

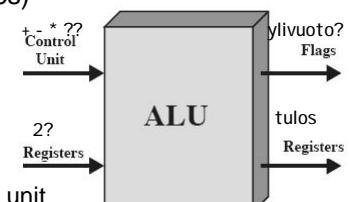


Lecture 5

Computer Arithmetic (Tietokonearitmetiikka)

Stallings: Ch 9

- Integer representation (*Kokonaislukuesitys*)
- Integer arithmetic (*Kokonaislukuaritmetiikka*)
- Floating-point representation (*Liukulukuesitys*)
- Floating-point arithmetic (*Liukulukuaritmetiikka*)



ALU

- ALU = Arithmetic Logic Unit (*Aritmeettis-looginen yksikkö*)
- Actually performs operations on data
 - Integer and floating-point arithmetic
 - Comparisons (*vertailut*), left and right shifts (*sivuttaissiirrot*)
 - Copy bits from one register to another
 - Address calculations (*Osoitelaskenta*): branch and jump (*hypyt*), memory references (*muistiviittaukset*)
- Data from/to internal registers (latches)
 - Input copied from normal registers (or from memory)
 - Output goes to reg (or memory)
- Operation
 - Based on instruction register, control unit

(Sta06 Fig 9.1)

Computer Organization II, Spring 2010, Tiina Niklander 29.1.2010 2



Computer Organization II

Integer representation (*kokonaislukujen esitys*)

Computer Organization II, Spring 2010, Tiina Niklander

29.1.2010

3



Integer Representation (*Kokonaislukuesitys*)

- Binary representation, bit sequence, only 0 and 1
- "Weight" of the number based on position

$$\begin{aligned} 57 &= 5 \cdot 10^1 + 7 \cdot 10^0 \\ &= 32 + 16 + 8 + 1 \\ &= 1 \cdot 2^5 + 1 \cdot 2^4 + 1 \cdot 2^3 + 0 \cdot 2^2 + 0 \cdot 2^1 + 1 \cdot 2^0 \\ &= 0011\ 1001 \\ &= \underline{0x39} \quad \text{hexadecimal} \\ &= 3 \cdot 16^1 + 9 \cdot 16^0 \end{aligned}$$

- Most significant bit, MSB (*eniten merkitsevä bitti*)
- Least significant bit, LSB (*vähiten merkitsevä bitti*)

Computer Organization II, Spring 2010, Tiina Niklander

29.1.2010

4

Integer Representation (Kokonaislukuesitys)

- Negative numbers?
 - Sign magnitude (*Etumerkki-suuruus*)
 - Twos complement (*2:n komplementtimuoto*)

$-57 = \underline{1}011\ 1001$ $-57 = \underline{1}100\ 0111$

Sign
(*etumerkki*)

- Computers use twos complement
 - Just one zero (no +0 and -0)
 - Comparison to zero easy
 - Math is easy to implement
 - No need to consider sign
 - Subtraction becomes addition
 - Simple hardware and circuit

$+2 = 0000\ 0010$
$+1 = 0000\ 0001$
$0 = 0000\ 0000$
$-1 = 1111\ 1111$
$-2 = 1111\ 1110$

Computer Organization II, Spring 2010, Tiina Niklander 29.1.2010 5

Twos complement (2:n komplementti)

- Example
 - 8-bit sequence, value -57

$57 = 0011\ 1001$	unsigned value (<i>itseisarvo</i>)
$1100\ 0110$	invert bit (ones complement)
$\begin{array}{r} 1100\ 0110 \\ \hline 1 \end{array}$	add 1 twos complement
	Reject overflow

- Easy to expand. As a 16-bit sequence

$57 = \underline{0}011\ 1001 = \underline{0}000\ 0000\ \underline{0}011\ 1001$	
$-57 = \underline{1}100\ 0111 = \underline{1}111\ 1111\ \underline{1}100\ 0111$	sign extension

Computer Organization II, Spring 2010, Tiina Niklander 29.1.2010 6



Twos complement

- Value range (arvoalue): $-2^{n-1} \dots 2^{n-1} - 1$

8 bits: $-2^7 \dots 2^7 - 1 = -128 \dots 127$

32 bits: $-2^{31} \dots 2^{31} - 1 = -2\ 147\ 483\ 648 \dots 2\ 147\ 483\ 647$

- Addition overflow (*yhteenlaskun ylivooto*) easy to detect
 - No overflow, if different signs in operands
 - Overflow, if same sign (*etumerkki*) and the results sign differs from the operands

$$\begin{array}{r} 57 = 0011\ 1001 \\ + 80 = 0101\ 0000 \\ \hline \end{array}$$

$$137 = \underline{1}000\ 1001 \quad \text{Overflow!}$$

Computer Organization II, Spring 2010, Tiina Niklander

29.1.2010

7



Twos complement

- Subtraction as addition (*vähennyslasku yhteenlaskuna*)!

