

Rechner Architektur I (RAI)

Informationsverarbeitung

Teil 2

Prof. Dr. Akash Kumar
Chair for Processor Design

Rechenoperationen

Arithmetic Operations

Addition Beispiel

3

dezimal ($n = 2$)

$$\begin{array}{r} a \quad \quad 5 \ 9_{10} \\ b \quad + \quad 9 \ 1_{10} \\ \hline s \quad = 1^1 5^1 0_{10} \end{array}$$

dual ($n = 8$)

$$\begin{array}{r} \quad \quad 0 \ 0 \ 1 \ 1 \ 1 \ 0 \ 1 \ 1_2 \\ + 0 \ 1 \ 0 \ 1 \ 1 \ 0 \ 1 \ 1_2 \\ \hline = 1^1 0^1 0^1 1^1 0^0 1^1 1^1 0_2 \end{array}$$

Schreibweise: Summe^{Übertrag der vorherigen Stelle}

Sukzessive stellenweise Addition der Summanden mit dem LSB beginnend bis hin zum MSB ($a_0 + b_0, \dots, a_{n-1} + b_{n-1}$) zur Summe (s_0, \dots, s_{n-1}) mit Übertragstechnik. Die bei der stellenweisen Addition auftretenden Überträge (carry) (c_0, \dots, c_{n-1}) werden bei der Addition der jeweils nächstfolgenden Stelle mit dazu addiert. Der letzte Übertrag c_{n-1} heisst auslaufender Übertrag (carry out).

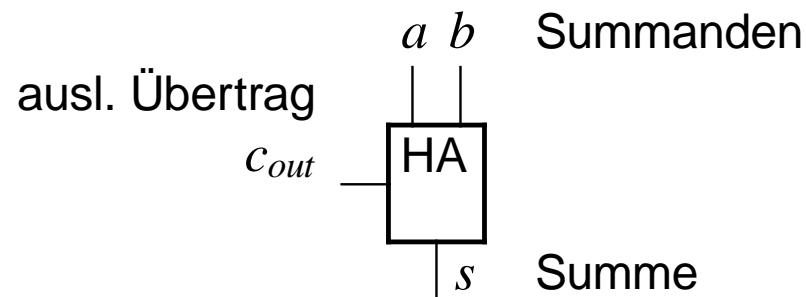
Hardware-Realisierung der Addition

4

1-Bit Addition

Halbaddierer (Half Adder (HA), ohne einlaufenden Übertrag)

Schaltsymbol



Modulo 2 Addition ohne Übertrag

Wertetabelle

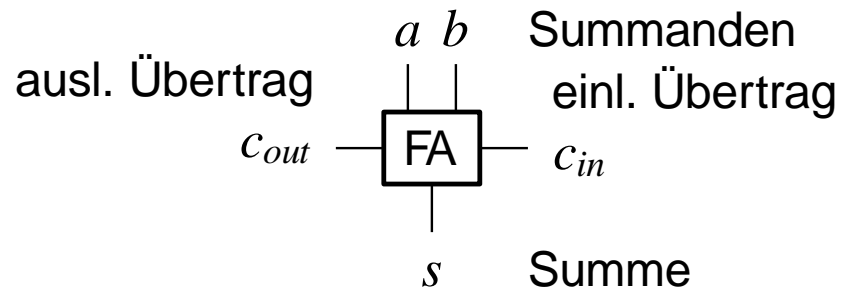
a	b	s	c_{out}
0	0	0	0
0	1	1	0
1	0	1	0
1	1	0	1

$$\begin{aligned} s &= a \wedge \bar{b} \vee \bar{a} \wedge b \\ \text{Schaltfunktion} &= a \oplus b \quad (a \text{ xor } b) \\ c_{out} &= a \wedge b \end{aligned}$$

Volladdierer (Full Adder (FA), mit einlaufenden Übertrag)

5

Schaltsymbol



Modulo 2 Addition mit Übertrag

Schaltfunktion

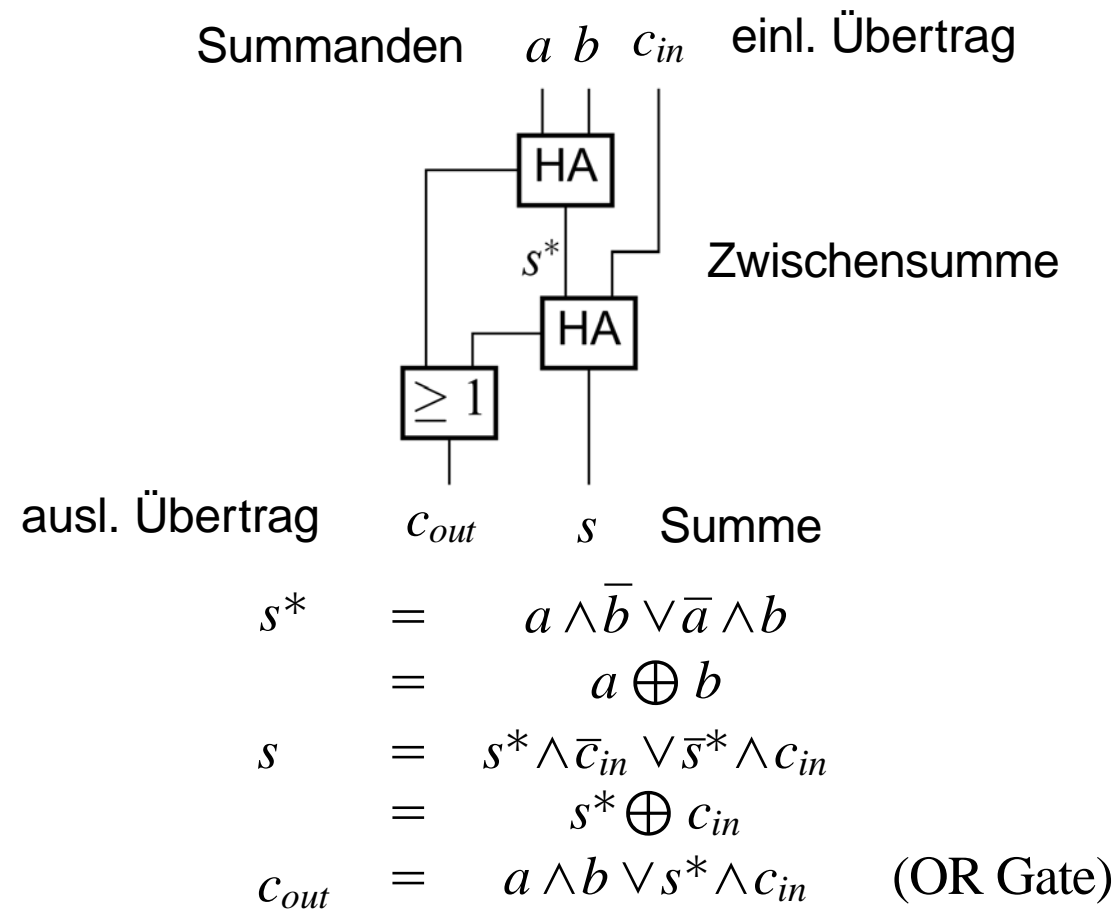
$$\begin{aligned}
 s &= \overline{a} \wedge b \wedge \overline{c_{in}} \vee a \wedge \overline{b} \wedge \overline{c_{in}} \vee \overline{a} \wedge \overline{b} \wedge c_{in} \vee a \wedge b \wedge c_{in} \\
 c_{out} &= a \oplus b \oplus c_{in} \\
 &= a \wedge b \vee a \wedge c_{in} \vee b \wedge c_{in}
 \end{aligned}$$

