



# Technische Grundlagen der Informatik 2

## Rechnerorganisation

### Kapitel 3: Rechnerarithmetik

Prof. Dr. Ben Juurlink

Fachgebiet: Architektur eingebetteter System  
Institut für Technische Informatik und Mikroelektronik  
Fak. IV – Elektrotechnik und Informatik

SS 2013

Nach dieser Vorlesung sollten Sie in der Lage sein:

- Binär-/Octal-/Hexadezimalzahlen zu Dezimalzahlen zu konvertieren und umgekehrt.
- 2-Komplement-Zahlen zu negieren.
- m-Bit 2-Komplement-Zahlen zu n-Bit zu konvertieren.
- die Unterschiede zwischen `addi` und `addiu`, `slt` und `sltu`, `lb` und `lbu`, etc. zu benennen.
- Überlauf (overflow) zu erklären. Erläutern wann er auftritt und wie er behandelt wird.
- Einen einfachen Addierer zu designen und zu erweitern.
- Einen „Carry-Look-Ahead“-Addierer zu beschreiben.
- Eine Fraktalzahl in „IEEE 754“ FP-Standard zu konvertieren und umgekehrt.
- MIPS-FP (Floating Point)-Befehlssatz zu nutzen.

- Vorzeichenbehaftete (signed) und vorzeichenlose (unsigned) Zahlen
- Addition und Subtraktion
- Multiplikation
- Division
- Gleitkommazahlen (Floating Point)
- Gleitkommazahlen in MIPS
- „Pentium bug“



Natürliche Zahlen können in jeder Basis repräsentiert werden:

$$a_{n-1} a_{n-2} \dots a_1 a_0_{(Basis B)} = a_{n-1} B^{n-1} + \dots + a_1 B^1 + a_0 B^0 = \sum_{i=0}^{n-1} a_i B^i$$

- B: **Basis (Radix)**, z.B. 10 (dezimal), 2 (binär), 8 (octal), 16 (hexadezimal)
- $B^i$ : Gewicht (**weight**) der i-ten Ziffer (**digit**)
- $a_i$ : i-te Ziffer aus der Menge  $\{0, 1, \dots, B-1\}$

Beispiel (binär):  $1011_B = 1 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 1 \times 2^0 = (8 + 0 + 2 + 1)_D = 11_D$

Diagram illustrating the binary expansion of 1011<sub>B</sub> to decimal 11<sub>D</sub>:

$a_3$	$a_1$	$B^3=8$	$B^2=4$	$B^1=2$	$B^0=1$
↓	↓	↓	↓	↓	↓
1	0	1	0	1	1
↑	↑	↑	↑	↑	↑
$a_2$	$a_0$	$a_3$	$a_2$	$a_1$	$a_0$

- Zerlegung nach **Horner-Schema**:

$$\sum_{i=0}^{n-1} a_i B^i = (((... (a_{n-1} B + a_{n-2}) B + ... + a_2) B + a_1) B + a_0$$

- Dezimal nach dual / binär:

167 <sub>D</sub> ->	167 / 2 = 83	Rest 1
	83 / 2 = 41	Rest 1
	41 / 2 = 20	Rest 1
	20 / 2 = 10	Rest 0
	10 / 2 = 5	Rest 0
	5 / 2 = 2	Rest 1
	2 / 2 = 1	Rest 0
	1 / 2 = 0	Rest 1

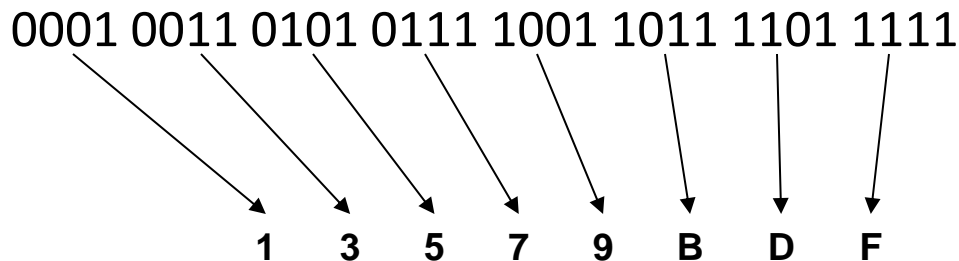
**Niederwertigstes Bit**  
(least significant bit (LSB))

**Höchstwertigstes Bit**  
(most significant bit (MSB))

- 167<sub>D</sub> = 10100111<sub>B</sub>



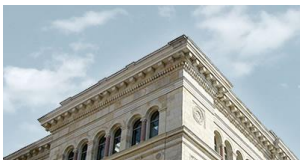
- Um lange Folgen mit Binärzahlen zu vermeiden, werden oft **Hexadezimalzahlen** (Basis 16) verwendet.
- Ziffernmenge: { 0, 1, 2, ... , 8, 9, A, B, ... ,F }
- Dezimal nach hexadezimal:
  - $167_D \rightarrow 167 / 16 = 10 \quad \text{Rest } 7$
  - $10 / 16 = 0 \quad \text{Rest } A$
  - $167_D = A7_H$
  - In C/Java: 0xa7
- Binär nach hexadezimal:





Einige Möglichkeiten negative Zahlen zu repräsentieren:

- *Vorzeichen-/Betrags-Zahlen* (Sign-magnitude numbers)
  - MSB zeigt Vorzeichen (*sign*) an (0: positiv, 1: negativ).
  - Die übrigbleibenden (n-1) Bits bilden den Betrag (*magnitude*).
- *1-Komplement-Zahlen* (One's complement numbers)
  - Zahl wird durch die Invertierung aller Bits negiert.
  - MSB *impliziert* das Vorzeichen.
- *2-Komplement-Zahlen* (Two's complement numbers)
  - MBS hat ein *negatives* Gewicht ( $-2^{n-1}$ ).
  - $b_{n-1}b_{n-2}...b_1b_0$  (binär) =  $-b_{n-1}2^{n-1}+b_{n-2}2^{n-2}+...+b_12^1+b_02^0$  (dezimal)
  - MSB *impliziert* das Vorzeichen.



Dezimal	Vorzeichen-/Betrags-Zahlen (Sign-Magnitude)	1-Komplement-Zahlen (One's complement)	2-Komplement-Zahlen (Two's complement)
-0	1000	1111	
-1	1001	1110	1111
-2	1010	1101	1110
-3	1011	1100	1101
-4	1100	1011	1100
-5	1101	1010	1011
-6	1110	1001	1010
-7	1111	1000	1001
-8			1000



- 2-Komplement-Zahlen
  - Arithmetik ist einfacher (identisch zu vorzeichenlos).
  - Es gibt nur eine Möglichkeit die 0 zu repräsentieren.
- MIPS benutzt 32-Bit 2-Komplement-Zahlen
- $b_{31}b_{30}...b_1b_0$  (binär) =  $-b_{31}2^{31}+b_{30}2^{30}+...+b_12^1+b_02^0$  (dezimal)
  - Was ist die kleinste von MIPS repräsentierbare Integerzahl?
  - Was ist die größte von MIPS repräsentierbare Integerzahl?
  - Was ist die binäre Repräsentation der 42? Hexadezimal?