- Forget the sign, handle as if unsigned!
- Complement 2nd term, the subtrahend, then add (*lisää 2:n komplementti vähentäjästä*)
- Simple hardware

$$\begin{array}{r} +1 = 0001 \\ -3 = 1101 \\ \hline -2 = 1110 \end{array}$$

$$3 = 0011$$

$$\begin{array}{r} 1100 \\ 1 \\ \hline 1101 \end{array}$$

-3 in two complement

- Check

- Overflow? (same rule as in addition)
- sign= 1, result is negative

(Sta06 Table 9.1)

Computer Organization II, Spring 2010, Tiina Niklander

29.1.2010

8



Computer Organization II

Integer arithmetics (kokonaislukuaritmetiikka)

Negation (*negaatio*)
 Addition (*yhteenlasku*)
 Subtraction (*vähennyslasku*)
 Multiplication (*kertolasku*)
 Division (*jakolasku*)

Computer Organization II, Spring 2010, Tiina Niklander

29.1.2010

9



Negation = Twos complement

- 1: invert all bits
- 2: add 1
- 3: Special cases
 - Ignore carry bit (*ylivuotobitti*)
 - Sign really changed?
 - Cannot negate smallest negative
 - Result in exception
- Simple hardware

$$\begin{array}{r} -57 = \underline{1}100\ 0111 \\ \quad 0011\ 1000 \\ \hline \quad \quad \quad 1 \\ \quad 0011\ 1001 \\ = 57 \end{array}$$

$$\begin{array}{r} -128 = \underline{1}000\ 0000 \\ \quad 0111\ 1111 \\ \hline \quad \quad \quad 1 \\ \quad 1000\ 0000 \end{array}$$

Computer Organization II, Spring 2010, Tiina Niklander

29.1.2010

10

Addition (and subtraction)

- Normal binary addition
 - In subtraction: complement the 2. operand, subtrahend (*vähentäjä*) and add to 1. operand, minuend (*vähennettävä*)
- Ignore carry
 - Check sign!

Overflow indication
- Simple hardware function
 - Two circuits:
 - Complement and addition

$$\begin{array}{r}
 1100 = -4 \\
 +1111 = -1 \\
 \hline
 1011 = ?
 \end{array}$$

Overflow

(Sta06 Fig 9.6)

Computer Organization II, Spring 2010, Tiina Niklander 29.1.2010 11

Integer multiplication

- "Just like" you learned at school
 - Easy with just 0 and 1!
- Hardware?
 - Complex
 - Several algorithms
- Overflow?
 - 32 b operands → result 64 b?
- Simpler, if only unsigned numbers
 - Just multiple additions
 - Or additions and shifts
 - Shift left = multiply by 2
 - esim: $5 * \Rightarrow$ add, shift, shift, add

$ \begin{array}{r} 1011 \\ \times 1101 \\ \hline 1011 \\ 0000 \\ 1011 \\ 1011 \\ \hline 10001111 \end{array} $	Multiplicand (11) Multiplier (13) Partial products Product (143)
--	---

(Sta06 Fig 9.7)