Wertetabelle

c_{in}	a	b	s	c_{out}
0	0	0	0	0
0	0	1	1	0
0	1	0	1	0
0	1	1	0	1
1	0	0	1	0
1	0	1	0	1
1	1	0	0	1
1	1	1	1	1

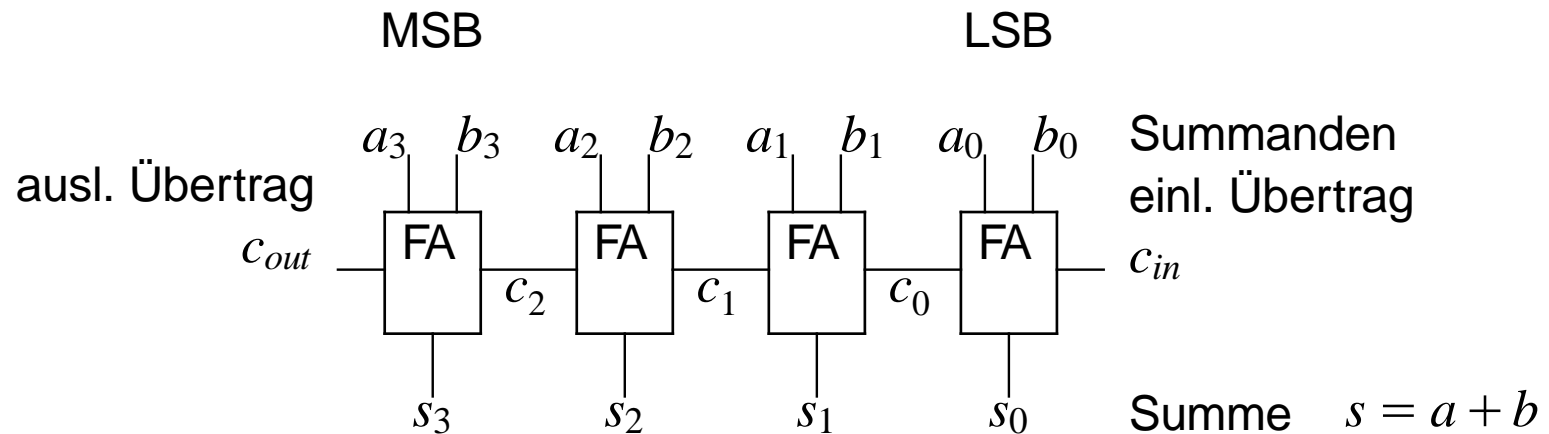
Volladdierer, realisiert durch Halbaddierer

6



Beispiel: 4-Bit Addition (Ripple-Carry-Adder (RCA))

7



Addition zweier n-Bit Zahlen durch stellenweise Addition mit Übertrag mit dem LSB beginnend und Weiterreichen des Übertrages an die jeweils nachfolgende Stelle.

$$s_v = a_v \oplus b_v \oplus c_{v-1} \quad \text{für } v = 0, 1, \dots, n-1$$

$$c_v = a_v \wedge b_v \vee a_v \wedge c_v \vee b_v \wedge c_v$$

$$c_{-1} = c_{in} \quad \text{kann auch bei } c_{in} = 0 \text{ entfallen}$$

$$c_{out} = c_{n-1}$$

Subtraktion (1) Beispiel

8

dezimal ($n = 2$)

$$\begin{array}{r} a \quad 9 \ 1_{10} \\ b \quad - 5 \ 9_{10} \\ \hline d \quad = 3^1 \ 2_{10} \end{array}$$

dual ($n = 8$)

$$\begin{array}{r} 0 \ 1 \ 0 \ 1 \ 1 \ 0 \ 1 \ 1_2 \\ - 0 \ 0 \ 1 \ 1 \ 1 \ 0 \ 1 \ 1_2 \\ \hline = 0^0 \ 0^1 \ 1^0 \ 0^0 \ 0^0 \ 0^0 \ 0^0 \ 0_2 \end{array}$$

Schreibweise: Differenz Übertrag (geborgt) der vorherigen Stelle

Sukzessive stellenweise Subtraktion mit dem LSB beginnend bis hin zum MSB mit Borgetechnik. Die bei der stellenweisen Subtraktion auftretenden (geborgten) Überträge (borrow) werden bei der Subtraktion der jeweils nächstfolgenden Stelle mit zum Subtrahenden addiert. Der letzte (geborgte) Übertrag c_{n-1} heisst auslaufender (geborgter) Übertrag (borrow out).

