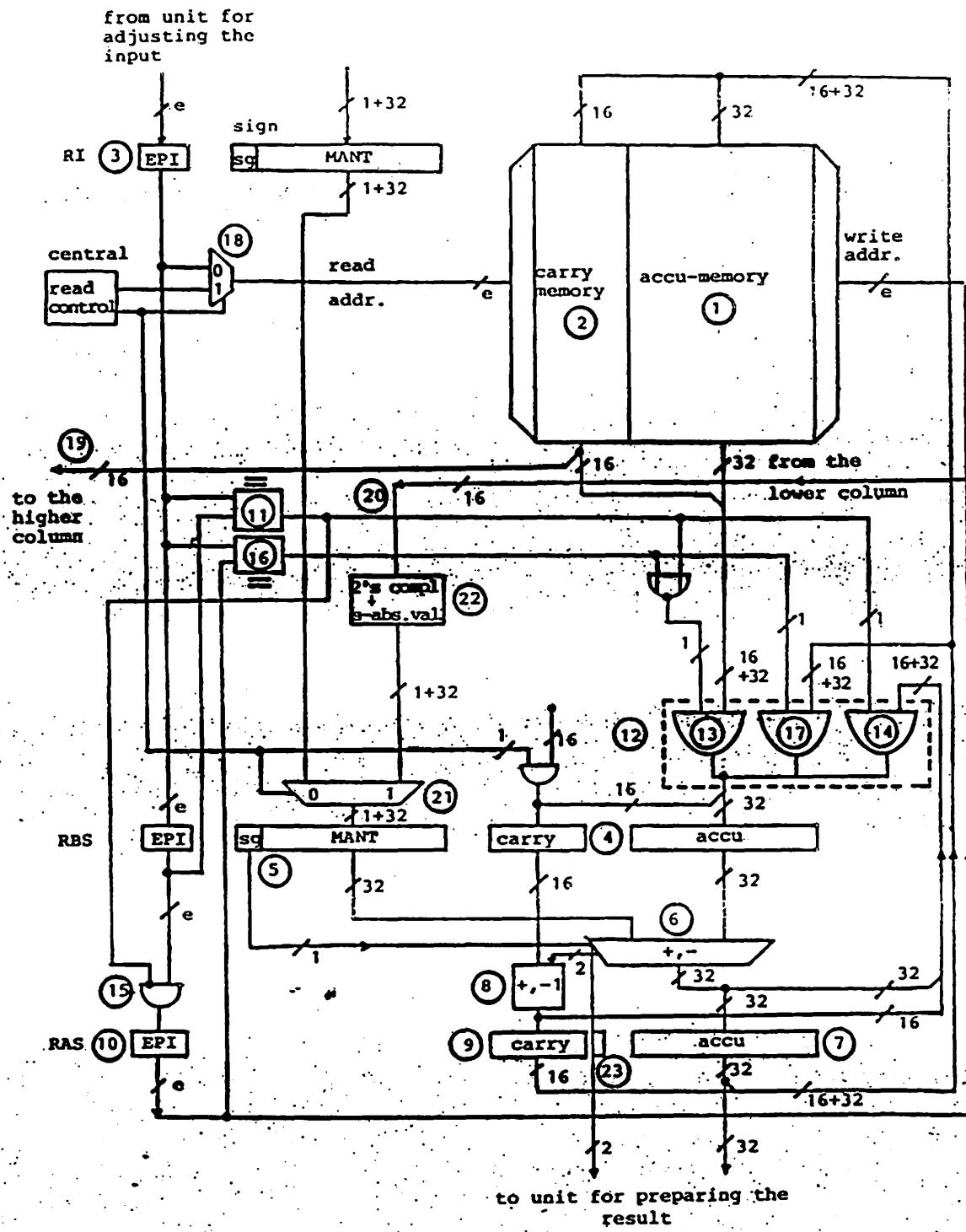


Figure 17: Structure of a "column" of the addition unit.