2 * 10011 => 100110

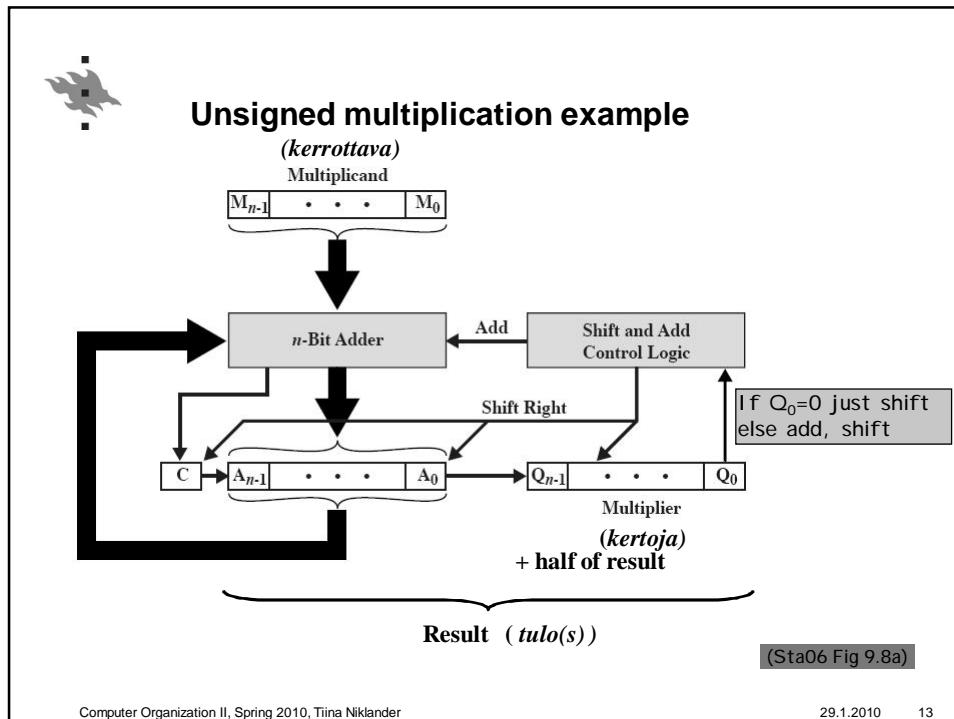
Example:

```

add=> 1011
shift=> 10110
shift=> 101100
add=>110111 (= 55)

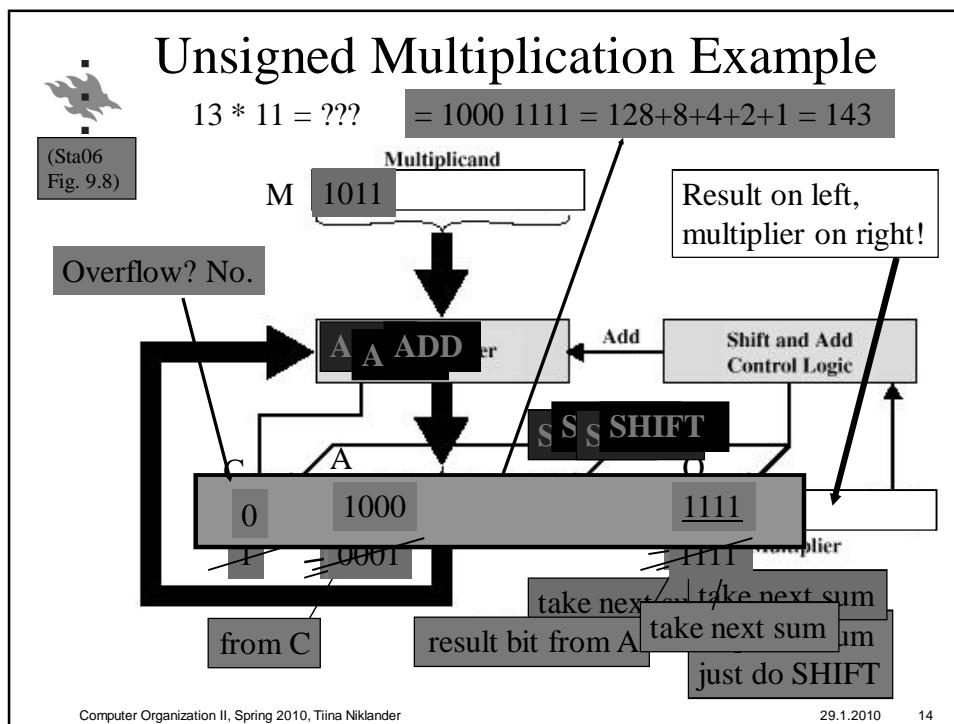
```

Computer Organization II, Spring 2010, Tiina Niklander 29.1.2010 12



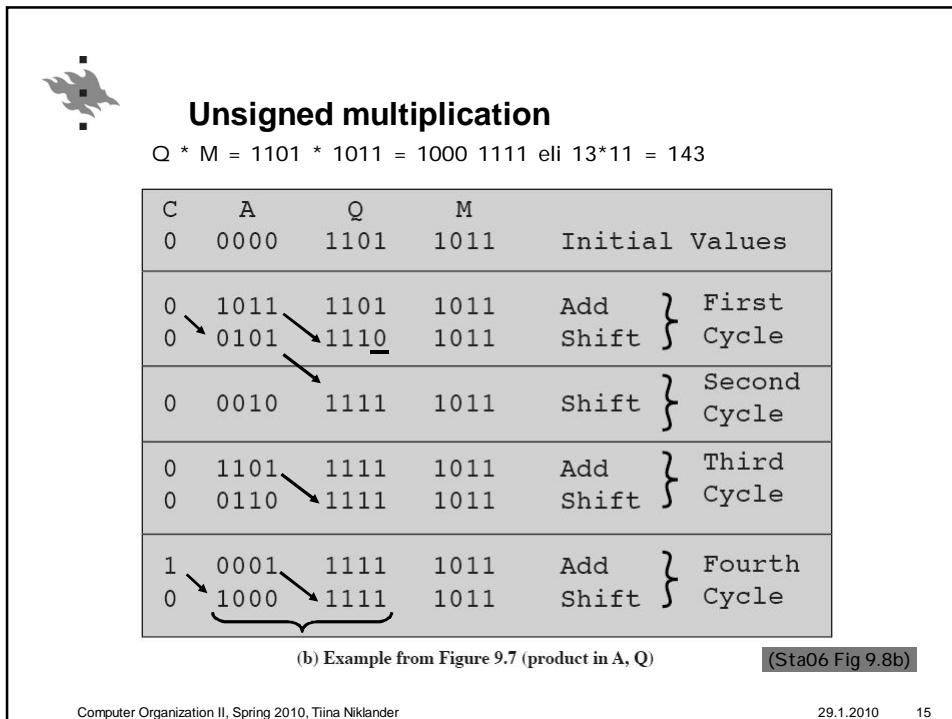
Computer Organization II, Spring 2010, Tiina Niklander