Subtraktion (2) Beispiel

9

dezimal ($n = 2$)

$$\begin{array}{r}
 a \quad 9 \ 1_{10} \\
 b \ (- \ 5 \ 9_{10} \) \\
 \hline
 {}^{(10)}b \ + \ 4 \ 1_{10} \\
 \hline
 d \ = \ 1 \ 3^0 \ 2_{10}
 \end{array}$$

dual ($n = 8$)

$$\begin{array}{r}
 0 \ 1 \ 0 \ 1 \ 1 \ 0 \ 1 \ 1_2 \\
 (- \ 0 \ 0 \ 1 \ 1 \ 1 \ 0 \ 1 \ 1_2 \) \\
 \hline
 + \ 1 \ 1 \ 0 \ 0 \ 0 \ 1 \ 0 \ 1_2 \\
 \hline
 = \ 1 \ 0^1 \ 0^0 \ 1^1 \ 0^1 \ 0^1 \ 0^1 \ 0_2
 \end{array}$$

Schreibweise: Summe^{Übertrag der vorherigen Stelle}

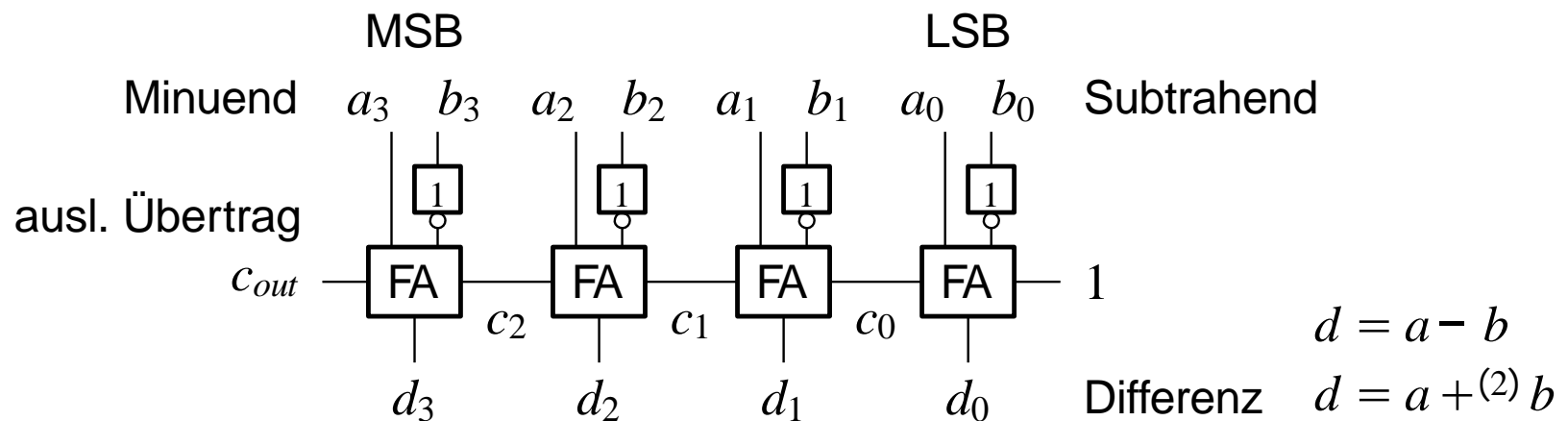
Rückführung der Subtraktion auf eine Addition mit dem B-Komplement des Subtrahenden bei fester Stellenzahl n : $d = a - b \rightarrow d = a + {}^{(2)}b$

Der dabei auftretende auslaufende Übertrag (carry out) c_{n-1} wird ignoriert.

Hardware-Realisierung der Subtraktion

10

Beispiel: 4-Bit Subtraktion (Ripple-Carry-Adder mit 2-Komplement)



$$d_v = a_v \oplus \bar{b}_v \oplus c_{v-1} \quad \text{für } v = 0, 1, \dots, n-1$$

$$c_v = a_v \wedge \bar{b}_v \vee a_v \wedge c_v \vee \bar{b}_v \wedge c_v$$

$$c_{-1} = 1 \quad \text{für 2-Komplement}$$

$$C_{out} = c_{n-1}$$

Multiplikation (1) Beispiel

11

dezimal ($n = 2$)

$$\begin{array}{r}
 1 \ 5_{10} \cdot 1 \ 3_{10} \\
 \hline
 4^1 \ 5 \\
 + 1^0 \ 5 \\
 \hline
 = 1^0 \ 9^0 \ 5_{10}
 \end{array}$$

dual ($n = 4$)

$$\begin{array}{r}
 1 \ 1 \ 1 \ 1_2 \cdot 1 \ 1 \ 0 \ 1_2 \\
 \hline
 1 \ 1 \ 1 \ 1 \\
 0 \ 0 \ 0 \ 0 \\
 1 \ 1 \ 1 \ 1 \\
 + 1 \ 1 \ 1 \ 1 \\
 \hline
 = 1 \ 1^{10} \ 0^{10} \ 0^{10} \ 0^1 \ 0^0 \ 1^0 \ 1_2
 \end{array}$$

Sukzessive stellenweise Multiplikation des Multiplikators (0 oder 1) mit dem LSB beginnend bis hin zum MSB mit dem gesamten Multiplikanten und stellenrichtige Aufsummierung der Teilergebnisse (Linksverschiebung um $n - 1$ Stellen):
 $(a \cdot b = a \cdot b_0 \cdot 2^0 + a \cdot b_1 \cdot 2^1 + \dots + a \cdot b_{n-1} \cdot 2^{n-1}).$

Multiplikation (2)

12

$$\begin{aligned}a \cdot b &= p \\ b &= b_0 + b_1 \cdot 2^1 + b_2 \cdot 2^2 + \dots + b_{n-1} \cdot 2^{n-1} \\ p &= a \cdot b_0 + a \cdot b_1 \cdot 2^1 + a \cdot b_2 \cdot 2^2 + \dots + a \cdot b_{n-1} \cdot 2^{n-1} \\ b_v &\in \{0, 1\}\end{aligned}$$

Die Multiplikation kann so auf eine fortlaufende Addition und eine entsprechende Stellenverschiebung zurückgeführt werden.

