



# 1 Zahlendarstellungen

---

Technische Grundlagen der Informatik 2  
(Rechnertechnologie 2)  
SS 2006

Andreas Koch

FG Eingebettete Systeme und ihre Anwendungen

Auf Basis von Material von  
R. Hoffmann

FG Rechnerarchitektur  
Technische Universität Darmstadt



# 1. Zahlendarstellungen

---

2

1.1 Zahl zur Basis  $b$

1.2 Vorzeichen-Betrag-Darstellung

1.3 Zweikomplementzahl

1.4 Einkomplementzahl

1.5 Binärcodierte Dezimalzahl

1.6 Gleitkommazahl


IEEE-32-Bit-Gleitkommaformat

IEEE-64-Bit-Gleitkommaformat

3.4 Konvertierung zwischen Zahlensystemen

# Zahlen

- [Abstrakte] Zahl = der **Wert** einer Zahl
- [Konkrete] Zahl = (**Darstellung**, **Wert**)

Darstellung				Wert
	/////	101		5
	zwei	iki (türkisch)		2

Acht Cova Zeichen des chinesischen Königs Fo-Hi vor mehr als 4000 Jahren

# Leibniz 1679

15 Martii 1679.

## De Progressione Dyadica - Pars I.

Numeratio															
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
1	10	11	100	101	110	111	1000	1001	1010	1011	1100	1101	1110	1111	10000
17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32
10001	10010	10011	10100	10101	10110	10111	11000	11001	11010	11011	11100	11101	11110	11111	100000

Ad iunctam progressionem facile continuari potest. ut patet a dextra finis horum, quibus subscribitur 0. donec occurrat in superstante etiam 0. cuius nullum est in alia progressione. Ita ex 1011000 esse  $2^6 + 2^4 + 2^3 + 2^2 + 2^1 + 2^0 = 64 + 16 + 8 = 88$ .  
 Nam 1 in quarta loco seu 1000 significat cubum fundamenti progressionis, et in communis progressionis significat cubum. Ita in nostra cubum a binario semper 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100, 101, 102, 103, 104, 105, 106, 107, 108, 109, 110, 111, 112, 113, 114, 115, 116, 117, 118, 119, 120, 121, 122, 123, 124, 125, 126, 127, 128, 129, 130, 131, 132, 133, 134, 135, 136, 137, 138, 139, 140, 141, 142, 143, 144, 145, 146, 147, 148, 149, 150, 151, 152, 153, 154, 155, 156, 157, 158, 159, 160, 161, 162, 163, 164, 165, 166, 167, 168, 169, 170, 171, 172, 173, 174, 175, 176, 177, 178, 179, 180, 181, 182, 183, 184, 185, 186, 187, 188, 189, 190, 191, 192, 193, 194, 195, 196, 197, 198, 199, 200, 201, 202, 203, 204, 205, 206, 207, 208, 209, 210, 211, 212, 213, 214, 215, 216, 217, 218, 219, 220, 221, 222, 223, 224, 225, 226, 227, 228, 229, 230, 231, 232, 233, 234, 235, 236, 237, 238, 239, 240, 241, 242, 243, 244, 245, 246, 247, 248, 249, 250, 251, 252, 253, 254, 255, 256, 257, 258, 259, 260, 261, 262, 263, 264, 265, 266, 267, 268, 269, 270, 271, 272, 273, 274, 275, 276, 277, 278, 279, 280, 281, 282, 283, 284, 285, 286, 287, 288, 289, 290, 291, 292, 293, 294, 295, 296, 297, 298, 299, 300, 301, 302, 303, 304, 305, 306, 307, 308, 309, 310, 311, 312, 313, 314, 315, 316, 317, 318, 319, 320, 321, 322, 323, 324, 325, 326, 327, 328, 329, 330, 331, 332, 333, 334, 335, 336, 337, 338, 339, 340, 341, 342, 343, 344, 345, 346, 347, 348, 349, 350, 351, 352, 353, 354, 355, 356, 357, 358, 359, 360, 361, 362, 363, 364, 365, 366, 367, 368, 369, 370, 371, 372, 373, 374, 375, 376, 377, 378, 379, 380, 381, 382, 383, 384, 385, 386, 387, 388, 389, 390, 391, 392, 393, 394, 395, 396, 397, 398, 399, 400, 401, 402, 403, 404, 405, 406, 407, 408, 409, 410, 411, 412, 413, 414, 415, 416, 417, 418, 419, 420, 421, 422, 423, 424, 425, 426, 427, 428, 429, 430, 431, 432, 433, 434, 435, 436, 437, 438, 439, 440, 441, 442, 443, 444, 445, 446, 447, 448, 449, 450, 451, 452, 453, 454, 455, 456, 457, 458, 459, 460, 461, 462, 463, 464, 465, 466, 467, 468, 469, 470, 471, 472, 473, 474, 475, 476, 477, 478, 479, 480, 481, 482, 483, 484, 485, 486, 487, 488, 489, 490, 491, 492, 493, 494, 495, 496, 497, 498, 499, 500, 501, 502, 503, 504, 505, 506, 507, 508, 509, 510, 511, 512, 513, 514, 515, 516, 517, 518, 519, 520, 521, 522, 523, 524, 525, 526, 527, 528, 529, 530, 531, 532, 533, 534, 535, 536, 537, 538, 539, 540, 541, 542, 543, 544, 545, 546, 547, 548, 549, 550, 551, 552, 553, 554, 555, 556, 557, 558, 559, 560, 561, 562, 563, 564, 565, 566, 567, 568, 569, 570, 571, 572, 573, 574, 575, 576, 577, 578, 579, 580, 581, 582, 583, 584, 585, 586, 587, 588, 589, 590, 591, 592, 593, 594, 595, 596, 597, 598, 599, 600, 601, 602, 603, 604, 605, 606, 607, 608, 609, 610, 611, 612, 613, 614, 615, 616, 617, 618, 619, 620, 621, 622, 623, 624, 625, 626, 627, 628, 629, 630, 631, 632, 633, 634, 635, 636, 637, 638, 639, 640, 641, 642, 643, 644, 645, 646, 647, 648, 649, 650, 651, 652, 653, 654, 655, 656, 657, 658, 659, 660, 661, 662, 663, 664, 665, 666, 667, 668, 669, 670, 671, 672, 673, 674, 675, 676, 677, 678, 679, 680, 681, 682, 683, 684, 685, 686, 687, 688, 689, 690, 691, 692, 693, 694, 695, 696, 697, 698, 699, 700, 701, 702, 703, 704, 705, 706, 707, 708, 709, 710, 711, 712, 713, 714, 715, 716, 717, 718, 719, 720, 721, 722, 723, 724, 725, 726, 727, 728, 729, 730, 731, 732, 733, 734, 735, 736, 737, 738, 739, 740, 741, 742, 743, 744, 745, 746, 747, 748, 749, 750, 751, 752, 753, 754, 755, 756, 757, 758, 759, 760, 761, 762, 763, 764, 765, 766, 767, 768, 769, 770, 771, 772, 773, 774, 775, 776, 777, 778, 779, 780, 781, 782, 783, 784, 785, 786, 787, 788, 789, 790, 791, 792, 793, 794, 795, 796, 797, 798, 799, 800, 801, 802, 803, 804, 805, 806, 807, 808, 809, 810, 811, 812, 813, 814, 815, 816, 817, 818, 819, 820, 821, 822, 823, 824, 825, 826, 827, 828, 829, 830, 831, 832, 833, 834, 835, 836, 837, 838, 839, 840, 841, 842, 843, 844, 845, 846, 847, 848, 849, 850, 851, 852, 853, 854, 855, 856, 857, 858, 859, 860, 861, 862, 863, 864, 865, 866, 867, 868, 869, 870, 871, 872, 873, 874, 875, 876, 877, 878, 879, 880, 881, 882, 883, 884, 885, 886, 887, 888, 889, 890, 891, 892, 893, 894, 895, 896, 897, 898, 899, 900, 901, 902, 903, 904, 905, 906, 907, 908, 909, 910, 911, 912, 913, 914, 915, 916, 917, 918, 919, 920, 921, 922, 923, 924, 925, 926, 927, 928, 929, 930, 931, 932, 933, 934, 935, 936, 937, 938, 939, 940, 941, 942, 943, 944, 945, 946, 947, 948, 949, 950, 951, 952, 953, 954, 955, 956, 957, 958, 959, 960, 961, 962, 963, 964, 965, 966, 967, 968, 969, 970, 971, 972, 973, 974, 975, 976, 977, 978, 979, 980, 981, 982, 983, 984, 985, 986, 987, 988, 989, 990, 991, 992, 993, 994, 995, 996, 997, 998, 999, 1000.