- Negation von 2-Komplement-Zahlen:
  - Invertiere alle Bits und addiere 1
- 8-Bit-Beispiel:  $0110\ 1001_B = +105_D$

invertieren: 1001 0110

1 addieren: 1001 0111 =  $-128 + 16 + 7 = -105$

**Rückwärts:**

invertieren: 0110 1000

1 addieren: 0110 1001

- Konvertierung einer n-Bit 2-Komplement-Zahl in eine m-Bit 2-Komplement-Zahl ( $m > n$ ):
  - MIPS' 16-Bit Immediats werden für die Arithmetik auf 32 Bit erweitert.
  - Kopiere das MBS (sign bit) in die anderen Bitpositionen.
  - Beispiel (4-Bit  $\rightarrow$  8-Bit):

0010  $\rightarrow$  0000 0010

1010  $\rightarrow$  1111 1010

$$=-8+2=-6_D$$

$$=-128+64+32+16+8+2=-6_D$$

- Dies wird Vorzeichenerweiterung (sign extension) genannt.



- Einige Sprachen (z.B. C) können mit vorzeichenbehafteten (signed) und vorzeichenlosen (unsigned) Zahlen arbeiten.

C Datentyp	MIPS Datentyp	MIPS Ladebefehl
<code>int</code>	32-Bit Wort	<code>lw</code>
<code>unsigned int</code>	32-Bit wort	<code>lw</code>
<code>short</code>	16-Bit Halbwort	<code>lh</code>
<code>unsigned short</code>	16-Bit Halbwort (unsigned)	<code>lhu</code>
<code>char</code>	Byte	<code>lb</code>
<code>unsigned char</code>	Byte (unsigned)	<code>lbu</code>

- Javas* primitive Datentypen (`byte`, `short`, `int`, `long`) sind signed.
- `lh` und `lb` machen für das zu ladene Byte/Halbwort eine Vorzeichenerweiterung (`lhu` und `lbu` nicht).

- MIPS bietet 2 Versionen für einen *set-on-less-than*-Befehl an:
  - `slt` und `slti` arbeiten mit vorzeichenbehafteten Integers.
  - `sltu` und `sltiu` arbeiten mit vorzeichenlosen Integers.
- Beispiel:

```
$s0 = 1111 1111 1111 1111 1111 1111 1111 1111  
     = -1 (signed),  $2^{32}-1$  (unsigned)
```

```
$s1 = 0000 0000 0000 0000 0000 0000 0000 0001  
     = 1 (signed und unsigned)
```

```
slt    $t0,$s0,$s1      # $t0 = 1 (wahr[true])  
sltu   $t1,$s0,$s1      # $t1 = 0 (falsch[false])
```



- Addition von rechts nach links mit Übertrag (carry) wie in der Grundschule
- Beispiele (4-Bit):

$$\begin{array}{r}
 \text{Carry} \rightarrow 1 \\
 0010 \\
 + 0011 \\
 \hline
 0101 \\
 \\
 \text{Carry} \rightarrow 1 \\
 0101 \\
 + 0110 \\
 \hline
 1011
 \end{array}$$

$$\begin{array}{r}
 2 \\
 + 3 \\
 \hline
 5
 \end{array}$$

$$\begin{array}{r}
 5 \\
 + 6 \\
 \hline
 -5
 \end{array}$$

## Rechenregeln:

- $0 + 0 = 0$
- $0 + 1 = 1$
- $1 + 0 = 1$
- $1 + 1 = 10$

Carry

??? Überlauf (overflow)

- Beispiele (4-Bit):

	0010	2
	- 0011	- 3
	<u>111</u>	
Borrows	1111	-1
	1101	-3
	- 0110	- 6
	<u>11</u>	
	0111	7

## Rechenregeln:

- $0 - 0 = 0$
- $0 - 1 = 1$  (Borrow)
- $1 - 0 = 1$
- $1 - 1 = 0$

**wieder Überlauf**

- Alternative:

Den 2. Operanden negieren und dann beide addieren. (Wir werden diesen Trick beim Erstellen unserer **Arithmetic Logic Unit (ALU)** nutzen.)

- **Überlauf(overflow):**
  - Das Ergebnis ist zu groß für ein endliches Computer-Wort.
  - z.B. Addition von zwei n-Bit-Zahlen muss keine n-Bit-Zahl ergeben.
- Kein Überlauf, wenn ...
  - Addition von 2 Zahlen mit entgegengesetztem Vorzeichen
    - Betrag der Summe ist immer  $\leq$  den Beträgen der beiden Operanden
    - Beispiel:  $-10 + 6 = -4$
  - Subtraktion von 2 Zahlen mit dem selben Vorzeichen
- Overflow tritt auf, wenn...

Operation	A	B	Ergebnis
A+B	$\geq 0$	$\geq 0$	$< 0$
A+B	$< 0$	$< 0$	$\geq 0$
A-B	$\geq 0$	$< 0$	$< 0$
A-B	$< 0$	$\geq 0$	$\geq 0$



- Vergleich der Operationen  $A + B$  und  $A - B$ 
  - Kann ein Überlauf auftreten, wenn  $B = 0$  ist?
  - Kann ein Überlauf auftreten, wenn  $A = 0$  ist?
- Überlauf tritt nur auf, wenn das Übertragsbit (carry in) für das MSB  $\neq$  dem entstehenden Übertragsbit (carry out) aus der Operation der beiden MSBs ist.
- Beispiele (4-Bit):

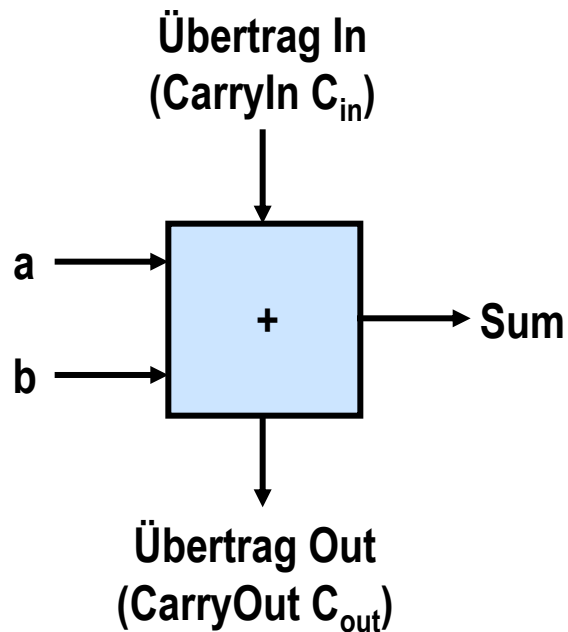
$$\begin{array}{r}
 \text{carry out} \longrightarrow 0 \quad 1 \longleftarrow \text{carry in} \\
 0101 \\
 + 0110 \\
 \hline
 1011
 \end{array}$$