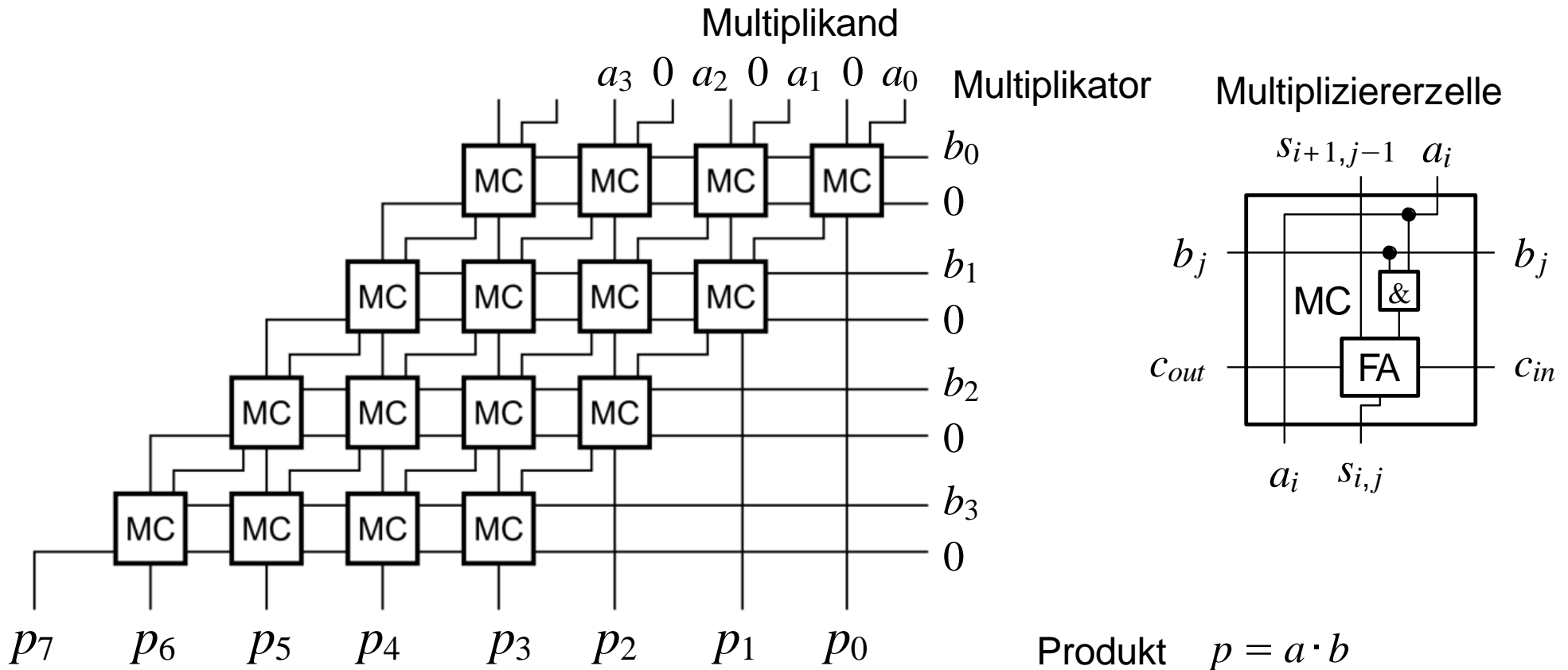
Eine Multiplikation mit 2 bedeutet eine Verschiebung des Multiplikanten um eine Stelle nach links, wobei das LSB mit 0 aufgefüllt wird (analog für eine Multiplikation mit 2^k eine Verschiebung um k Stellen).

$$a \cdot 2^k \rightarrow (a_l = a_{l-k}, \text{ für } l = n-1, \dots, k) \text{ und } (a_l = 0, \text{ für } l = k-1, \dots, 0)$$

Hardware-Realisierung der Multiplikation

13

Beispiel: 4-Bit Multiplikation (Ripple-Carry-Addition mit 1-Bit Multiplikation)



Division Beispiel

14

dezimal ($n = 2$)

$$\begin{array}{r} 59_{10} : 3_{10} = 19_{10} \\ - 3 \\ \hline 29 \\ - 27 \\ \hline 2 \end{array}$$

dual ($n = 8$) ohne führende 0

$$\begin{array}{r} 111011_2 : 11_2 = _2 \\ - 11 \\ \hline 01011 \rightarrow 11 \leq 11 \rightarrow 1 \uparrow \\ 11 \rightarrow 11 > 01 \rightarrow 0 \uparrow \\ 11 \rightarrow 11 > 010 \rightarrow 0 \uparrow \\ - 11 \rightarrow 11 < 0101 \rightarrow 1 \uparrow \\ \hline 101 \\ - 11 \rightarrow 11 < 101 \rightarrow 1 \uparrow \\ \hline 10 \end{array}$$

Division (2)

15

$$\begin{aligned}a : b &= q \\ q &= q_0 + q_1 \cdot 2^1 + q_2 \cdot 2^2 + \dots + q_{n-1} \cdot 2^{n-1} \\ a &= b \cdot q \\ a &= b \cdot q_0 + b \cdot q_1 \cdot 2^1 + b \cdot q_2 \cdot 2^2 + \dots + b \cdot q_{n-1} \cdot 2^{n-1} \\ q_v &\in \{0, 1\}\end{aligned}$$

Bei der Division wird stellenweise mit dem MSB von q beginnend das Produkt $b \cdot [q_v] \cdot 2^v$ von a (durch Verschiebung) subtrahiert. Ist die Differenz positiv, so ist $q_v = 1$, sonst 0. Mit der verbleibenden Differenz von a wird analog weiter verfahren bis alle Stellen abgearbeitet sind. Als letzte Differenz bleibt der Divisionsrest übrig.

Die Division kann so auf eine fortlaufende Stellenverschiebung und eine entsprechende Subtraktion (Addition mit 2-Komplement) zurückgeführt werden.

Division (3)

16

Eine Division durch 2 bedeutet eine Verschiebung des Dividenden um eine Stelle nach rechts, wobei das MSB mit 0 aufgefüllt wird (analog für eine Division durch 2^k eine Verschiebung um k Stellen).

$$a/2^k \rightarrow (a_{l-k} = a_l, \text{ für } l = k, \dots, n-1) \text{ und } (a_l = 0, \text{ für } l = n-1, \dots, n-k-1)$$

Die Division ist in Hardware aufwendiger zu realisieren als die Multiplikation und benötigt daher allgemein auch mehr Rechenzeit bei der Abarbeitung.

Eine Alternative zu diesem Divisionsverfahren stellen verschiedene iterative Verfahren dar, die die Division ebenfalls auf Stellenverschiebung, Multiplikation und Addition zurückführen.

Division (4)

17

Iterative Division mit dem Newton-Verfahren

$$q = \frac{a}{b}$$

$$f(x) = \frac{1}{x} - b \quad \text{Ansatzfunktion mit Nullstelle bei } x = \frac{1}{b}$$

$$x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)} \quad \text{Iterationsvorschrift nach Newton}$$

$$f'(x_i) = \frac{-1}{x_i^2} \quad \text{Ableitung von } f(x)$$

$$x_{i+1} = x_i \cdot (2 - x_i \cdot b) \quad \text{Iterationsvorschrift, } x_{i+1} - \text{verbesserte Näherung}$$

$$\frac{a}{b} = a \cdot x_n \quad \text{für } n \rightarrow \infty$$

x_0 - Startwert, x_i - Näherungswert nach i Iterationen

(Alternatives Verfahren: Iterative Division mit dem Goldschmidt-Algorithmus)

Division (5)

18

Beispiel: $q = \frac{\bar{a}}{b} = \frac{\bar{15}}{7} = 15 \cdot \frac{\bar{1}}{7} = 2.142857142857 \dots$