15. März 1679

## Das dyadische Zahlensystem. Teil I

Nebenstehende Folge kann leicht fortgesetzt werden, wenn man, von rechts nach links gehend, unter die Eins der oberen Zahl jeweils 0 schreibt, bis bei der oberen auch 0 vorkommt, worunter man dann 1 schreibt; weiter braucht man nicht zugehen, da die übrigen Ziffern gleich bleiben wie bei der oberen Zahl.

So wird aus: 1010111 87

1011000 88

Das ist dasselbe, wie wenn man

sagte: 1011000 ist	64
$2^6 + * + 2^4 + 2^3 + * * *$	16
64 + 16 + 8	8
	88

# 1.1 Zahl zur Basis b

- b-näre Zahl, b-adische Zahl, polyadische Zahl, Stellenwertzahl, Positionssystem
- Folge von Ziffern  $X_i \in \{0, 1, \dots, b-1\}$
- Darstellung  $(X_n \dots X_1)_b$  besitzt den Wert  $\sum X_i * b^{i-1}$ 
  - wird oft als Gleichung geschrieben, obwohl es keine ist

$$(X_n \dots X_1)_b = \sum_{i=1}^n X_i * b^{i-1}$$

- Erweiterung für gebrochene Zahlen

$$(X_n \dots X_1, X_0 \dots X_{-k})_b = \sum_{i=-k}^n X_i * b^{i-1}$$

- b=2 : Dualzahl
- b=10: Dezimalzahl
- b=16 : Hexadezimalzahl
- b=8 : Oktalzahl
- Spezielle Formen
  - b ist negativ
  - b ist komplexe Zahl



# Dualzahl, Oktalzahl, Hexadezimalzahl

7

Dualzahl	Oktalzahl	Hexadezimalzahl	Wert, dargestellt als Dezimalzahl
0000	0	0	0
0001	1	1	1
0010	2	2	2
0011	3	3	3
0100	4	4	4
0101	5	5	5
0110	6	6	6
0111	7	7	7
1000	10	8	8
1001	11	9	9
1010	12	A	10
1011	13	B	11
1100	14	C	12
1101	15	D	13
1110	16	E	14
1111	17	F	15

# Notation

- Mit  $X[n] = X_n X_{n-1} \dots X_1$  ist zum einen die **Darstellung** einer Zahl als Bitfolge gemeint
- Mit  $X[n]$  ist zum anderen die **n-stellige Dualzahl** gemeint, die den folgenden[dualen/binären] Wert besitzt

$$\sum_{i=1}^n X_i * 2^{i-1}$$

- Mit kleinem Buchstaben **x** ist der **Wert** der darzustellenden Zahl gemeint.
- Wenn die Stellenzahl n vereinbart ist, dann wird verkürzt nur X anstelle von X[n] benutzt
- Zur Darstellung der Dualziffern wird hier oft
  - o, 1 benutzt anstelle von
  - 0, 1

$$X[n] = X_n X_{n-1} \dots X_1 = \sum_{i=1}^n X_i * 2^{i-1}$$

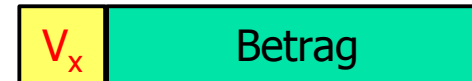


# Beispiel: Darstellung und Wert

X[4] Darstellung	X dualer oder binärer Wert (als Dezimalzahl dargestellt)	z.B. Wert x einer 2K-Zahl (als Dezimalzahl dargestellt)
0000	0	0
0001	1	1
0010	2	2
0011	3	3
0100	4	4
0101	5	5
0110	6	6
0111	7	7
1000	8	-8
1001	9	-7
1010	10	-6
1011	11	-5
1100	12	-4
1101	13	-3
1110	14	-2
1111	15	-1

# 1.2 Vorzeichen-Betrag-Darstellung

- Vorzeichenbit  $V_x$  und Betrag  $|x|$ 
  - $V_x=0$  für positive Werte,  $V_x=1$  für negative Werte
- Eine positive Dualzahl  $X[n]$  mit Vorzeichenbit  $V_x$  nennen wir Vorzeichenzahl  $(V_x, X[n])$
- Zwei Darstellung für die Null sind möglich
- Darstellung ist
  - gut für Multiplikation, Division
    - getrennte Behandlung von Vorzeichen und Betrag
  - schlecht für Addition und Subtraktion
    - Hardware zu kompliziert



Wertebereich

$$|x| \leq 2^n - 1$$

Abbildung

$$X[n] = |x| \quad \text{und} \quad V_x = \begin{cases} 0 & \text{für } x \geq 0 \\ 1 & \text{für } x \leq 0 \text{ bzw. } x < 0 \end{cases}$$