- Wie werden Überläufe behandelt?
  - Eine **Exception (Interrupt)** wird ausgelöst.
  - Sprung zu einer vordefinierten Adresse, um die Ausnahme (Exception) zu behandeln.
    - **Interrupt Handler** (Teil des Betriebssystems (Operating System [OS]))
    - Register **\$k0** (26) und **\$k1** (27) sind fürs OS reserviert.
  - Unterbrochene Adresse wurde für einen möglichen Rücksprung im **exception program counter** (\$epc) gespeichert.
    - Neuer MIPS-Befehl: *move from coprocessor 0* (`mfc0 $k0, $epc`)
- Details basieren auf dem Software-System / der Software-Sprache
  - C/Java ignorieren Überlauf, Fortran nicht
- Nicht immer soll ein Überlauf angezeigt werden.
  - neue MIPS-Befehle: `addu`, `addiu`, `sltu`, `sltiu`, `subu`
  - **`addiu` und `sltiu` “sign extend” ihre 16-Bit Immediats bereits!**



- Wir werden jetzt ein **Addierer** basteln.
- Dazu brauchen wir einen 1-Bit **Volladdierer** (*1-bit Full Adder*).

Blockbild:



Wahrheitstabelle (Truth table):

a	b	$C_{in}$	$C_{out}$	Sum
0	0	0	0	0
0	0	1	0	1
0	1	0	0	1
0	1	1	1	0
1	0	0	0	1
1	0	1	1	0
1	1	0	1	0
1	1	1	1	1



## Boolsche Gleichungen (Summe von Produkten):

$$\text{Sum} = \bar{a}\bar{b}c_{\text{in}} + \bar{a}b\bar{c}_{\text{in}} + a\bar{b}\bar{c}_{\text{in}} + abc_{\text{in}}$$

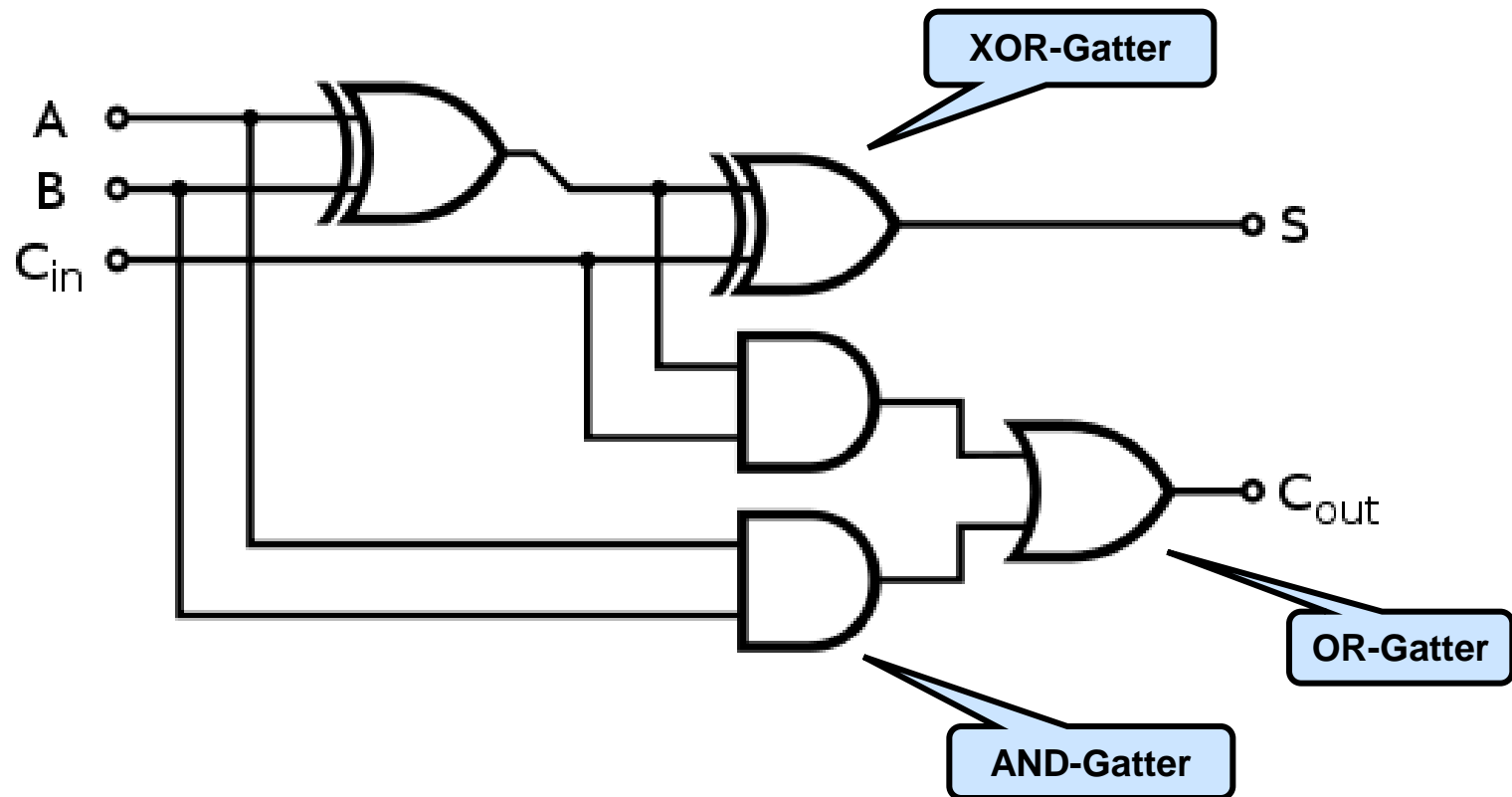
$$c_{\text{out}} = \bar{a}bc_{\text{in}} + a\bar{b}c_{\text{in}} + ab\bar{c}_{\text{in}} + abc_{\text{in}}$$