Startwert: $x_0 = 0.1$

$x_{i+1} =$	$x_i \cdot (2 - x_i \cdot b)$	Genauigkeit
$x_1 =$	$0.1 \cdot (2 - 0.1 \cdot 7) = \underline{0.13}$	10^{-1}
$x_2 =$	$0.13 \cdot (2 - 0.13 \cdot 7) = \underline{0.1417}$	10^{-2}
$x_3 =$	$0.1417 \cdot (2 - 0.1417 \cdot 7) = \underline{0.14284777}$	10^{-4}
$x_4 =$	$\dots = \underline{0.142857142242}$	10^{-8}
$x_5 =$	$\dots = \underline{0.142857142857} \dots$	$< 10^{-10}$

Lösung: $q = 15 \cdot x_5 = 2.142857142857 \dots$

Fest/Gleit-kommazahlen

Fixed/Floating-point numbers



Zahlendarstellung im Computer, Datenformate

20

Die Zahlendarstellung im Computer (Datenformat) erfolgt in der Regel binär codiert in Formaten fester Länge n . Es wird das Dualzahlen-Stellenwertsystem verwendet. Es gelten die Rechenregeln für Dualzahlen.

Folgende Darstellungsvarianten werden allgemein unterschieden:

- vorzeichenlose natürliche Zahlen
- vorzeichenbehaftete ganze Zahlen
- Festkommazahlen (Festpunktzahlen)
- Gleitkommazahlen (Gleitpunktzahlen)

Durch die Darstellungsvariante selbst und die begrenzte Stellenzahl (Datenformate fester Länge n) sind die Gültigkeitsbereiche des Assoziativgesetzes und des Distributivgesetzes zu beachten.

Gebrochene Zahlen im Stellenwertsystem

21

Die Zifferndarstellung einer allgemeinen gebrochenen Zahl G_B lautet im Stellenwertsystem der Basis B (Festkommadarstellung):

$$G_B = \pm(z_{n-1} \dots z_1 z_0, z_{-1} z_{-2} \dots z_{-m})_B = \pm((z_{n-1} \dots z_1 z_0)_B + (0, z_{-1} z_{-2} \dots z_{-m})_B)$$

Der Wert von G_B bestimmt sich durch :

$$G_B = \pm \sum_{i=-m}^{n-1} z_i \cdot B^i = \sum_{i=-m}^{-1} z_i \cdot B^i + \sum_{i=0}^{n-1} z_i \cdot B^i$$

$$\begin{aligned} G_B = \pm & (((((z_{n-1} \cdot B + z_{n-2}) \cdot B \dots) \cdot B + z_1) \cdot B + z_0) \quad \leftarrow \text{ganzzahlig} \\ & + ((z_{-m} \cdot B^{-1} + z_{-m+1}) \cdot B^{-1} \dots) \cdot B^{-1} + z_{-1}) \cdot B^{-1} \quad \leftarrow \text{echt gebrochen} \end{aligned}$$

Gebrochene Zahlen im Stellenwertsystem

22

Die Zifferndarstellung einer allgemeinen gebrochenen Zahl G_B lautet im Stellenwertsystem der Basis B (Gleitkommadarstellung):

$$G_B = \pm(0, z_{-1}z_{-2} \dots z_{-m})_B \cdot B^{(\pm e_{k-1} \dots e_1 e_0)_B} \quad \text{mit} \quad M_B = \sum_{i=1}^m z_{-i} B^{-i}, \quad E_B = \sum_{j=0}^{k-1} e_j B^j.$$

Der Wert von G_B bestimmt sich durch :

$$G_B = \pm \left(\sum_{i=1}^m z_{-i} \cdot B^{-i} \right) \cdot B^{(\pm \sum_{j=0}^{k-1} e_j B^j)} = \pm M_B \cdot B^{(\pm E_B)}$$

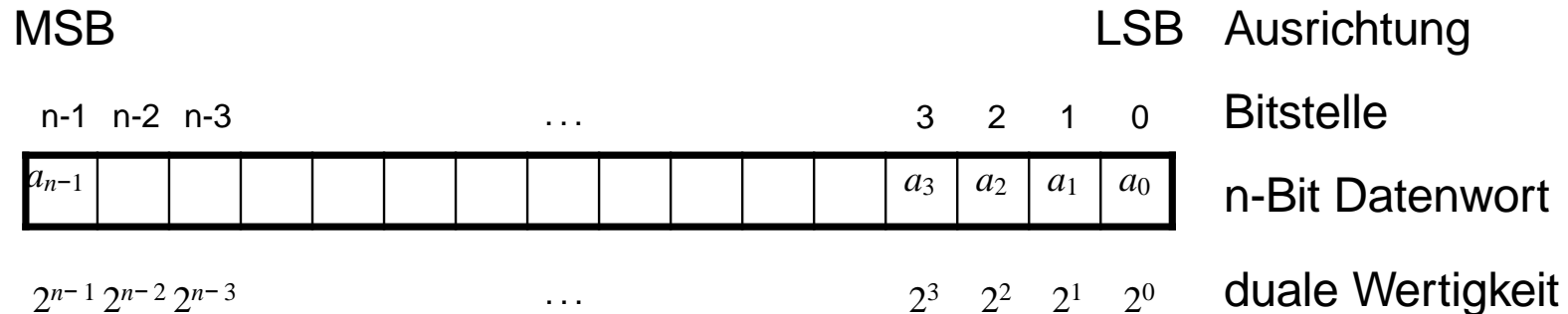
M_B - ist eine echt gebrochene Zahl und wird als Mantisse bezeichnet. (Hidden bit?)

E_B - ist eine ganze Zahl und wird als Exponent bezeichnet.

Der Exponent gibt an, um wieviele Stelle das Komma nach rechts (positiv) bzw. nach links (negativ) verschoben werden muss.

N-Bit Datenformat (Dualzahlen-Stellenwertsystem)

23



Wert Z bei Interpretation als reine Dualzahl:

$$Z = a_{n-1} \cdot 2^{n-1} + a_{n-2} \cdot 2^{n-2} + \dots + a_2 \cdot 2^2 + a_1 \cdot 2^1 + a_0 \cdot 2^0$$

Übliche Datenformate: 8 Bit, 16 Bit, 32 Bit, 64 Bit, ...

Je nach Datenformat haben einzelne Bitstellen spezielle Funktionen, z.B. Vorzeichen, Exponent, ...

Jede Zahlendarstellung hat daher ein gesondertes Datenformat.