Rückabbildung

$$x = (1 - 2V_x) * X[n] = (\overline{V_x} - V_x) * X[n]$$

# Offset Binary, Exzess-Darstellung

- Durch Addition der Basis  $2^{n-1}$  wird der Zahlenbereich so nach oben verschoben, daß nur noch Werte = 0 entstehen.
- Definitionsbereich  
 $x = -2^{n-1} \dots (2^{n-1} - 1)$
- Abbildung  $X[n] = x + 2^{n-1}$
- Rückabbildung  $x = X[n] - 2^{n-1}$
- MSB gibt Vorzeichen an
  - = 0: negative Zahl
  - = 1: positive Zahl
- Anwendungen
  - Digital/Analog-Umsetzer
  - zur Darstellung des Exponenten in Gleitkommazahlen, als sogen. Charakteristik

Dezimal $x$	Offset Binary $X[4]$
+7	1111
+6	1110
+5	1101
+4	1100
+3	1011
+2	1010
+1	1001
0	1000
-1	0111
-2	0110
-3	0101
-4	0100
-5	0011
-6	0010
-7	0001
-8	0000



## 1.3 Zweikomplementzahl

---

12

- auch 2K-Zahl, signed number
- Standard für Festkommarithmetik
- positive Werte werden nicht verändert
  - $x=5 \rightarrow X[4]=5 = 0101$
- auf negative Werte wird  $2^n$  addiert
  - $x=-5 \rightarrow X[4]=-5 + 2^4 = 11 = 1011$
  - an dem MSB kann man erkennen, daß der dargestellte Wert negativ ist, dieses Bit wird auch Sign genannt.
- Rückabbildung
  - falls  $X_n=1$  dann muß  $2^n$  abgezogen werden
  - $X[4]=1011 \rightarrow x = 11 - 2^4 = -5 = 1011 - 10000 = 011 - 1000$



# Zweikomplementzahl

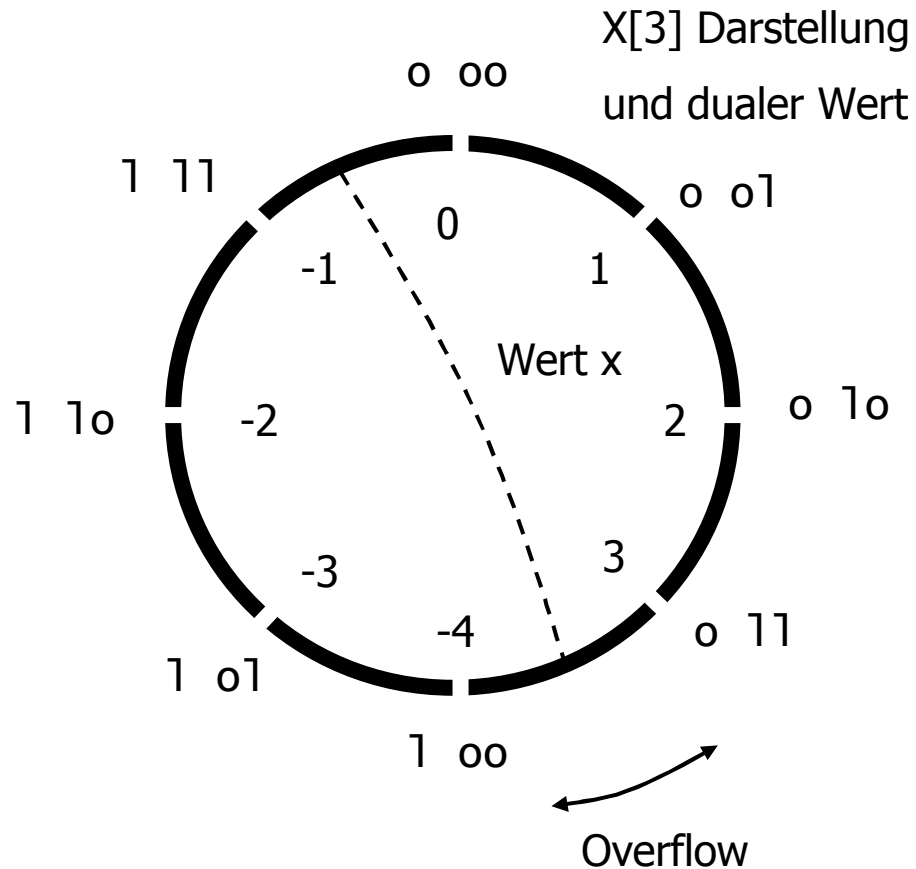
Abbildung V2K

$$V2K : \{-2^{n-1} : 2^{n-1} - 1\} \rightarrow \{0 : 2^n - 1\}$$
$$x \mapsto X[n].$$

$$X[n] = x \bmod 2^n = \begin{cases} x & \text{für } x \geq 0 \\ 2^n + x & \text{für } x < 0 \end{cases}$$

Rückabbildung V2K<sup>-1</sup>

$$x = X[n] - X_n * 2^n = \begin{cases} X[n] & \text{für } X_n = 0 \\ X[n] - 2^n & \text{für } X_n = 1 \end{cases}$$
$$x = X[n-1] - X_n * 2^{n-1}$$



- Definitionsbereich ist unsymmetrisch zur Null
  - $-2^{n-1}$  darstellbar aber nicht  $+2^{n-1}$
  - bei der Addition kann das Ergebnis größer als  $2^{n-1}-1$  oder kleiner als  $-2^{n-1}$  werden → Überlauf (Overflow)
  - auch Komplementbildung von  $2^{n-1}$  führt zu Overflow!



# Wert x der 2K-Zahl X[n]

15

X[4] Darstellung	X dualer Wert	Wert x der 2K-Zahl
0000	0	0
0001	1	1
0010	2	2
0011	3	3
0100	4	4
0101	5	5
0110	6	6
0111	7	7
1000	8	-8
1001	9	-7
1010	10	-6
1011	11	-5
1100	12	-4
1101	13	-3
1110	14	-2
1111	15	-1

## Arithmetische Erweiterung

- **Aufgabe:** Die Anzahl der Stellen einer 2K-Zahl soll von  $m$  auf  $n > m$  vergrößert werden,  $X[m] \rightarrow X[n]$
- **Lösung:** wenn das Sign  $X_m$ 
  - 0 ist, dann werden alle weiteren Stellen links davon auf 0 gesetzt
  - 1 ist, dann werden alle weiteren Stellen links davon auf 1 gesetzt
  - allgemein:  $X[n] = X_m \dots X_m X[m]$
- **Beispiel**
  - $X[4] = 0101 \rightarrow X[8] = 00000101$
  - $X[4] = 1011 \rightarrow X[8] = 11111011$
- 2K-Zahlen, die kleine Werte darstellen, beginnen mit 00...0 oder 11...1





# Arithmetischer Schift von 2K-Zahlen

17

- **Nach rechts**, Vorzeichenstelle nachziehen
  - Entspricht einer Division durch 2
  - Das Vorzeichen bleibt erhalten
  - Eine Stelle geht verloren (falls sie nicht als Nachkommastelle aufgehoben wird).

$$x=-7, X2K[4] = 1001$$

$$\text{Nach arith-R-Shift: } 1100, x = -4$$

$$x=7, X2K[4] = 0111$$

$$\text{Nach arith-R-Shift: } 0011, x=3$$

Beachte Rundung!