## Vereinfachung:

$$\text{Sum} = a \text{ XOR } b \text{ XOR } c_{\text{in}}$$

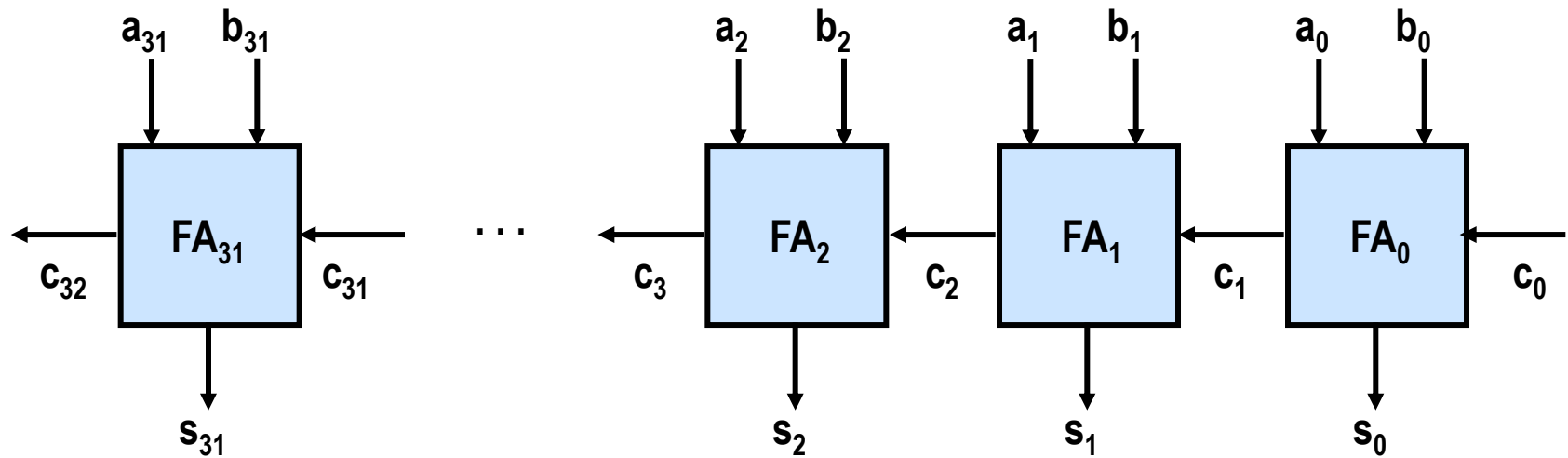
$$c_{\text{out}} = ab + ac_{\text{in}} + bc_{\text{in}}$$

a	b	c <sub>in</sub>	c <sub>out</sub>	Sum
0	0	0	0	0
0	0	1	0	1
0	1	0	0	1
0	1	1	1	0
1	0	0	0	1
1	0	1	1	0
1	1	0	1	0
1	1	1	1	1





## Ripple-Carry-Addierer:



FA = Full Adder

- Problem: **Lange Durchlaufzeit für den Übertrag** bei großer Stellenzahl, z.B. bei 32-Bit- oder 64-Bit-Addierer.
- Lösung: Parallelisierung der Übertragsbildung:  
→ Carry-Look-Ahead-Addierer (CLA-Addierer)



- Gleichungen (Summe von Produkten):

$$c_1 = b_0c_0 + a_0c_0 + a_0b_0$$

$$c_2 = b_1c_1 + a_1c_1 + a_1b_1 = b_1b_0c_0 + b_1a_0c_0 + b_1a_0b_0 + a_1b_0c_0 + a_1a_0c_0 + a_1a_0b_0 + a_1b_1$$

$$c_3 = b_2c_2 + a_2c_2 + a_2b_2 = \dots$$

$$c_4 = b_3c_3 + a_3c_3 + a_3b_3 = \dots$$

- Nicht umsetzbar! Warum?

- $VA_i$  wird immer ein Übertrag generieren, wenn  $g_i = a_i b_i (= 1)$

- $VA_i$  propagiert ein Übertrag, wenn  $p_i = a_i + b_i (= 1)$

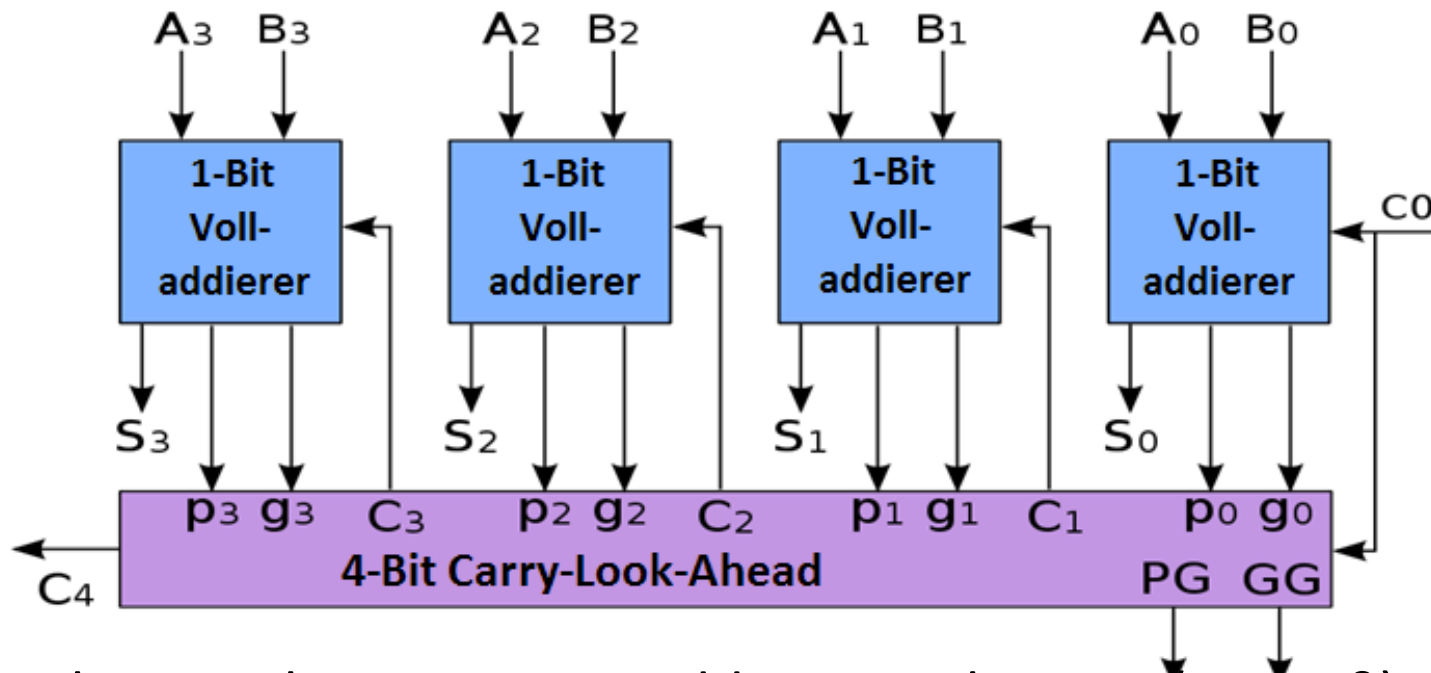
$$c_1 = g_0 + p_0c_0$$

$$c_2 = g_1 + p_1c_1 = g_1 + p_1g_0 + p_1p_0c_0$$

$$c_3 = g_2 + p_2c_2 = g_2 + p_2g_1 + p_2p_1g_0 + p_2p_1p_0c_0$$

$$c_4 = g_3 + p_3c_3 = \dots$$

- Umsetzbar! Warum?