Festkommazahlen (fixed-point numbers)

24

Festkommazahlen werden zur Darstellung von gebrochenen Zahlen verwendet:

- Die Darstellung analog der Vorzeichen-Wert- oder 2-Komplement-Darstellung.
- Die Anzahl der Nachkommastellen wird vorher fest definiert (f -fracional part).
- Das Komma steht implizit immer an der selben Stelle (nicht abgespeichert).
- Die Zahl der Vorkommastellen ergibt sich aus der Differenz von Gesamtstellenzahl und Nachkommastellenzahl (Vorzeichen im MSB).

n-stellige Festkommazahl mit m Nachkommastellen:

$$z = d, f \rightarrow d_{n-1}d_{n-2} \dots d_{n-m}f_{m-1}f_{m-2} \dots f_0$$

kleinste Zahl für $|z| > 0$: 2^{-m} größte Zahl für $|z|$: $(2^{n-1} - 1) \cdot 2^{-m}$

Anwendung: Real-Arithmetik

Konvertierung von Festkommazahlen (1)

25

Beispiel: Konvertierung einer Hexadezimalzahl $X_{16} = A45,8F$ (2 hexadezimale Nachkommastellen) in die wertgleiche Dezimalzahl X_{10}

$$X_{16} = A \quad 4 \quad 5, \quad 8 \quad F_{16}$$

$$X_{10} = A_{16} \cdot 16^2 + 4_{16} \cdot 16^1 + 5_{16} \cdot 16^0 + 8_{16} \cdot 16^{-1} + F_{16} \cdot 16^{-2}$$

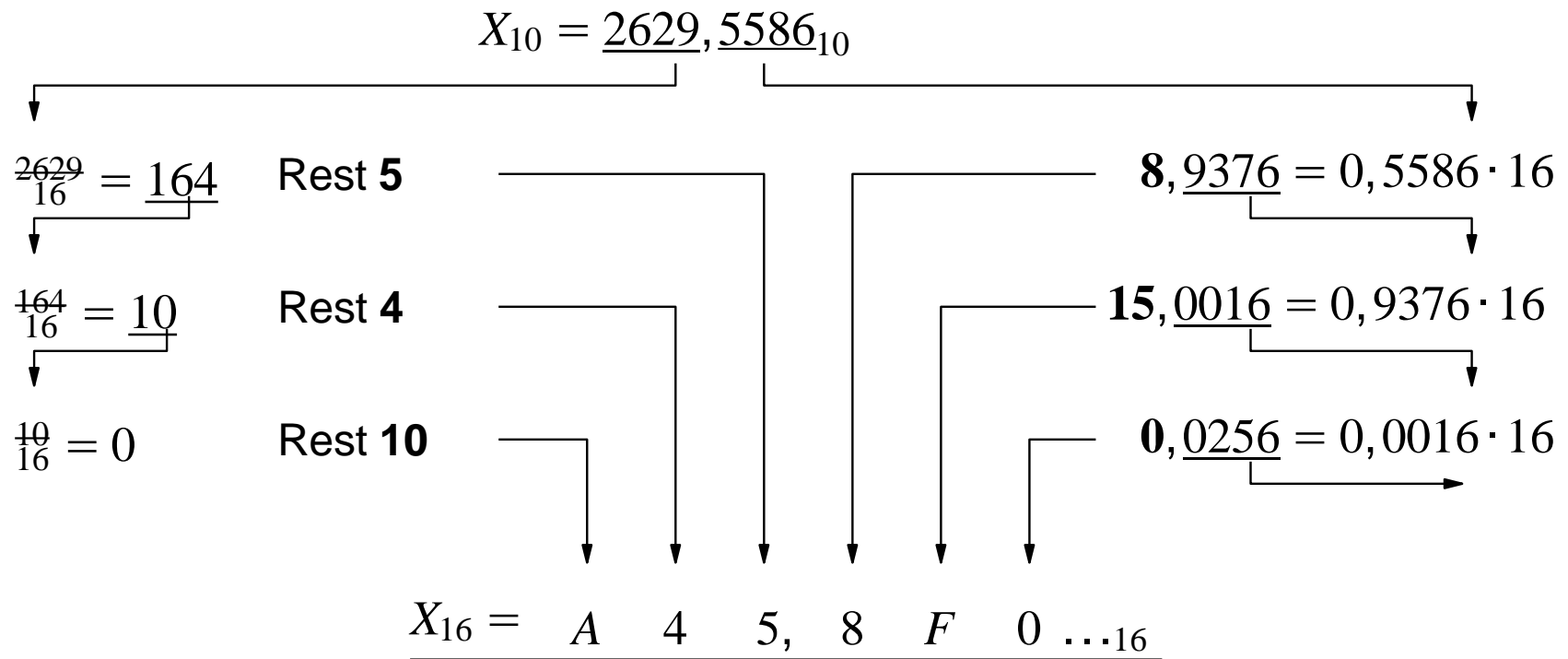
$$X_{10} = 10_{10} \cdot 16^2 + 4_{10} \cdot 16^1 + 5_{10} \cdot 16^0 + 8_{10} \cdot 16^{-1} + 15_{10} \cdot 16^{-2}$$

$$\underline{X_{10} = 2629,5586 \dots_{10}}$$

Konvertierung von Festkommazahlen (2)

26

Beispiel: Konvertierung einer Dezimalzahl $X_{10} = 2629,5586$ (4 dezimale Nachkommastellen) in die wertgleiche Hexadezimalzahl X_{16}



Gleitkommazahlen (floating-point numbers)

27

Halblogarithmische Darstellung gebrochener Zahlen durch Mantisse und Exponent:

$$z = M \cdot 2^E$$

- V - Vorzeichen der Mantisse
- M - Mantisse (Signifikant S)
- E - Exponent (Charakteristik)

Vorteil: wesentlich größerer Wertebereich und flexibler Nachkommabereich

Nachteil: geringere relative Genauigkeit

Standard: ANSI/IEEE 754 von 1985

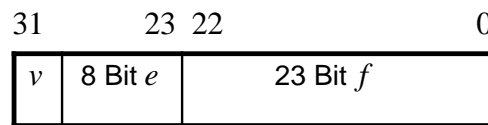
(Institutes for Electrical and Electronic Engineers)

Anwendung: Real-Arithmetik

Datenformate für Gleitkommazahlen

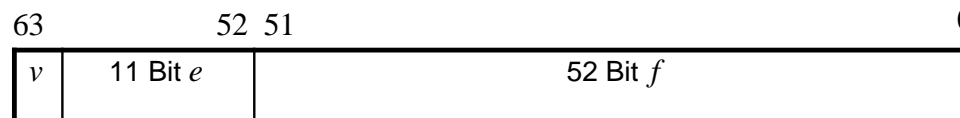
28

32-Bit Format (single)