- **Nach links**
  - Entspricht Multiplikation mit 2
  - Das Ergebnis ist u.U. mit n Stellen nicht mehr darstellbar (Überlauf)

$$x=-2, X2K[4] = 1110$$

$$\text{Nach arith-L-Shift: } 1100, x = -4$$

$$x=3, X2K[4] = 0011$$

$$\text{Nach arith-R-Shift: } 0110, x=6$$

# Komplementbildung

- **Problem:** Anstelle von  $x$  soll  $-x$  dargestellt werden
  - $X[n] \rightarrow X'[n]$
- **Lösung:**  $X'[n] = 2^n - X[n]$
- $X'[n]$  nennt man das Komplement von  $X[n]$  gegen  $2^n$
- Wie kann es einfach berechnet werden?

- Es gilt

$$X[n] + \overline{X[n]} = \sum_{i=1}^n (X_i + \overline{X_i}) * 2^{i-1} = 2^n - 1 = 11\dots 1$$

$$X'[n] = 2^n - X[n]$$

- einsetzen in
- ergibt

$$X'[n] = \overline{X[n]} + 1$$

- Beispiel

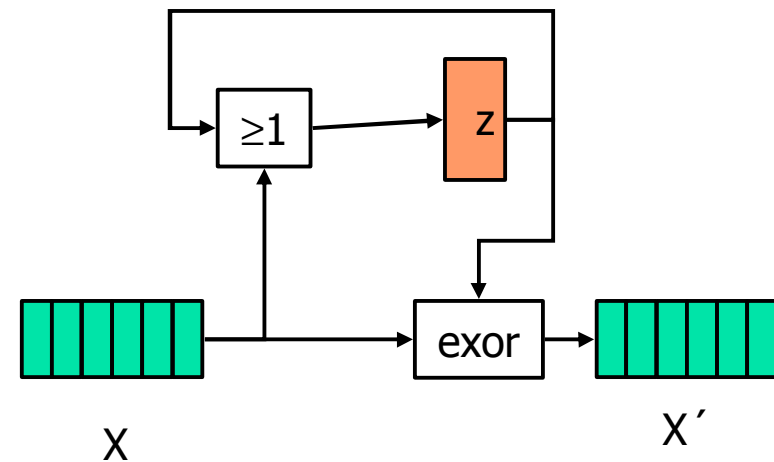
010100  $\rightarrow$  alle Bits negieren  $\rightarrow$  101011 + 1 = 101100

# Serielle Komplementbildung

- Betrachte alle Bits  $X_i$  von rechts nach links,  $i=1 \dots n$
- Solange das Bit  $X_i = 0$  wird es unverändert übernommen.
- Das erste Bit  $X_i = 1$  wird auch noch übernommen.
- Alle weiteren Bits werden negiert.
- Beispiel

■  $010100 \rightarrow 101100$

$\underbrace{\hspace{2em}} \quad \underbrace{\hspace{2em}}$



## Schaltwerk zur seriellen Komplementbildung

z ist am Anfang 0 und wird mit der ersten 1 auf 1 gesetzt, wodurch die folgenden Bits über exor negiert werden



## 1.4 Einskomplementzahl

---

- Auch 1K-Zahl, inzwischen kaum noch benutzt.
- Unterschied zur 2K-Zahl: Es wird das Komplement gegen  $2^n - 1$  gebildet.
- positive Werte werden nicht verändert
  - $x=5 \rightarrow X[4]=5 = 0101$
- auf negative Werte wird  $2^n - 1$  addiert
  - $x=-5 \rightarrow X[4]= -5 + 2^4 - 1 = 10 = 1010$
- Zahlenbereich ist symmetrisch zur Null.
- für die Null gibt es zwei Darstellungen
  - $00 \dots 0$
  - $11 \dots 1$  (falls nicht verboten)



# Einskomplementzahl

Abbildung

$$V1K : \{-2^{n-1} + 1 : 2^{n-1} - 1\} \rightarrow \{0 : 2^n - 2\}$$
$$x \mapsto X[n].$$

Abbildungsvorschrift

$$X[n] = x \bmod (2^n - 1) = \begin{cases} x & \text{für } x \geq 0 \\ 2^n - 1 + x & \text{für } x < 0 \end{cases}$$

Rückabbildungsvorschrift

$$x = X[n] - X_n * (2^n - 1) = \begin{cases} X[n] & \text{f. } X_n = 0 \\ X[n] - 2^n + 1 & \text{f. } X_n = 1 \end{cases}$$



# Komplementbildung

- Besonders einfach, alle Bits werden negiert:

$$X[n] + X'[n] = 2^n - 1 = 11 \dots 1 = X[n] + \overline{X[n]}$$
$$X'[n] = \overline{X[n]}.$$

- Warum wird die 1K-Darstellung nicht benutzt?
  - Bei der Addition muß im Falle eines Übertrages eine 1 addiert werden (Hardware- und Zeitaufwand)
  - Für die Addition von Dualzahlen und 2K-Zahlen kann derselbe Addierer benutzt werden.



# Vergleich

Wert $x$ dezimal	Vorzeichen zahl $V_x X[2]$	2K-Zahl $X[3]$	1K-Zahl $X[3]$
-4	-----	1 00	-----
-3	1 11	1 01	1 00
-2	1 10	1 10	1 01
-1	1 01	1 11	1 10
0	0 00 (1 00)	0 00	0 00 (1 11)
1	0 01	0 01	0 01
2	0 10	0 10	0 10
3	0 11	0 11	0 11

# Festkommazahlen in Programmiersprachen

Typ Delphi	C	JAVA	Bereich	Bits	
Byte	unsigned char		0..255	8	un- signed
Word	unsigned int		0..65535	16	
Longword	unsigned long int		0..4294967295	32	
Shortint	char	byte	-128..127	8	signed 2K
Smallint	int	short	-32768..32767	16	
Longint	long	int	-2147483648..2147483647	32	
Int64	long long int	long	$-2^{63}..2^{63}-1$	64	
Currency			-922337203685477.5808.. 922337203685477.5807	64	





## 1.5 Binärcodierte Dezimalzahl

---

25

- Eine Codierung von Dezimalzahlen im Rechner
- Gründe zur Verwendung
  - Keine Konvertierung für die externe Darstellung erforderlich
    - z.B. in Taschenrechnern
  - Vermeidung von Approximationsfehlern
    - $0,1 = 0,0\ 0011\ 0011\ 0011\ \dots$  (nicht endlich!)
    - In der kommerziellen Datenverarbeitung
- Jede BCD-Ziffer wird durch 4 Bits dual codiert.
- In ungepackter Form werden 8 Bits zur Codierung benutzt. Dabei werden die höchstwertigen Bits auf 0000 oder auf 0010 (ASCII) gesetzt, oder der EBCDIC wird verwendet.