29.1.2010 13



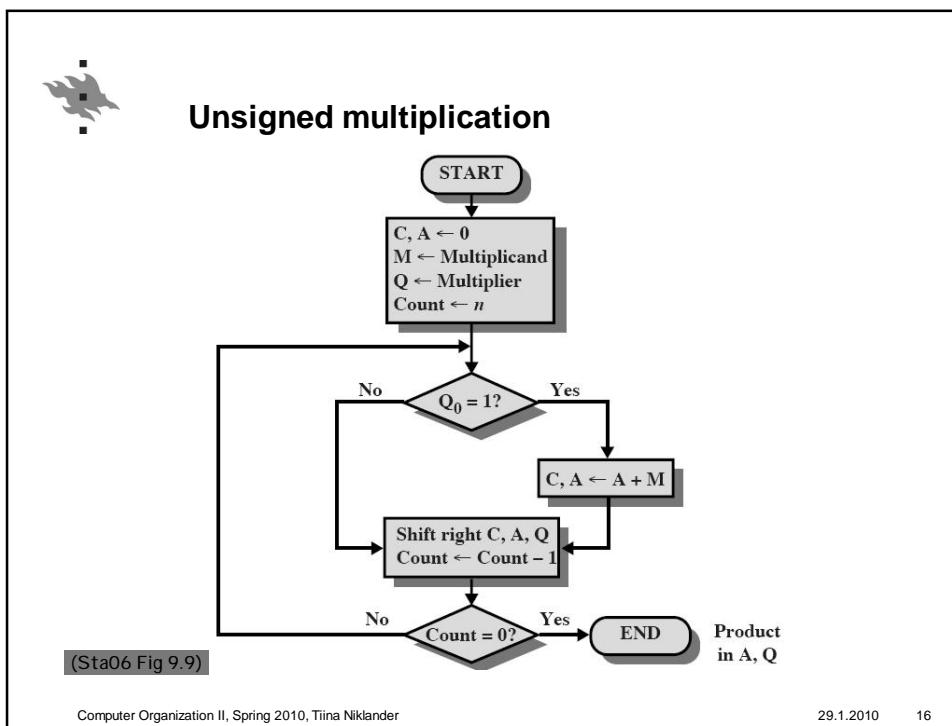
Computer Organization II, Spring 2010, Tiina Niklander

29.1.2010 14



Computer Organization II, Spring 2010, Tiina Niklander

29.1.2010 15





Multiplication with negative values?

- The preceding algorithm for unsigned numbers does NOT work for negative numbers
- Could do with unsigned numbers
 - ① Change operands to positive values
 - ② Do multiplication with positive values
 - ③ Check signs and negate the result if needed
- This works, but there are better and faster mechanisms available