## Mehrfach genaue Operationen

lassen sich mit optimalem Skalarprodukt schnell und fehlerfrei ausführen:

### 1. doppelt lange Arithmetik

#### 1.1 Summe und Differenz

$$a+b$$

$$a+b+c+\dots+z$$

Summen von Matrizen ebenso

#### 1.2 Produkt

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

$$a \cdot b \cdot c \cdot d = (a \cdot b) \cdot (c \cdot d) = \sum_{i=1}^2 a_i \cdot \sum_{j=1}^2 c_j = \sum_{i=1}^2 \sum_{j=1}^2 a_i \cdot c_j$$

Produkte von Matrizen

### 2. dreifach lange Arithmetik

### 3. vierfach lange Arithmetik

#### 3.1 Summe und Differenz

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

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

Summen von Matrizen ebenso

#### 3.2 Produkt

$$a \cdot b = (a_1 + a_2 + a_3 + a_4)(b_1 + b_2 + b_3 + b_4) = \sum_{i=1}^4 \sum_{j=1}^4 a_i \cdot b_j$$

$$a \cdot b \cdot c \cdot d = (a \cdot b) \cdot (c \cdot d) = (\sum_{i=1}^4 \sum_{j=1}^4 a_i \cdot b_j) (\sum_{i=1}^4 \sum_{j=1}^4 c_i \cdot d_j) =$$

$$= (\sum_{i=1}^{32} a^i) (\sum_{j=1}^{32} c^j) = \sum_{i=1}^{32} \sum_{j=1}^{32} a^i \cdot c^j$$

Produkte von Matrizen

## Polynomial and Arithmetic Expression Evaluation:

$$p(t) = a_3 t^3 + a_2 t^2 + a_1 t + a_0 = ((a_3 t + a_2) t + a_1) t + a_0$$

$$a_0, a_1, a_2, a_3, t \in \mathbb{R}$$

$$\begin{aligned} de_1 &= a_3 \\ de_2 &= de_1 t + a_2 & -t de_1 + de_2 &= a_3 \\ de_3 &= de_2 t + a_1 & -t de_2 + de_3 &= a_2 \\ de_4 &= de_3 t + a_0 & -t de_3 + de_4 &= a_1 \\ && de_4 &= a_0 \end{aligned}$$

$$A de = b, \text{ with } A = \begin{pmatrix} 1 & 0 & 0 & 0 \\ -t & 1 & 0 & 0 \\ 0 & -t & 1 & 0 \\ 0 & 0 & -t & 1 \end{pmatrix}, de = \begin{pmatrix} de_1 \\ de_2 \\ de_3 \\ de_4 \end{pmatrix}, b = \begin{pmatrix} a_3 \\ a_2 \\ a_1 \\ a_0 \end{pmatrix}$$

$p(t) = de_4$

Arithmetic Expressions:  $a, b, c, d, e \in \mathbb{R}$

$$(a+b) \cdot c - d/e$$

$$(a+b)(c+d)$$

$$\begin{aligned} de_1 &= a \\ de_2 &= de_1 + b \\ de_3 &= c \cdot de_2 \\ de_4 &= d \\ e \cdot de_5 &= de_4 \\ de_6 &= de_3 - de_5 \end{aligned}$$

$$\begin{aligned} de_1 &= a \\ de_2 &= de_1 + b \\ de_3 &= c \\ de_4 &= de_3 + d \\ de_5 &= de_2 \cdot de_4 \end{aligned}$$

All such systems: simple form;  
 can be solved by non linear system solve. technique  
 or transferred into a linear system by an  
 algebraic transformation process

Step: diadic to multiadic operations  
with maximum accuracy

## Linear Systems of Equations      $A \cdot x = b$

$\hat{x}$ : solution;  $\tilde{x}$ : approximation;  $e = \hat{x} - \tilde{x}$ : error;

$b - A\hat{x} = d$  : defect can be computed with full accuracy  
 $b - A\tilde{x} = 0$

$$Ae = d$$

$$\text{If } e = \hat{x} - \tilde{x} \in E \Rightarrow \hat{x} \in \tilde{x} + E.$$

Interval iteration scheme:

$$E_{n+1} := (\underbrace{I - RA}_{\uparrow} \underbrace{E_n + Rd}_{\uparrow}) \quad (*)$$

converges for every  $E_0 \in V_n \cap \mathbb{R}$  to the unique fixed point

iff  $\rho(|I - RA|) < 1$  (contraction)

not easy to verify.

Retraction easier to verify

$$E_{n+1} \subset \overset{\circ}{E}_n$$

$\Rightarrow R$  and  $A$  not singular and  $e \in E_{n+1}$

$$\Rightarrow \hat{x} \in \tilde{x} + E_{n+1}$$

Choose  $R$  as an approximate inverse of  $A$ ;

then  $\rho(|I - RA|) < 1$  practically always holds.

$E_0$  is obtained by adding a small interval to  $\tilde{x}$ ;

then usually  $E_{n+1} \subset \overset{\circ}{E}_n$  after one or two steps.

(\*) is very sensitive towards roundings;

round as little as possible; apply opt. scalar product