# Beispiele

---

- Gepackte BCD-Zahlen
  - $59_{10} = 0101\ 1001_{\text{BCD}}$
  - $81_{10} = 1000\ 0001_{\text{BCD}}$
- Ungepackte BCD-Zahlen
  - $59_{10} = 0000101\ 00001001_{\text{BCD}}$
  - $81_{10} = 00001000\ 00000001_{\text{BCD}}$



# Negative BCD-Zahlen

---

27

- Möglichkeiten
  - Vorzeichen-Betrag-Darstellung
  - 9-Komplement
    - Komplement gegen  $10^m-1$
    - bei m Dezimalstellen
  - 10-Komplement
    - Komplement gegen  $10^m$

# 9-Komplement

- $m$  = Anzahl der Stellen
- $|x| < 499\dots 9$
- Vorzeichenstelle (4 Bits)
  - 0 bis 4: positive Zahlen
  - 5 bis 9: negative Zahlen
- Komplementbildung
  - $(X_{9K}, x) \rightarrow (X'_{9K}, -x)$
  - Für jede Ziffer wird Differenz gegen 9 gebildet
    - durch Addition von 6 = 0110 und anschließender Negation aller Bits

$$X_{9K} = \begin{cases} x & \text{für } x \geq 0 \\ (10^m - 1) + x & \text{für } x < 0 \end{cases}$$

$$X'_{9K} = (10^m - 1) - X_{9K} = 99\dots 9 - X_{9K}$$

$$\begin{aligned} X_{9K}(i) &= 15 - (6 + X_{9K}(i)) \\ &= 1111 - (0110 + X_{9K}(i)) = \overline{0110 + X_{9K}(i)} \end{aligned}$$

# 10-Komplement

- $m$  = Anzahl der Stellen
- $-500\dots0 \leq x \leq 499\dots9$
- Vorzeichenstelle (4 Bits)
  - 0 bis 4: positive Zahlen
  - 5 bis 9: negative Zahlen
- Komplementbildung
  - $(X1_{10K}, x) \rightarrow (X'_{10K}, -x)$
  - Komplement von  $X_{10K} = 500\dots0$  ist verboten (Überlauf)
  - 1.) Negiere alle Bits und addiere 1, merke Überträge zwischen den Tetraden
  - 2.) Addiere lolo (-6) auf alle Tetraden, die keinen Übertrag erzeugt haben

$$X_{10K} = \begin{cases} x & \text{für } x \geq 0 \\ 10^m + x & \text{für } x < 0 \end{cases}$$

$$\begin{aligned} X'_{10K} &= 10^m - X_{10K} = (10^m - 1) - X_{10K} + 1 \\ &= (99\dots9 - X_{10K}) + 1. \end{aligned}$$

$$\begin{array}{ll} X_{10K} = 509400 & x = -490600 \\ X'_{10K} = 490600 & x = +490600 \end{array}$$



# Beispiele 9er und 10er Komplement

30

Number	10's complement	9s' compleme
1849	8151	8150
2067	7933	7932
100	9900	9899
7	9993	9992
8151	1849	1848
0	10000 (= 0)	9999

- Beachte: 9er Komplement hat zwei Nullen
  - +0 und -0 = 9999

# BCD-Darstellung in der IntelArchitecture-32

31

- 80 Bits
- Vorzeichen S und 18 Dezimalstellen (Magnitude)
- Kompatibel zum COBOL-Standard





# 1.6 Gleitkommazahl, Motivation

32

- Festkommazahlen können nur einen relativ kleinen Zahlenbereich überstreichen
- n Bit Wort
  - positive Dualzahl:  $0 \dots 2^n - 1$
  - Darstellung größerer Zahlen durch Verschiebung des fiktiven Kommas nach rechts um m Stellen: Zahlenbereich wird mit  $2^m$  multipliziert
    - Abstand zwischen den Zahlen ist  $2^m$
  - Verschiebung des Kommas um m Stellen nach links: Zahlenbereich wird durch  $2^m$  dividiert
    - Abstand zwischen den Zahlen ist  $2^{-m}$
- Problem
  - Der Programmierer muß sich für die Anwendung den Zahlenbereich genau überlegen und das fiktive Komma programmtechnisch mitführen (die Festkommazahl wird auf  $2^m$  normiert.) Ein Überlauf wird durch geeignete Normierung vermieden
- Lösung
  - Codierung der Kommastellung und automatische Normierung
  - Zusätzliche Bits für die Kommastellung
  - Gleitendes Komma: floating point





# Eigenschaften der Gleitkommazahl

33

## ■ Vorteile

- größerer Wertebereich darstellbar
  - seltener: Überlauf (Overflow), Unterlauf (Underflow)
  - Programmierung vereinfacht (keine vorausgeplante Normierung)
- Relative Darstellungsgenauigkeit in etwa konstant

## ■ Allgemein

- Anzahl der GK-Zahlen ist endlich
- Es gibt keine
  - beliebig großen/kleine Zahlen
  - beliebig dicht benachbarte GKZ
  - Der Abstand zwischen benachbarten GKZ ist nicht konstant.