Computer Organization II, Spring 2010, Tiina Niklander

29.1.2010 17



Booth's Algorithm

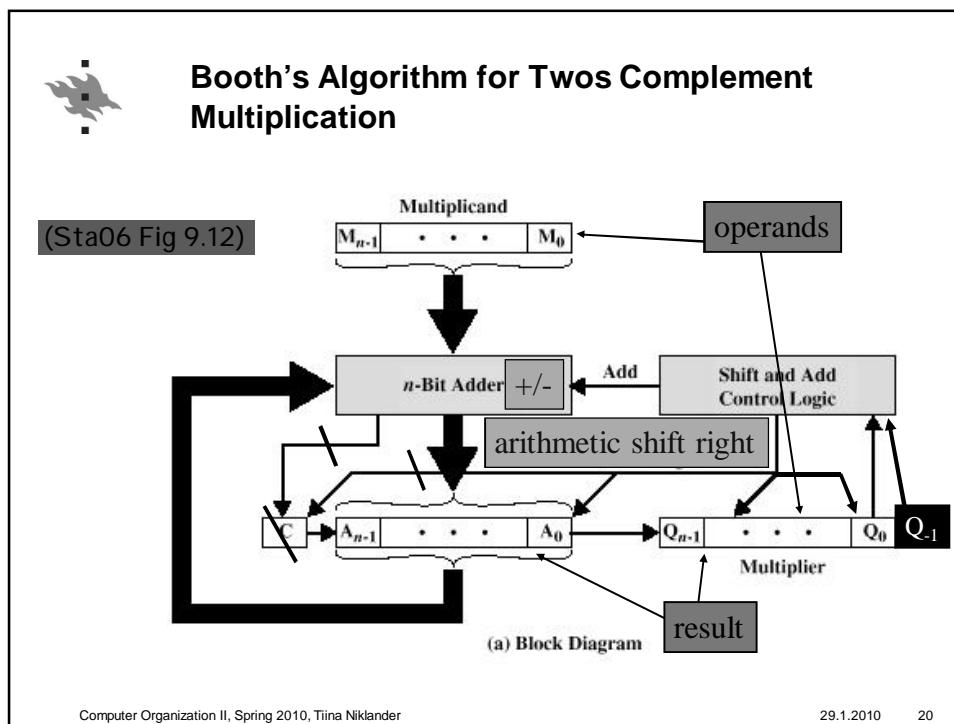
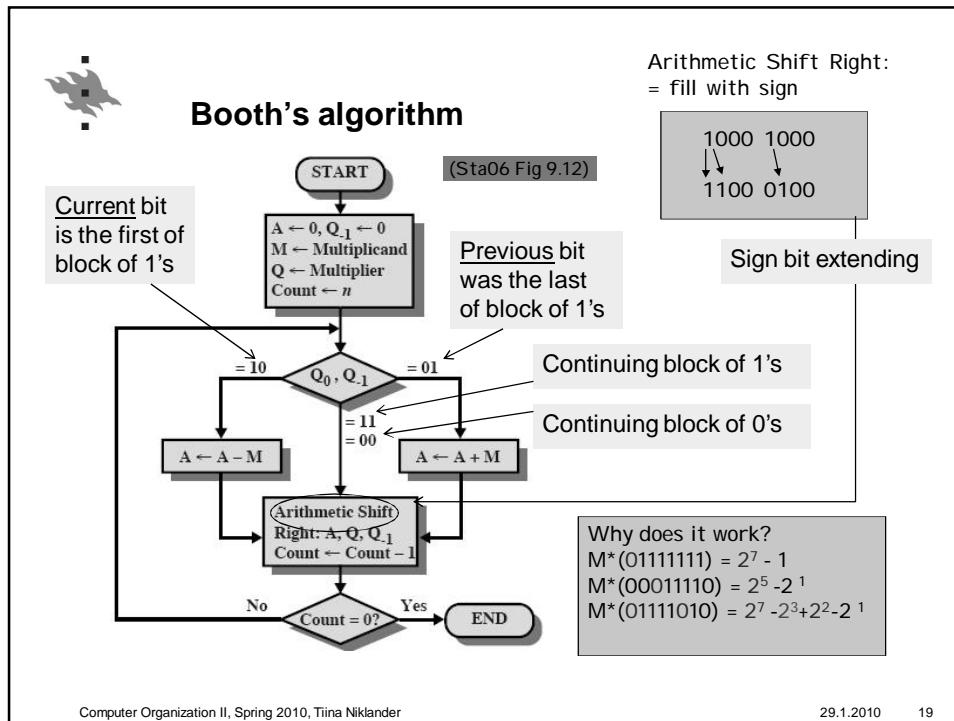
- Unsigned multiplication:
 - Addition (only) for every "1" bit in multiplier (*kertoja*)
- Booth's algorithm (improvement)
 - Combine all adjacent 1's in multiplier together,
 - Replace all additions by one subtraction and one addition
 - Example: $7 \times x = 8 \times x + (-x)$
 - $111 \times x = 1000 \times x + (-x) =$
 - add, shift, shift, shift, complement, add
(in reality, the complement would be first)