- Man kann so keinen 16-Bit Addierer realisieren (zu groß).
- Man könnte den Übertrag durch ein Hintereinanderschalten von 4-Bit CLA-Addierern realisieren (Ripple-Carry-Prinzip).
- Besser: erneutes Anwenden des CLA -Prinzips!





- Für den Moment betrachten wir nur positive Zahlen.
- Schulmathematik:

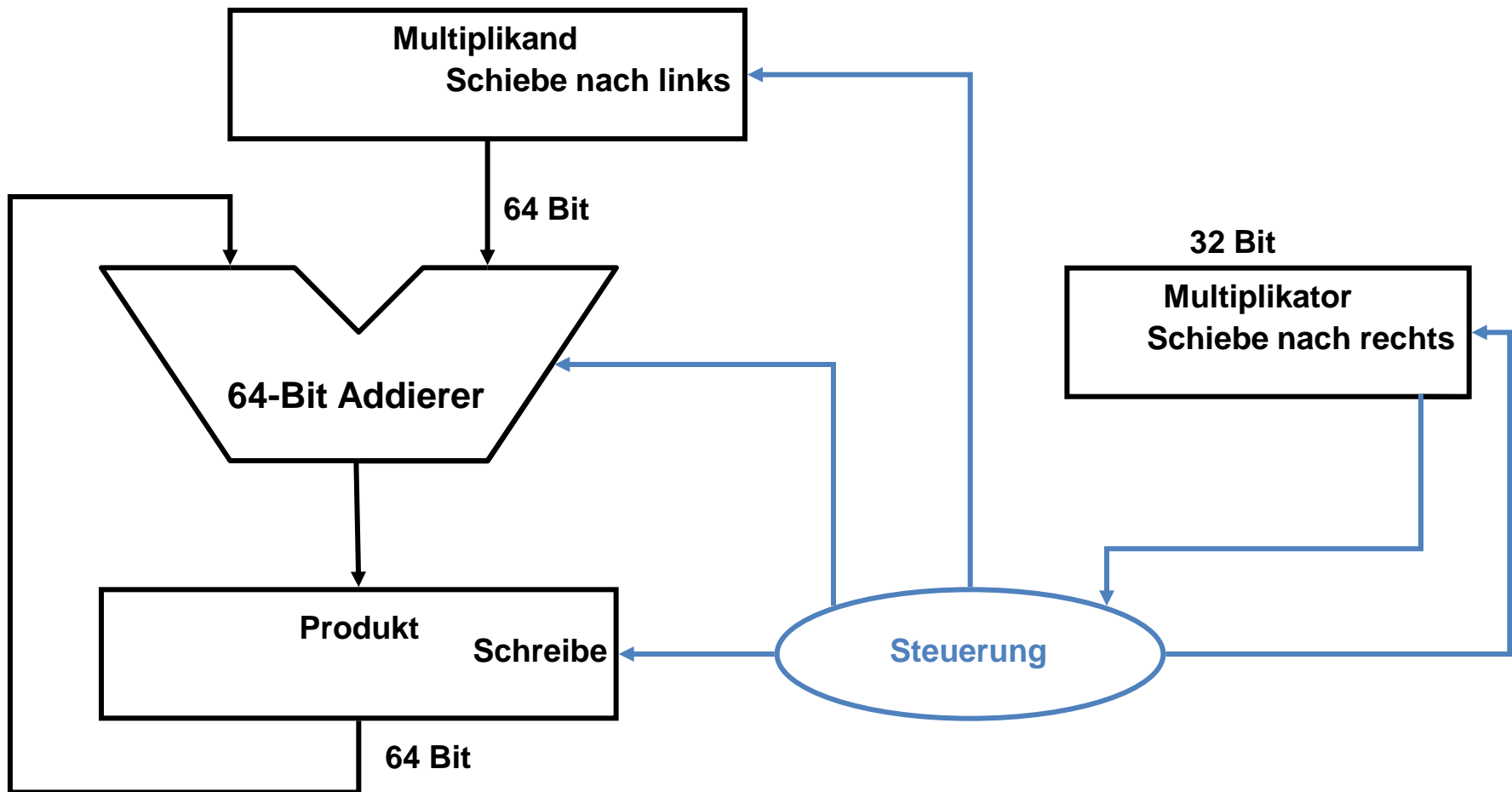
Multiplikand		1101		13
Multiplikator	x	1011		x 11
		1101		
		1101		
		0000		
		1101		
Produkt		10001111		143

- Produkt erfordert doppelte Stellenanzahl.

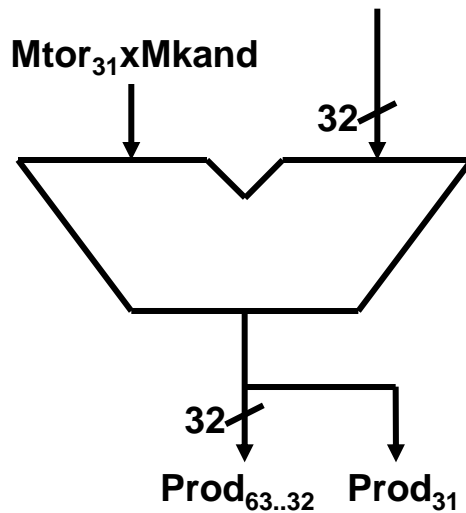
## C-Algorithmus für einen sequentiellen Multiplizierer

- Multiplikand:  $2n$ -Bit Register
- Multiplikator:  $n$ -Bit Register
- Produkt:  $2n$ -Bit Register
- Verfügen über  $2n$ -Bit Addierer