## ■ Nachteile

- Hardware-Implementierung deutlich komplexer als für Festkommazahlen
- Rundungsfehler entstehen, insbesondere bei Addition/Subtraktion von Zahlen deren Kommastellungen sich stark voneinander unterscheiden. Genauigkeit wird vorgetäuscht.
- Je nach der Auswertungsart von Formeln ergeben sich unterschiedlich große Fehler  
 $(a-b) + (c-d) \neq (a+c) - (b+d)$

# Normalisierung

- Wert  $x$
- Mantisse  $mx$
- Basis  $b$
- Exponent  $ex$
- Vorzeichen  $v_x = +1/-1$
- $k$  wird festgelegt (meist 0 oder 1)
- **Allgemein:** Mantisse wird normalisiert, um die Darstellung eindeutig zu machen und um immer mit der höchsten Genauigkeit zu rechnen,
  - $mx_{norm}$  hat nur  $r$  Stellen hinter dem Komma, die erste Stelle hinter dem Komma ist  $>0$
- **Speziell beim IEEE 754 - Standard** wurde  $b=2$  und  $k=1$  gewählt, d.h. die normalisierte Mantisse (mit  $\hat{\phantom{x}}$ ) liegt zwischen 1 und  $2-\varepsilon$

$$x = mx * b^{ex}$$

$$x = v_x * \underbrace{mx_{norm}}_{\hat{mx}_{norm}} * b^k * b^{ex}$$

$$\frac{1}{b} \leq mx_{norm} \leq 1 - \frac{1}{b^r} < 1$$

$$x = v_x * \hat{mx}_{norm} * b^{ex}$$

$$1 = \hat{mx}_{norm} = 1,11 \dots 1$$

**IEEE** is an abbreviation for Institute of Electrical and Electronic Engineers

# Beispiel Normalisierung bei $b=10$

- Darstellung von 300
  - $mx=3, ex=2$
  - $mx=30, ex=1$
  - $mx=300, ex=0$
- Jetzt mit  $k=1, r=1$  normieren
  - Normierte Mantisse
    - $mxnorm$  Minimal = 0,1
    - $mxnorm$  Maximal = 0,9
    - $mx^{norm}$  Minimal=1
    - $mx^{norm}$  Maximal=9
  - Eindeutige Darstellung von 300
    - $mxnorm=0,3 \quad ex=2$

$$x = v_x * \underbrace{mx_{norm}}_{\widehat{mxnorm}} * b^k * b^{ex}$$

$$\frac{1}{b} \leq mx_{norm} \leq 1 - \frac{1}{b^r} < 1$$

# Relativer Fehler

- Abstand zwischen zwei Werten (r=Stellen der Mantisse hinter dem Komma) ergibt sich aus dem LSB der Mantisse und dem Exponenten

$$\Delta x = \frac{1}{b^r} * b^{exp}$$

- Relativer Fehler (LSB unsicher)

- höchstens

- $b^{-r}/b^{k-1}$

- mindestens

- $b^{-r}/(b^k - b^{-r})$

**IEEE, k=1**

1,00 ... 0x  $\rightarrow 2^{-r}/1$

1,11 ... 1x  $\rightarrow 2^{-r}/(2-2^{-r}) \approx 2^{-r}/2$

Wie ungenau ist die Darstellung von 8 Millionen € wenn r=23 Nachkommastellen benutzt werden?

Maximaler Relativer Fehler:  $2^{-23} \rightarrow$  etwa 1 € ungenau



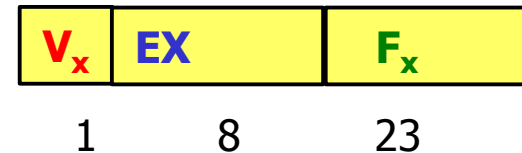
# Darstellung des Exponenten

37

- Charakteristik (biased exponent)
- Auf den Exponenten (true exponent) wird eine Basis (bias) addiert, so daß nur noch positive Werte entstehen.
- 1. Möglichkeit:  
Addition von  $2^{n-1}$ 
  - **$EX[n] = ex + 2^{n-1}$**  ,
  - verwendet bei DEC-Rechnern, IBM/Siemens-Großrechner mit  $b=16$
- 2. Möglichkeit:  
Addition von  $2^{n-1} - 1$ 
  - **$EX[n] = ex + 2^{n-1} - 1$**
  - IEEE-Standard
- Durch Verwendung der Charakteristik-Darstellung und der Anordnung (Exponent steht vor der Mantisse) kann der Größenvergleich zweier GK-Zahlen relativ leicht durchgeführt werden:
  - 1. Vorzeichen prüfen
  - 2. Dualzahlen-Vergleich der Charakteristiken
  - 3. Dualzahlen-Vergleich der Mantissen

# IEEE-32-Bit-Gleitkommataformat

- **Vorzeichen  $V_x$**  (Sign)
- **Charakteristik  $EX = ex+127$** 
  - für  $ex = -126 \dots +127$
- Die **Mantisse**  $MX = 1,FX$  wird auf Werte zwischen 1 und 2 normalisiert
  - $1 = MX = 2 - 2^{-23}$
  - Dargestellt werden nur die 23 Binärstellen hinter dem Komma (fraction FX), da die Stelle vor dem Komma immer gleich 1 ist.
- **Wert  $x$**
- Betragsmäßig größter und kleinster Wert (normalisiert):



$$V_x = \begin{cases} 1 & \text{für negative Zahlen} \\ 0 & \text{für positive Zahlen} \end{cases}$$

fmann, TUD

$$x = (-1)^{V_x} * 2^{EX-127} * (1, FX)$$

$$= \begin{cases} +2^{EX-127} * (1, FX) & \text{für } V_x = 0 \\ -2^{EX-127} * (1, FX) & \text{für } V_x = 1 \end{cases}$$

$$x_{\max} = 2^{127} * (2 - 2^{-23}) \approx 1,7 * 10^{38}$$

$$x_{\min} = 2^{-126} * (1,0) \approx 1,17 * 10^{-38}$$

# Gültige und ungültige Exponenten

$ex$	EX	Interpretation
128	255 = 1111 1111	kein gültiger Exponent; zur Darstellung von $\infty$ und Kontrollcodes
127	254 = 1111 1110	gültige Exponenten für normalisierte Gleitkommazahlen
126	253 = 1111 1101	
⋮		
0	127 = 0111 1111	
-1	126 = 0111 1110	
⋮		
-126	1 = 0000 0001	
-127	0 = 0000 0000	kein gültiger Exponent; zur Darstellung von Null bzw. kleiner Zahl



# Beispiele

---

$$\begin{aligned}x &= 3,5 = 1,75 * 2^1 = 1,75 * 2^{128-127} \\ &= 1,110\ 0000\ 0000\ 0000\ 0000\ 0000 * 2^{128-127} \\ &= (0, 128, 0,75) \\ &= (0, 1000\ 0000, 110\ 0000\ 0000\ 0000\ 0000\ 0000)\end{aligned}$$

$$\begin{aligned}x &= -11,375 = -1011,011 * 2^0 = -1,011011 * 2^3 \\ &= (1, 130, 0,421875) \\ &= (1, 1000\ 0010, 011\ 0110\ 0000\ 0000\ 0000\ 0000)\end{aligned}$$





# Sonderfälle (1)

## a) Die Null

Darstellung der positiven Null:  $(0, 0, 0)$

Darstellung der negativen Null:  $(1, 0, 0)$

## b) Kleine Zahl (denormalized number)

Darstellung:  $(V_{\infty}, 0, \neq 0)$  Betrag:  $0 < |x| < 2^{-126}$

Kleine Zahlen repräsentieren sehr kleine Werte, die kleiner als  $x_{\min}$  und größer als Null sind. Der Wert ergibt sich zu

$$x = (-1)^{V_{\infty}} * 2^{EX-126} * (0, FX) .$$

## c) Große Zahl, $\infty$

Darstellung:  $(V_{\infty}, 255, 0)$  Betrag:  $|x| > x_{\max}$

$\infty$  kann sowohl ein positives als auch ein negatives Vorzeichen besitzen.