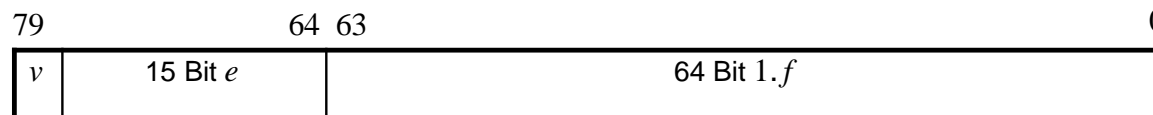


v -Vorzeichen
 e -biased exponent
 f -fractional part

64-Bit Format (double)



80-Bit Format (double-extended)



Mantissendarstellung (1)

29

Darstellung der Mantisse M als Festkommazahl (vorzeichenbehaftete gebrochene Zahl in Vorzeichen-Wert-Darstellung) mit Vorzeichen im MSB.

- **normalisierte Form:** Die führende 1 und das Komma werden nicht explizit gespeichert, sondern nur der gebrochene Anteil f (fractional part), (Ausnahme: double-extended). → Vor der Verarbeitung ist die 1 zu ergänzen (entpacken). Das Komma steht immer rechts neben der 1.

$$M = (-1)^v \cdot 1, f$$

- **denormalisierte Form:** Die führende 1 wird durch 0 ersetzt. Dient zur Darstellung sehr kleiner Zahlen (Kennzeichnung im Exponenten) oder zur Exponentenanpassung bei Berechnungen. → geringere Genauigkeit als normalisierte Darstellung.

$$M = (-1)^v \cdot 0, f$$

Mantissendarstellung (2)

30

Wertebereich: $1,0 \leq |M| < 2,0$ mit $0,0 \leq f < 1,0$

Genauigkeit der Mantisse (binär - dezimal):

Datenformat	Binärstellen	max. rel. Genauigkeit	Dezimalstellen
single	23	$2^{-24} = 5,96 \dots 10^{-8}$	7
double	52	$2^{-53} = 1,11 \dots 10^{-16}$	15
double-extended	64	$2^{-64} = 5,42 \dots 10^{-20}$	19

Exponentendarstellung

31

Darstellung des Exponenten E als vorzeichenbehaftete ganze Zahl in Basiswert-Darstellung) mit dem Basiswert B_e .

$$E = e - B_e \quad \text{bzw.} \quad e = E + B_e$$

Wertebereich (ohne Ausnahmesituationen):

Datenformat	b_e	B_e	Wertebereich
single	8	127	$-126 \leq E \leq 127$
double	11	1023	$-1022 \leq E \leq 1023$
double-extended	15	16383	$-16382 \leq E \leq 16383$

(b_e - Binärstellen (Exponent), B_e - Basiswert (biased))

Ausnahmesituationen der Gleitkommadarstellung (1)

32

Größter ($e = 2^{b_e} - 1$) und kleinster ($e = 0$) Exponentenwert sind für Ausnahmesituationen reserviert (\rightarrow 4 Darstellungsvarianten):

Interpretation	V	e	f	Zahlendarstellung
Null	v	0	0	$(-1)^v \cdot 0$
1. denormalisiert	v	0	$\neq 0$	$(-1)^v \cdot 0, f \cdot 2^{(1-B_e)}$
2. normalisiert	v	$0 < e < 2^{b_e} - 1$	$\neq 0$	$(-1)^v \cdot 1, f \cdot 2^{(e-B_e)}$
3. Unendlich	v	$2^{b_e} - 1$	0	$(-1)^v \cdot \infty$
4. Not-a-Number	v	$2^{b_e} - 1$	$\neq 0$	NaN

- zwei verschiedene Darstellungen der Null (+0 und -0).
- die Darstellung der Null erfolgt als denormalisierte Zahl.
- denormalisierte Zahlen werden mit $e = 1$ interpretiert.
- Bereichsüberschreitungen sind nicht möglich.
- Sonderdarstellungen für Unendlich, ($-\infty$ und $+\infty$).

Ausnahmesituationen der Gleitkommadarstellung (2)

33

Beispiele (single):

v	e	f	Zahlenwert	Interpretation
0	00000000	00 ... 00	+0	pos. Null
1	00000000	00 ... 00	-0	neg. Null
0	00000000	00 ... 01	$+0,00 \dots 01 \cdot 2^{-126}$	denormalisiert
1	00000000	00 ... 01	$-0,00 \dots 01 \cdot 2^{-126}$	denormalisiert
0	00000001	00 ... 00	$+1,00 \dots 00 \cdot 2^{-126}$	normalisiert
1	00000001	00 ... 00	$-1,00 \dots 00 \cdot 2^{-126}$	normalisiert
0	11111110	11 ... 11	$+1,11 \dots 11 \cdot 2^{+127}$	normalisiert
1	11111110	11 ... 11	$-1,11 \dots 11 \cdot 2^{+127}$	normalisiert
0	11111111	00 ... 00	$+\infty$	pos. unendlich
1	11111111	00 ... 00	$-\infty$	neg. unendlich
0	11111111	00 ... 01	NaN	Not-a-Number

Addition und Subtraktion von Gleitkommadarstellung (1)

34

$$z_1 = M_1 \cdot 2^{E_1} = (-1)^{v_1} \cdot 1, f_1 \cdot 2^{(e_1 - B_e)}$$

$$z_2 = M_2 \cdot 2^{E_2} = (-1)^{v_2} \cdot 1, f_2 \cdot 2^{(e_2 - B_e)}$$

$$\begin{aligned} z_1 + z_2 &= (M_1 + M_2 \cdot 2^{-(E_1 - E_2)}) \cdot 2^{E_1} \\ &= ((-1)^{v_1} \cdot 1, f_1 + (-1)^{v_2} \cdot 1, f_2 \cdot 2^{-(e_1 - e_2)}) \cdot 2^{(e_1 - B_e)} \end{aligned}$$

Bei ungleichen Exponenten e_1, e_2 muß die Zahl mit dem kleineren Exponenten denormalisiert werden (Erhöhung des Exponenten bei gleichzeitiger Rechtsverschiebung des Signifikanten).