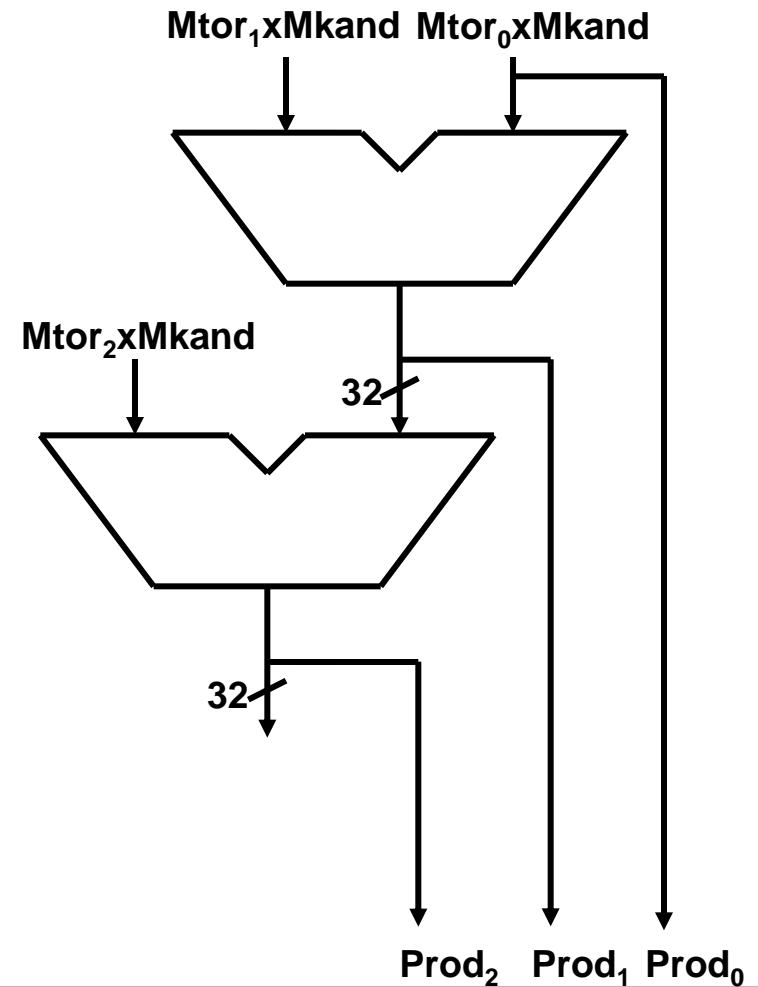
```
for (i=0; i<n; i++)  
{  
    if (Multiplikator&1)           // Multiplikator0==1  
        Produkt += Multiplikand;  
    Multiplikator >>= 1;           // Schiebe rechts 1 Bit  
    Multiplikand <<= 1;           // Schiebe links 1 Bit  
}
```



Hardware, die die Schleife „ausrollt“



...



- Das Multiplizieren von zwei 32-Bit Zahlen kann ein 64-Bit Produkt ergeben.
  - Neue Register: **Hi** und **Lo**
    - **Hi** beinhaltet die 32 MSBs des 64-Bit Produkts.
    - **Lo** beinhaltet die 32 LSBs.
  - MIPS-Befehle:
    - `mult $s2,$s3`      # `Hi#Lo = $s2x$s3`
    - `multu $s2,$s3`      # `Hi#Lo = $s2x$s3 (unsigned)`
    - *Move from lo*: `mflo $s1`      # `$s1 = Lo`
    - *Move from hi*: `mfhi $s1`      # `$s1 = Hi`
- Pseudo-Instruktion: `mul $s1,$s2,$s3`
- Reale Umsetzung:    `mult $s2,$s3`  
                           `mflo $s1`



- Schulmathematik:

Dividend

1001010

/

1000

Divisor

=

1

Quotient

1000

1

- Schulmathematik:

Dividend

1001010

/

1000

=

10

1000



10

Divisor

Quotient



- Schulmathematik:

Dividend		Divisor		Quotient
1001010	/	1000	=	100
1000				
<hr/>				
101				

*Note: A blue arrow points from the 4th bit of the dividend (1001010) to the 4th bit of the remainder (101).*





- Schulmathematik:

Dividend

$$\begin{array}{r}
 1001010 \\
 \underline{1000} \\
 1010 \\
 \underline{1000} \\
 10 \rightarrow \text{Rest} \\
 \underline{\phantom{10}} \\
 =2
 \end{array}$$

Divisor

$$\begin{array}{r}
 1001 \\
 \hline
 =9
 \end{array}$$

Quotient

$$74/8 = 9 \text{ Rest } 2$$



## C-Algorithmus für Division

- Quotient:  
n-Bit Register
- Divisor:  
2n-Bit Register, Divisor wird zu Beginn in der linken Hälfte abgelegt
- Rest:  
2n-Bit Register, wird mit Dividenden initialisiert
- Verfügen über 2n-Bit Subtrahierer

```
Quotient = 0;  
Rest = Dividend;  
Divisor <<= n;  
for (i=0; i<=n; i++){  
    Quotient <<= 1;  
    if (Rest>=Divisor){  
        Rest -= Divisor;  
        Quotient |= 1;  
    }  
    Divisor >>= 1;  
}
```



- $7_D$  durch  $2_D$ , d. h.  $0111_B$  durch  $0010_B$

Iteration	Quot .	Divisor	Rest	Anm.
Anfangsw.	0000	0010 0000	0000 0111	$R < D$
0	0000	0001 0000	0000 0111	$R < D$
1	0000	0000 1000	0000 0111	$R < D$
2	0000	0000 0100	0000 0111	$R \geq D$
3	0001	0000 0010	0000 0011	$R \geq D$
4	0011	0000 0001	0000 0001	

```

Quotient = 0;
Rest = Dividend;
Divisor <= n;
for (i=0; i<=n; i++){
    Quotient <= 1;
    if (Rest>=Divisor){
        Rest -= Divisor;
        Quotient |= 1;
    }
    Divisor >>= 1;
}

```

- Dividieren in MIPS:
  - Benutzt auch die Register **Hi** und **Lo**
    - Quotient befindet sich im **Lo**, Rest im **Hi**.
    - Rest wird als Nebeneffekt der Division mitproduziert.
- MIPS-Befehle:
  - `div $s2,$s3`      # `Lo = $s2/$s3, Hi = $s2%$s3`
  - `divu $s2,$s3`      # `idem unsigned`
- Pseudo-Befehle:
  - `div $s1,$s2,$s3`
  - `rem $s1,$s2,$s3`



- Darstellung gebrochener Zahlen:

$$a_{n-1}a_{n-2} \dots a_0, a_{-1}a_{-2} \dots a_{-m}$$

$$= a_{n-1}B^{n-1} + \dots + a_0B^0 + a_{-1}B^{-1} + a_{-2}B^{-2} + \dots + a_{-m}B^{-m}$$

$$= \sum_{i=0}^{n-1} a_i B^i + \sum_{i=-m}^{-1} a_i B^i$$

- Beispiel (binär):

$$\begin{aligned} 11.1010_B &= 3 + 2^{-1} + 2^{-3} \\ &= 3 + 0,5 + 0,125 = 3,625_D \end{aligned}$$

- Dezimal nach Dual

- Beispiel:  $0,24_D$

$$0,24_D \rightarrow 0.24 \cdot 2 = 0.48 + 0$$

$$0.48 \cdot 2 = 0.96 + 0$$

$$0.96 \cdot 2 = 0.92 + 1$$

$$0.92 \cdot 2 = 0.84 + 1$$

$$0.84 \cdot 2 = 0.68 + 1$$

$$0.68 \cdot 2 = 0.36 + 1$$

$$0.36 \cdot 2 = 0.72 + 0$$

$$0.72 \cdot 2 = 0.44 + 1$$

MSb!