## d) Ungültige Zahl (Not a Number, *NAN*)

Darstellung:  $(V_{\infty}, 255, \neq 0)$

Ungültige Zahlen werden als Kontrollcodes interpretiert. Entweder entstehen bestimmte Kontrollcodes bei der Ausführung unzulässiger Operationen (z. B.  $0 * \infty$ ) oder sie werden dazu benutzt, Statuswerte durch eine Serie von Operationen zu schleusen.

# Sonderfälle (2)

## e) Überlauf (Overflow)

Überlauf kann bei der Addition, Multiplikation und Division auftreten, wenn nach dem Runden ein Ergebnis  $|x| > x_{\max}$  entsteht.

## f) Unterlauf (Underflow)

Unterlauf kann bei der Addition, Multiplikation, Division und Reziprokwertbildung ( $1/x$ ) auftreten, wenn nach dem Runden ein Ergebnis  $0 < |x| < x_{\min}$  entsteht. Als Ergebnis wird entweder eine „Kleine Zahl“ erzeugt oder  $\pm 0$ .

## g) Besondere Operationen

Besondere Operationen können für die Verarbeitung von  $\infty$ , kleinen Zahlen und durch die Unterscheidung von positiver und negativer Null definiert werden.

## h) Rundung bei einem nicht darstellbaren Rest $R$

Vier Rundungsarten stehen zur Auswahl:

- Rundung zur nächsten darstellbaren Zahl:

$$\begin{aligned} |x|, R &\rightarrow |x| + 1 && \text{für } R \geq 0,5 \quad (\text{MSB von } R \text{ ist } 1) \\ &\rightarrow |x| && \text{für } R < 0,5 \quad (\text{MSB von } R \text{ ist } 0) \end{aligned}$$

- Rundung in Richtung  $+\infty$ :

$$\begin{aligned} +|x|, R &\rightarrow +|x| + 1 && \text{für } R \neq 0 \\ -|x|, R &\rightarrow -|x| \end{aligned}$$

zur nächsten

Nach oben

- Rundung in Richtung  $-\infty$ :

$$\begin{aligned} +|x|, R &\rightarrow +|x| \\ -|x|, R &\rightarrow -|x| - 1 && \text{für } R \neq 0 \end{aligned}$$

- Rundung in Richtung 0

$$|x|, R \rightarrow |x|$$

Nach unten

Nach Null



# IEEE-64-Bit-Gleitkommaformat

---

43

- Das IEEE-64-Bit-Gleitkommaformat (long real) entspricht dem IEEE-32-Bit-Gleitkommaformat (short real) in seiner Aufbaustruktur.
- Für das Vorzeichenbit wird 1 Bit, für die Charakteristik werden 11 Bits (bias = 1023) und für Fraction werden 52 Bits benutzt.
- Für noch höhere Ansprüche an die Rechengenauigkeit kann das IEEE-Gleitkommaformat auf (1, = 15, = 63) Bits erweitert werden.

# Gleitkommazahlen in Anwendungen

	Bereich +/-	Dezimalstellen-Genauigkeit	Intel-Architecture 32	C	JAVA	Delphi
32-Bit IEEE Short Real	$1.2 \times 10^{-38} .. 1.7 \times 10^{38}$	7-8	Single Real	float	float	Single
64-Bit IEEE Long Real	$2.2 \times 10^{-324} .. 2.0 \times 10^{308}$	15-16	Double Real	double	double	Double
80-Bit IEEE			Extended Real			
96-Bit IEEE				long double		



## 3.4 Konvertierung zwischen Zahlensystemen

45

- **Gegeben:** m-stellige ganze Zahl A zur Basis a
- **Gesucht:** n-stellige ganze Zahl zur Basis b
- **Lösungsansatz:** stelle A oder B als Potenzreihe oder nach dem Hornerschema dar, setze A=B.

$$A = (A_m A_{m-1} \dots A_1)_a \rightarrow (B_n B_{n-1} \dots B_1)_b = B$$

$$A = A_m * a^{m-1} + \dots + A_2 * a + A_1$$

$$B = B_n * b^{n-1} + \dots + B_2 * b + B_1$$

$$A = (\dots (A_m * a + A_{m-1}) * a + \dots + A_2) * a + A_1$$

$$B = (\dots (B_n * b + B_{n-1}) * b + \dots + B_2) * b + B_1$$

# Summationsmethode

- $A_a \rightarrow B_b$
- Erstelle eine Tabelle, in der für jede Ziffer  $A_i$  und die Stellenwertigkeit  $a^{i-1}$  der äquivalente Summand  $A_i * a^{i-1}$  (codiert im Zielsystem  $b$ ) enthalten ist.
- Lies für Ziffer  $A_i$  (Zeile der Tabelle) mit der Stellenwertigkeit  $a^{i-1}$  (Spalte) den äquivalenten Summanden 😊.
- Addiere alle Summanden im Zielsystem.

$$\sum_{i=1}^m A_i * a^{i-1} = B$$

	Äquivalenter Summand, für Stellenwertigkeit $a^{i-1}$ , Basis $b$				
Ziffer $A_i$ Basis $a$	$i=0$	1	2	...	$m-1$
0					
1			😊		
...					
$a-1$					

# Hexadezimalzahl → Dezimalzahl

Hexadezimal- ziffer	Stellenwertigkeit					
	*16 <sup>0</sup>	*16 <sup>1</sup>	*16 <sup>2</sup>	*16 <sup>3</sup>	*16 <sup>4</sup>	*16 <sup>5</sup>
0	00	000	0000	00000	000000	00000000
1	01	016	0256	04096	065536	01048576
2	02	032	0512	08192	131072	02097152
3	03	048	0768	12288	196608	03145728
4	04	064	1024	16384	262144	04194304
5	05	080	1280	20480	327680	05242880
6	06	096	1536	24576	393216	06291456
7	07	112	1792	28672	458752	07340032
8	08	128	2048	32768	524288	08388608
9	09	144	2304	36864	589824	09437184
A	10	160	2560	40960	655360	10485760
B	11	176	2816	45056	720896	11534336
C	12	192	3072	49152	786432	12582912
D	13	208	3328	53248	851968	13631488
E	14	224	3584	57344	917504	14680064
F	15	240	3840	61440	983040	15728640
Äquivalenter dezimaler Summand						