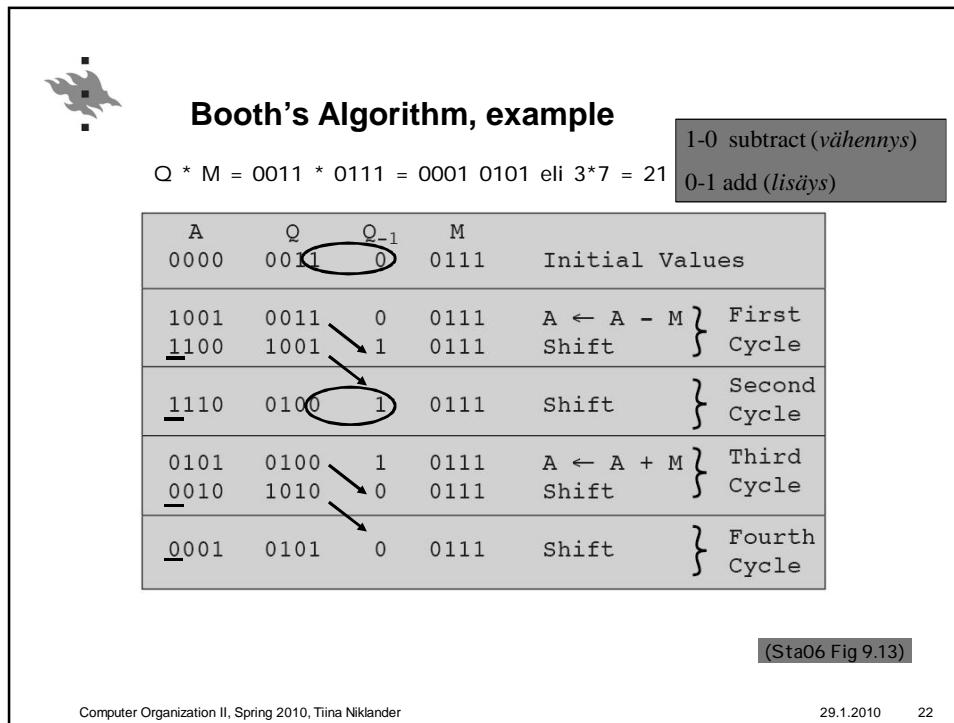
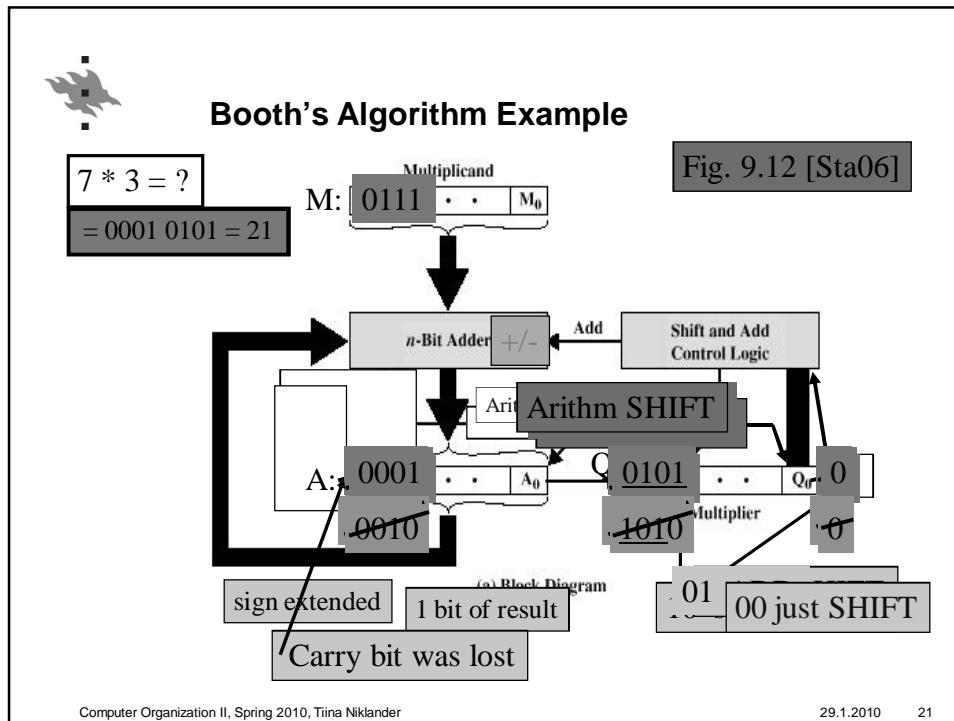
$$\begin{array}{l}
 5 * 7 = 0101 * 0111 \\
 = 0101 * (1000-0001) \quad \text{(circled 1000-0001)}
 \end{array} \rightarrow
 \begin{array}{r}
 00101000 \quad 40 \\
 -1111011 \quad -5 \\
 \hline
 100100011 = 35
 \end{array}$$

- Works for two's complement! Also negative values!

Computer Organization II, Spring 2010, Tiina Niklander

29.1.2010 18





Integer division (*kokonaislukujen jakolasku*)