LSb!

$\rightarrow .00111101_B$

- Abbruch nach 8 Stellen (Näherung mit  $0,238\dots$ )

- Näherung für reelle Zahlen: 
$$(-1)^s \times 1.f \times 2^E$$
  - $s$ : **Vorzeichen** (*sign*): 0  $\rightarrow$  positiv, 1  $\rightarrow$  negativ
  - $1.f$ : **Mantisse** (Betrag, *significand*) als **normalisierte Zahl**
    - Zahl wird so lange geschoben, bis sie führende 1 aufweist
    - Binärpunkt wird rechts von dieser 1 festgelegt ( $1.0 \leq 1.f < 2.0$ )
  - $f$ : nur der **Bruch  $f$**  (*fraction*) wird gespeichert, **führende 1** ist implizit (wird von Recheneinheit ergänzt)
  - $E$ : vorzeichenbehafteter Exponent, wird als **transformierter Exponent  $e$**  gespeichert
  - $e$ :  $e = E + \text{bias}$ 
    - **bias** wird so gewählt, dass 2-Komplement-Zahl  $E$  zur vorzeichenlosen Dualzahl  $e$  wird



- Einfache Genauigkeit (*single precision*, 32 Bit)



- bias = 127
- C/Java: `float`

- Doppelte Genauigkeit (*double precision*, 64 Bit)

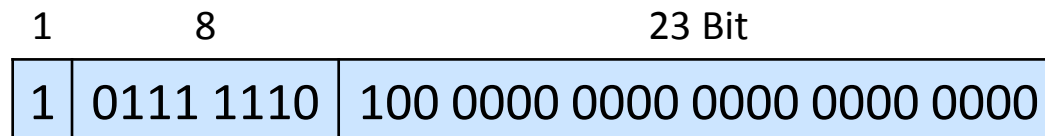


- bias = 1023
- C/Java: `double`



- $-0,75_D$  mit einfacher Genauigkeit
  - $s = 1$
  - $0,75_D$  als gebrochene Dualzahl ist  $0.11_B$
  - Normalisiere:  $0.11 = 1.1 \times 2^{-1}$ 
    - führende 1 ist implizit  $\rightarrow$  Bruch = 10000....
  - transformierter Exponent  $e$ 

$$e = E + \text{bias} = -1 + 127 = 126 = 0111\ 1110_B$$



- Neben **normalisierten Zahlen** sind außerhalb des Zahlenraums definiert:
  - $\pm$  Null
  - $\pm$  Unendlich: z. B. Division durch Null
  - $\pm$  unnormalisierte (*unnormalized*) Zahlen: winzige (*tiny*) Zahlen
  - Nichtzahlen (Not a Number, NaN): Ergebnis ungültiger Operation wie  $0/0$
- Codiert durch den größten und kleinsten Exponentwert  $e$  und  $f$ :
  - normal.:  $(-1)^s \times 1.f \times 2^{e-127/-1023}$   $1 \leq e \leq 254/2046$
  - Null:  $(-1)^s \times 0$   $e = 0, f = 0$
  - Unendlich:  $(-1)^s \times \infty$   $e = 255/2047, f = 0$
  - unnorm.:  $(-1)^s \times 0.f \times 2^{-126/-1022}$   $e = 0, f \neq 0$ ; interpretiert mit  $e = 1$
  - Nichtzahl: NaN  $e = 255/2047, f \neq 0$

- Beispiel basiert auf 16-Bit **Minifloat** Format:



bias = 15

Exponentenbereich:  $-14 \leq E \leq 15$

- $Z = X + Y$  mit
  - $X = 2,35_D = 10.0101\ 1001\ 1001\ 1001 \dots_B$
  - $Y = 10,17_D = 1010.0010\ 1011\ 1000\ 0101 \dots_B$
- 1. Schritt:**
  - Normalisieren** und **Anpassung** an 16-Bit-Format:
  - $X = 1.0010\ 1100\ 11 \cdot 2^1$
  - $Y = 1.0100\ 0101\ 01 \cdot 2^3$
  - Verlust von signifikanten Stellen!

- **2. Schritt:**

- **Vergleichen der beiden Exponenten  $E$ .**
- Bei Ungleichheit kleinere Exponent an den größeren anpassen
- $X = 0.0100\ 1011\ 00\ \textcolor{red}{11} \cdot 2^3$   
 → **rot** dargestellten Stellen gehen verloren

- **3. Schritt:**

- **Addieren der Mantissen:**

$$\begin{array}{r}
 0.0100\ 1011\ 00\ (X) \\
 +\ 1.0100\ 0101\ 01\ (Y) \\
 \hline
 1.1001\ 0000\ 01\ (Z)
 \end{array}$$

- **Ergebnis**

- muss ggf. noch normalisiert werden (hier nicht)
- $Z = 1.1001\ 0000\ 01 \cdot 2^3 = 12,500\ 975\ 656\ 2510$  (korrekt wäre: 12,52).

- MIPS hat
  - 32 Floating-Point-Register mit einfacher Genauigkeit (single-precision)  
[\$f0, \$f1, ..., \$f31] **oder**
  - 16 Register mit doppelter Genauigkeit (double-precision)  
[\$f0, \$f2, ..., \$f30]

single-precision

double-precision

\$f0		\$f0
\$f1		
\$f2		\$f2
\$f3		
. . .		. . .
\$f30		\$f30
\$f31		



FP add single	<code>add.s \$f0, \$f1, \$f2</code>	$\$f0 = \$f1 + \$f2$
FP sub. single	<code>sub.s \$f0, \$f1, \$f2</code>	$\$f0 = \$f1 - \$f2$
FP mult. single	<code>mul.s \$f0, \$f1, \$f2</code>	$\$f0 = \$f1 * \$f2$
FP div. single	<code>div.s \$f0, \$f1, \$f2</code>	$\$f0 = \$f1 / \$f2$
FP add double	<code>add.d \$f0, \$f2, \$f4</code>	$\$f0, \$f1 = \$f2, \$f3 + \$f4, \$f5$
FP sub. double	<code>sub.d \$f0, \$f2, \$f4</code>	$\$f0, \$f1 = \$f2, \$f3 - \$f4, \$f5$
FP mult. double	<code>mul.d \$f0, \$f2, \$f4</code>	$\$f0, \$f1 = \$f2, \$f3 * \$f4, \$f5$
FP div. double	<code>div.d \$f0, \$f2, \$f4</code>	$\$f0, \$f1 = \$f2, \$f3 / \$f4, \$f5$
load word coproc. 1	<code>lwc1 \$f0, 100(\$s0)</code>	$\$f0 = \text{Mem}[\$s0 + 100]$
store word copr. 1	<code>swc1 \$f0, 100(\$s0)</code>	$\text{Mem}[\$s0 + 100] = \$f0$
branch on coproc.1 true	<code>bc1t 25</code>	if (cond) goto PC+4+100
branch on coproc.1 false	<code>bc1f 25</code>	if (!cond) goto PC+4+100
FP compare single	<code>c.lt.s \$f0, \$f1</code>	cond = ( $\$f0 < \$f1$ )
FP compare double	<code>c.ge.d \$f0, \$f2</code>	cond = ( $\$f0, \$f1 \geq \$f2, \$f3$ )