Beispiel

$$\begin{array}{r}
 = \quad \text{F18C}_{16} \\
 + \quad \text{F000} \quad 61440 \\
 + \quad 100 \quad 256 \\
 + \quad 80 \quad 128 \\
 + \quad \text{C} \quad 12 \\
 \hline
 61836_{10}
 \end{array}$$

# Dezimalzahl → Hexadezimalzahl

Dezimal- ziffer	Stellenwertigkeit						
	*10 <sup>0</sup>	*10 <sup>1</sup>	*10 <sup>2</sup>	*10 <sup>3</sup>	*10 <sup>4</sup>	*10 <sup>5</sup>	*10 <sup>6</sup>
0	0	00	000	0000	00000	000000	0000000
1	1	0A	064	03E8	02710	186A0	0F4240
2	2	14	0C8	07D0	04E20	30D40	1E8480
3	3	1E	12C	0BB8	07530	493E0	2DC6C0
4	4	28	190	0FA0	09C40	61A80	3D0900
5	5	32	1F4	1388	0C350	7A120	4C4B40
6	6	3C	258	1770	0EA60	927C0	5B8D80
7	7	46	2BC	1B58	11170	AAE60	6ACFC0
8	8	5D	320	1F40	13880	C3500	7A1200
9	9	5A	384	2328	15F90	DBBA0	895440
Äquivalenter hexadazimaler Summand							

Beispiel

$$\begin{array}{r}
 61836_{10} \\
 = 60000 \quad EA60 \\
 + 1000 \quad 3E8 \\
 + 800 \quad 320 \\
 + 30 \quad 1E \\
 + 6 \quad 6 \\
 \hline
 F18C_{16} = 1111 \ 0001 \ 1000 \ 1100
 \end{array}$$



# Divisionmethode

- $A_a \rightarrow B_b$
- Bei der sukzessiven ganzzahligen Division durch  $b$  stellen die Reste die gesuchten Ziffern  $B_1, B_2, \dots, B_n$  dar.
- Division und Restbildung erfolgt im **Quellsystem**
- Ziffernumcodierung ist erforderlich wenn  $b > a$  ist
- Für den Menschen geeignet; im Rechner sollte die Division vermieden werden.

$$A = (\dots (B_n b + B_{n-1})b + \dots + B_2)b + B_1$$

Beispiel: Dezimalzahl  $\rightarrow$  Dualzahl

196/2 = 98	R	0	} 1100 0100
98/2 = 49	R	0	
49/2 = 24	R	1	
24/2 = 12	R	0	
12/2 = 6	R	0	
6/2 = 3	R	0	
3/2 = 1	R	1	
1/2 = 0	R	1	

LSB

# Divisionsmethode, Dualzahl → Dezimalzahl

- Dieses Beispiel demonstriert
  - den Aufwand im Rechner wegen der Division von Dualzahlen
  - die Ziffernumcodierung, weil  $b > a$
  - → Wahl einer Methode mit weniger Aufwand

196		10		19		R		6		
/		/		/		/		/		
1100	0100	/	1010	=	10	011	R	0110		}
1	0011	/	1010	=		1	R	1001		
	1	/	1010	=		0	R	0001		
			⏟							
			10							



# Divisionsmethode für gebrochene Zahlen

51

- **Gegeben:** Gebrochene Dezimalzahl  $0,A_{-1} A_{-2} \dots A_{-m}$
- **Gesucht:** Gebrochene Dualzahl  $B_b = 0,B_{-1} B_{-2} \dots B_{-n}$
- **Ansatz:**  $0,B_{-1} B_{-2} \dots B_{-n} = A$
  
- for  $i:=1$  to  $n$  do (im Quellsystem)
  - $A:=A * b$
  - $B_{-i}:=$  Ganzzahliger Teil von  $A$
  - $A:=A - B_{-i}$
  
- Die Methode kann zu den Divisionsmethoden gerechnet werden, weil eine Division durch  $1/b$  im Quellsystem durchgeführt wird, die einer Multiplikation mit  $b$  entspricht.
- Für gebrochene Zahlen ist die Divisionsmethode im Rechner geeignet.

# Beispiel: Dezimal $\rightarrow$ Dual

52

- **Beispiel**

- $0,421875 * 2 = 0,84375$

- $0,84375 * 2 = 1,6875$

- $0,6875 * 2 = 1,3750$

- $0,3750 * 2 = 0,75$

- $0,75 * 2 = 1,5$

- $0,5 * 2 = 1,0$

- $0 * 2 = 0,0$

- $B = 0,0110110$

- Zum Nachrechnen

- $0,1 = 0,001100110011 \dots$  nicht endlich!

- Multiplikation mit 2 kann durch BCD-Addition ( $x+x=2x$ ) erfolgen.

# Beispiel: Dual $\rightarrow$ Dezimal

## Beispiel

$$0,0110110 * 1010 =$$

$$00,110110 \quad \text{um 1 Stelle verschoben}$$

$$0011,0110 \quad \text{um 3 Stellen verschoben}$$

$$0100,001110 * 1010$$

$$00,01110$$

$$\underline{0001,110}$$

$$0010,00110 * 1010$$

$$00,0110$$

$$\underline{0001,10}$$

$$0001,1110$$

$$0001,110$$

$$\underline{0111,00}$$

$$1000,110 * 1010$$

$$01,10$$

$$\underline{0110,}$$

$$0111,10 * 1010$$

$$01,0$$

$$\underline{0100,}$$

$$0101,0$$

$$B = 0,421875$$

Die Multiplikation mit  $10=1010$  wird auf die stellenversetzte Addition zurückgeführt

Es ergeben sich die Dezimalziffern im BCD-Code

# Multiplikationsmethode

- $A_a \rightarrow B_b$
- 1. Die Ziffern  $A_i$  werden umcodiert in die Basis  $b$  (Rechnung im Zielsystem)
- 2. for  $i:=m$  downto 1 do
  - $B:=B * a + A_i$

geeignet zur Konvertierung ganzer Zahlen  
im Rechner

$$(\dots(A_m * a + A_{m-1}) * a + \dots + A_2) * a + A_1 = B$$

Beispiel: Dezimalzahl 196 → Dualzahl

$$\overbrace{(0001 * 1010 + 1001)}^1 * 1010 + \overbrace{0110}^6 = 1100 \ 0100$$

Beispiel: Dualzahl 1100 0100 → Dezimalzahl

$$\begin{matrix} & 1 & & 1 & & 0 & & 0 & & 0 & & 1 & & 0 & & 0 \\ ((((((1 * 2 & + & 1) * 2 & + & 0) * 2 & + & 0) * 2 & + & 0) * 2 & + & 1) * 2 & + & 0) * 2 & + & 0 = & 196 \end{matrix}$$