- Like in school algorithm
 - Easy: new quotient digit always 0 or 1

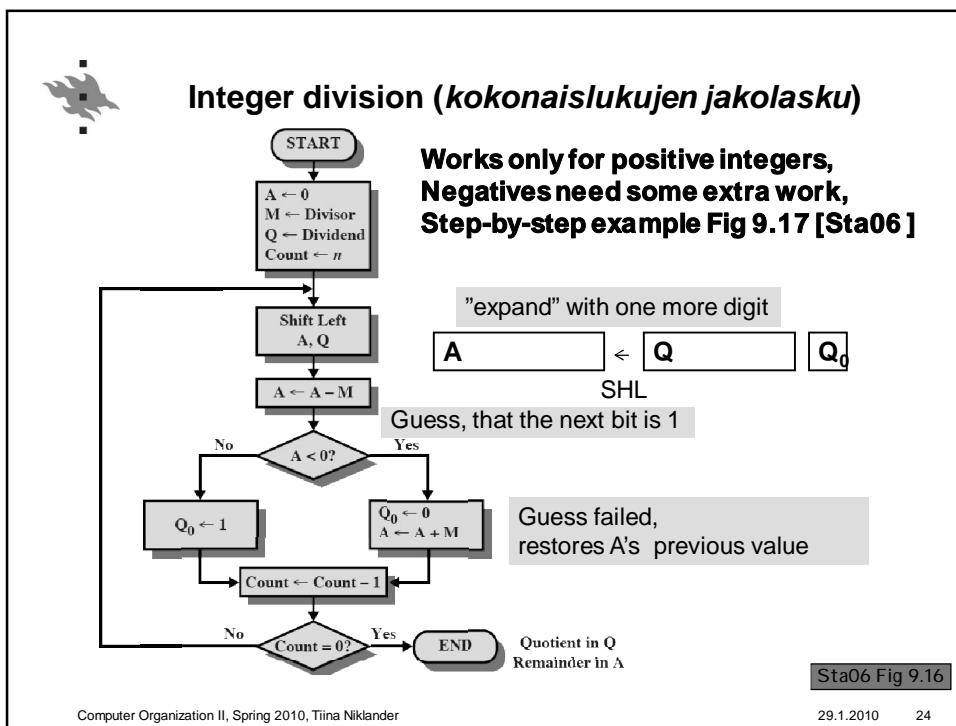
The diagram illustrates the integer division process. It shows the quotient (00001101), dividend (10010011), divisor (1011), and remainder (100). Partial remainders are shown as circled numbers: 001110, 001111, and 1011. Arrows indicate the shift left operation used to consider new digits.

- Hardware needs as in multiplication
 - Shift left = consider new digit

(Sta06 Fig 9.15)

Computer Organization II, Spring 2010, Tiina Niklander

29.1.2010 23



Computer Organization II, Spring 2010, Tiina Niklander

29.1.2010 24

Example: twos complement division

■ Division: 7/3 A+ Q = 7 = 0000 0111 M= 3 = 0011

A	Q	
0000	0111	initial value
0000	1110	shift left
1101	1110	subtract M
0000	1110	restore
0001	1100	shift left
1110	1100	subtract M
0001	1100	restore
0011	1000	shift left
0000	1000	subtract M
0000	1001	set $Q_0=1$
0001	0010	shift
1110	0010	subtract M
0001	0010	restore

Sta09 Fig 9.17

Subtract M = Add (-M)
 $-M = -3 = 1101$

First try, if you can do the subtraction
 (or add if different signs).
 If the sign changed, subtraction failed
 and A must be restored, $Q_0 = 0$

If subtraction successful, $Q_0 = 1$

$Q = \text{quotient} = 2$
 $A = \text{remainder} = 1$

Repeat as many times as Q has bits.

Computer Organization II, Spring 2010, Tiina Niklander

29.1.2010 25

Computer Organization II

Floating Point Representation *(Liukulukuesitys)*

Computer Organization II, Spring 2010, Tiina Niklander

29.1.2010 26

Floating Point Representation

The diagram illustrates the floating-point representation of a 32-bit word. It is divided into three fields: a 1-bit sign of significand, an 8-bit biased exponent, and a 23-bit significand or mantissa.

- Sign of significand: 1 bit
- Biased exponent: 8 bits
- Significand or mantissa: 23 bits

Key points:

- Significant digits (*Merkitsevät numerot*) and exponent (*suuruusluokka*)
- Normalized number (*Normeerattu muoto*)
 - Most significant digit is nonzero >0
 - Commonly just one digit before the radix point (*desim. pilkku*)

```

-0.000 000 000 123 = -1.23 * 10^-10
0.123 = +1.23 * 10^-1
123.0 = +1.23 * 10^2
123 000 000 000 000 = +1.23 * 10^14
  
```

Computer Organization II, Spring 2010, Tiina Niklander 29.1.2010 27

IEEE 754 (floating point) formats

Parameter	Single	Single Extended	Double	Double Extended
Word width (bits)	32	≥ 43	64	≥ 79
Exponent width (bits)	8	≥ 11	11	≥ 15
Exponent bias	127	unspecified	1023	unspecified
Maximum exponent	127	≥ 1023	1023	≥ 16383
Minimum exponent	-126	≤ -1022	-1022	≤ -16382
Number range (base 10)	$10^{-38}, 10^{+38}$	unspecified	$10^{-308}, 10^{+308}$	unspecified
Significand width (bits)*	23	≥ 31	52	≥ 63
Number of exponents	254	unspecified	2046	unspecified
Number of fractions	2^{23}	unspecified	2^{52}	unspecified
Number of values	1.98×2^{31}	unspecified	1.99×2^{63}	unspecified