Addition und Subtraktion von Gleitkommadarstellung (2)

35

Algorithmus für vorzeichenbehaftete Addition $z_s = z_1 + z_2$:

$$z_s = (M_1 + M_2 \cdot 2^{-(e_1 - e_2)}) \cdot 2^{E_1}$$

1. Wenn $e_1 < e_2$, dann vertausche z_1 und z_2 .
2. Verschiebe M_2 um $e_1 - e_2$ Stellen nach rechts.
3. Addiere bzw. Subtrahiere M_1 und M_2 analog zu vorzeichenbehafteten ganzen Zahlen (Zweierkomplement). Bei negativem Ergebnis Konvertierung in die Vorzeichen-Wert-Darstellung.
4. Normalisierung und Rundung des Ergebnisses.

Addition und Subtraktion von Gleitkommadarstellung (3)

36

Beispiel (single): (⁰ - Zusatzstelle wegen 2-Komplement Rechnung)

z	v	e	s	Bemerkung	Dez.
z_1	0	1000 0001	⁰ 1,0000...	normalisiert	4,0
z_2	0	0111 1111	⁰ 1,1000...	normalisiert	1,5
z_2	0	1000 0001	⁰ 0,0110...	denormalisiert	1,5
$z_1 + z_2$	0	1000 0001	⁰ 1,0110...	normalisiert	5,5
$-z_2$	1	1000 0001	¹ 1,1010... ...	denormalisiert 2-Komplement	-1,5
$z_1 - z_2$	0	1000 0001	⁰ 0,1010...	denormalisiert	2,5
$z_1 - z_2$	0	1000 0000	⁰ 1,0100...	normalisiert	2,5

Multiplikation und Division von Gleitkommazahlen (1)

37

$$z_1 = M_1 \cdot 2^{E_1} = (-1)^{v_1} \cdot 1, f_1 \cdot 2^{(e_1 - B_e)}$$

$$z_2 = M_2 \cdot 2^{E_2} = (-1)^{v_2} \cdot 1, f_2 \cdot 2^{(e_2 - B_e)}$$

$$z_1 \cdot z_2 = M_1 \cdot M_2 \cdot 2^{(E_1 + E_2)}$$

$$= (-1)^{v_1} \cdot (-1)^{v_2} \cdot 1, f_1 \cdot 1, f_2 \cdot 2^{((e_1 + e_2 - B_e) - B_e)}$$

$$\frac{z_1}{z_2} = \frac{M_1}{M_2} \cdot 2^{(E_1 - E_2)}$$

$$= \frac{(-1)^{v_1}}{(-1)^{v_2}} \cdot \frac{1, f_1}{1, f_2} \cdot 2^{((e_1 - e_2 + B_e) - B_e)}$$

Multiplikation und Division von Gleitkommazahlen (2)

38

Algorithmus für $z_p = z_1 \cdot z_2$ bzw. $\frac{z_1}{z_2}$:

$$z_p = M_1 \cdot M_2 \cdot 2^{((e_1 + e_2 - B_e) - B_e)} \text{ bzw. } z_p = M_1 / M_2 \cdot 2^{((e_1 - e_2 + B_e) - B_e)}$$

1. Bestimmung des Vorzeichens ($v_1 \neq v_2 \rightarrow$ negativ, $v_1 = v_2 \rightarrow$ positiv).
2. Multiplikation bzw. Division von M_1 und M_2 analog zu ganzen Zahlen.
3. Addition von e_1 und e_2 und Subtraktion von B_e bzw. Subtraktion von e_1 und e_2 und Addition von B_e analog zu vorzeichenbehafteten ganzen Zahlen.
4. Normalisierung und Rundung des Ergebnisses.

Multiplikation und Division von Gleitkommazahlen (3)

39

Beispiel (single):

z	v	e	s	Bemerkung	Dez.
z_1	0	1000 0001	1,0000...	normalisiert	4,0
z_2	0	0111 1111	1,1000...	normalisiert	1,5
$z_1 \cdot z_2$	0	1000 0001	1,1000...	normalisiert	6,0
$\frac{z_2}{z_1}$	0	0111 1101	1,1000...	normalisiert	0,375

Rundung von Gleitkommazahlen (1)

40

Der IEEE 754 - Standard fordert:

Das Ergebnis einer arithmetischen Operation soll das selbe sein, als würde exakt gerechnet und anschließend gerundet (Runden zum nächsten Wert, die Hälfte auf den geraden Wert (round-to-even)).

Verschiedene Rundungsmodi:

- nach 0,
- nach $+\infty$ oder $-\infty$,
- zum nächsten Wert, bei 0,5 auf den geraden Wert (round-to-even)

Relevante Bitstellen für die Rundung (auf q Stellen, vom MSB gezählt):

- guard Bit g - Bitstelle auf die gerundet wird (q Stelle vom MSB),
- round Bit r - Rundungsbit (q+1 Stelle vom MSB),
- sticky Bit s - Zusatzstellen (q+2, q+3, ... Stellen vom MSB bis zum LSB)

Rundung von Gleitkommazahlen (2)

41

Beispiel (Rundung von $n = 8$ -Bit Zahlen auf $q = 4$ Stellen):

Nr.	n-Bit Zahl			gerundet
1	1, 0 0 0	1 0 0 0	→	1, 0 0 0
2	1, 0 0 0	1 0 0 1	→	1, 0 0 1
3	1, 0 0 1	1 0 0 0	→	1, 0 1 0
4	1, 0 0 1	1 1 1 1	→	1, 0 1 0
5	1, 1 1 1	1 0 0 0	→	1 0, 0 0
		$g \mid r \text{ --- } s \text{ ---}$		

Ausschlaggebend für die Rundung ist, ob in den s -Bits (restlichen Stellen) überhaupt eine 1 steht, nicht die Anzahl. Für eine genaue Rundung sind mehrere s -Bits erforderlich (z.B. $s = 3$).

Probleme beim Rechnen mit Gleitkommazahlen

42

- Berechnungen mit denormalisierten Zahlen ($e = 0$) (Ungenauigkeit),
- Überlauf (Wertebereichsüberschreitung),
- Unterlauf (begrenzte Darstellung betragsmässig kleiner Werte),
- Rundung.

Assoziativ- und Distributivgesetze gelten nur eingeschränkt:

Ursachen für die Differenzen sind die endliche Stellenzahl zur Zahlendarstellung und die Rundung der Ergebnisse.

$$a + (b + c) \neq (a + b) + c$$

$$a \cdot (b \cdot c) \neq (a \cdot b) \cdot c$$

$$(a + b) \cdot c \neq (a \cdot c) + (b \cdot c)$$

Zusammenfassung

43

- Floating point representation common for large numbers
- Hardware requirements for floating point processing – especially Division – can be high
 - one needs to evaluate the trade-offs involved in having one