- Vergleich-Befehle setzen den **condition code (cc)**
- Sukzessive Branch-Befehle testen, ob cc erfüllt ist (*true*) oder nicht (*false*)
- Beispiel: Suche nach kleinstem  $n$ , so dass gilt  $0.5^n \leq 1.0 \cdot 10^{-9}$

```
int n = 1;
float exp = 0.5;
while (exp > 1e-9) {
    exp = exp*0.5;
    n++;
}
```

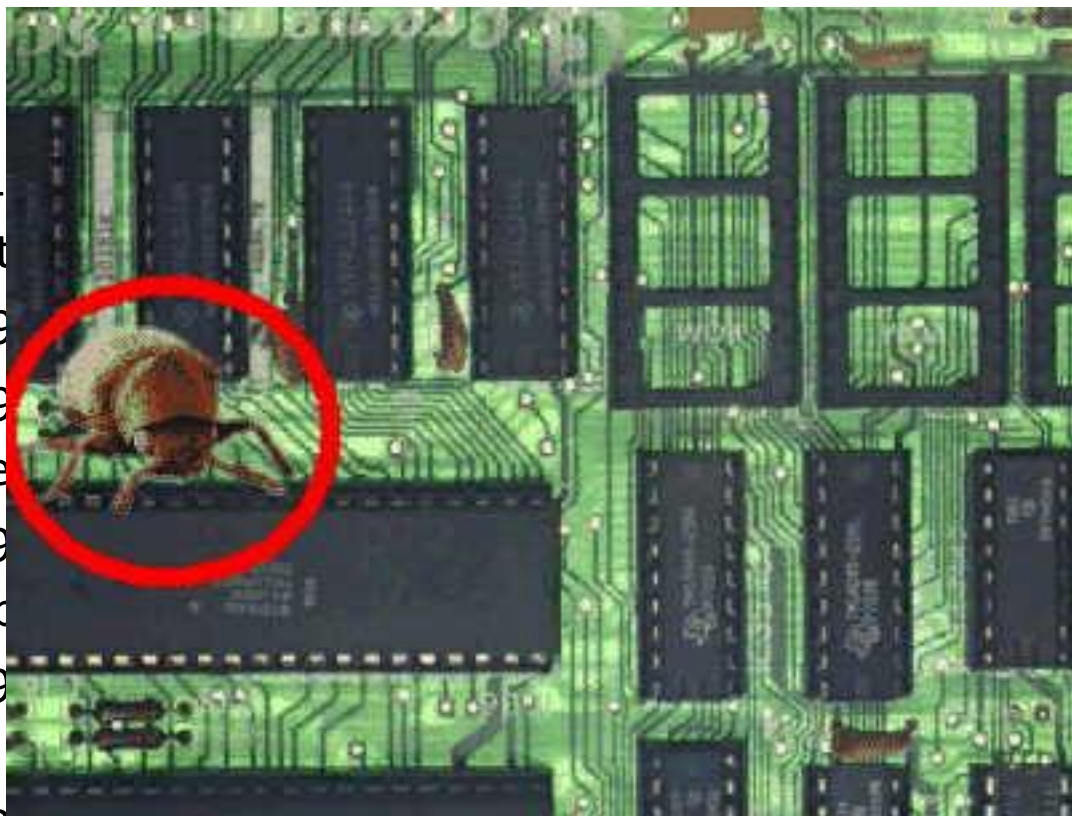
```
int n = 1;
float exp = 0.5;
while (exp > 1e-9) {
    exp = exp*0.5;
    n++;
}
```

```
addi    $t0,$zero,1        # n = 1
lwc1    $f0,fphalf($gp)    # exp = 0.5
lwc1    $f1,fptiny($gp)    # $f1 = 1e-9
lwc1    $f2,fphalf($gp)    # $f2 = 0.5
while:
    c.gt.s    $f0,$f1        # cc = exp>1e-9
    bclf     endwhile       # if (!cc) goto endwhile
    mul.s    $f0,$f0,$f2     # exp = exp*0.5
    addi    $t0,$t0,1        # n++
    j        while          # goto while
endwhile:    ...
```



## Fehler im Divisionsalgorithmus für Gleitpunktzahlen

- Juli 1994: Intel entdeckt Fehler: **einige** 100K\$
- Sept. 1994: Intel gibt offizielle Stellungnahme: keine Internet
- 7. Nov. 1994: Intel gibt offizielle Stellungnahme: an 9. Stelle
- 22. Nov. 1994: Intel gibt offizielle Stellungnahme: 1000 Jahren
- 5. Dez. 1994: Intel gibt offizielle Stellungnahme: auftreten b
- 12. Dez. 1994: Intel gibt offizielle Stellungnahme: 24 Tage Fehler
- 21. Dez. 1994: Intel gibt zu: Jeder Besitzer darf Pentium austauschen. Geschätzte Kosten: **500 M\$!**



$4195835.0/3145727.0 = 1.333\ 820\ 449\ 136\ 241\ 002$  (korrekter Wert)

$4195835.0/3145727.0 = 1.333\ 739\ 068\ 902\ 037\ 589$  (fehlerhaften Pentium)

- Computer-Arithmetik ist beschränkt durch **limitierte Genauigkeit**.
  - *Overflow* (Überlauf)
  - Underflow (floating point)
- Bitmuster haben keine inherente Bedeutung, es gibt jedoch **Standards**
  - 2-Komplement
  - IEEE 754 Floating Point
- Computer-Befehle bestimmen die “Bedeutung” der Bitmuster.
- Leistung und Genauigkeit sind wichtig. Daraus ergeben sich viele Komplexitäten für reale Maschinen.
  - z.B. Algorithmen und Implementierungen
- Was kommt als Nächstes?
  - Wir werden einen Prozessor implementieren.
  - Zuerst müssen wir jedoch “Leistung” verstehen.