* not including implied bit

(Sta06 Table 9.3)

Computer Organization II, Spring 2010, Tiina Niklander 29.1.2010 28



32-bit floating point

- 1 b sign
 - 1 = “-”, 0 = “+”
- 8 b exponent
 - Biased representation, no sign (*Ei etumerkkiä, vaan erillinen nollataso*)
 - Exp=5 → store 127+5, Exp=-5 → store 127-5 (bias127)
- 23 b significant (*mantissa*)
 - In normalized form the radix point is preceded with 1, which is not stored. (hidden bit, Zuse Z3 1939)
- The binary value of the floating point representation

$$\text{-1 Sign} * \text{1.Mantissa} * \text{2}^{\text{Exponent-127}}$$



Example

$$23.0 = +10111.0 * 2^0 = +1.0111 * 2^4 = ?$$

$$127+4=131$$

0	1000 0011	011 1000 0000 0000 0000 0000
sign	exponent	mantissa

$$1.0 = +1.0000 * 2^0 = ?$$

$$0+127=127$$

0	0111 1111	000 0000 0000 0000 0000 0000
sign	exponent	mantissa

Example

$X = ?$

$X = (-1)^0 * 1.1111 * 2^{(128-127)}$

$= 1.1111_2 * 2$

$= (1 + 1/2 + 1/4 + 1/8 + 1/16) * 2$

$= (1 + 0.5 + 0.25 + 0.125 + 0.0625) * 2$

$= 1.9375 * 2$ $= 3.875$

Computer Organization II, Spring 2010, Tiina Niklander

29.1.2010 31

Accuracy (*tarkkuus*) (32b)

- Value range (*arvoalue*)
 - 8 b exponent $\rightarrow 2^{126} \dots 2^{127} \sim -10^{-38} \dots 10^{38}$
- Not exact value
 - 24 b mantissa $\rightarrow 2^{24} \sim 1.7 \cdot 10^{-7} \sim 6$ decimals
- Balancing between range and precision

Numerical errors: Patriot Missile (1991), Ariane 5 (1996)
<http://ta.twi.tudelft.nl/nw/users/vuik/wi211/disasters.html>

Computer Organization II, Spring 2010, Tiina Niklander

29.1.2010 32



Interpretation of IEEE 754 Floating-Point Numbers

Single Precision (32 bits)				
	Sign	Biased exponent	Fraction	Value
positive zero	0	0	0	0
negative zero	1	0	0	-0
plus infinity	0	255 (all 1s)	0	∞
minus infinity	1	255 (all 1s)	0	$-\infty$
quiet NaN	0 or 1	255 (all 1s)	$\neq 0$	NaN
signaling NaN	0 or 1	255 (all 1s)	$\neq 0$	NaN
positive normalized nonzero	0	$0 < e < 255$	f	$2^{e-127}(1.f)$
negative normalized nonzero	1	$0 < e < 255$	f	$-2^{e-127}(1.f)$
positive denormalized	0	0	$f \neq 0$	$2^{e-126}(0.f)$
negative denormalized	1	0	$f \neq 0$	$-2^{e-126}(0.f)$

Not a Number

Double Precision similarly

(Sta06 Table 9.4)

Computer Organization II, Spring 2010, Tiina Niklander

29.1.2010 33



NaN: Not a Number

Operation	Quiet NaN Produced by
Any	Any operation on a signaling NaN
Add or subtract	Magnitude subtraction of infinities: $(+\infty) + (-\infty)$ $(-\infty) + (+\infty)$ $(+\infty) - (+\infty)$ $(-\infty) - (-\infty)$
Multiply	$0 \times \infty$
Division	$\frac{0}{0}$ or $\frac{\infty}{\infty}$
Remainder	$x \text{ REM } 0$ or $\infty \text{ REM } y$
Square root	\sqrt{x} where $x < 0$

(Sta06 Table 9.6)

Computer Organization II, Spring 2010, Tiina Niklander

29.1.2010 34



Computer Organization II

Floating Point Arithmetics (Liukulukuaritmetiikka)

IEEE-754 Standard
Addition
Subtraction
Multiplication
Division

Computer Organization II, Spring 2010, Tiina Niklander

29.1.2010 35



Floating point arithmetics

- Calculations need wide registers
 - Guard bits - pad right end of significand
 - More bits for the significand (mantissa)
 - Using Denormalized formats
- Addition and subtraction
 - More complex than multiplication
 - Operands must have same exponent
 - Denormalize the smaller operand (alignment!)
 - Loss of digits (less precise and missing information)
 - Result (must) be normalised
- Multiplication and division
 - Significand and exponent handled separately

Computer Organization II, Spring 2010, Tiina Niklander

29.1.2010 36

Floating point arithmetics

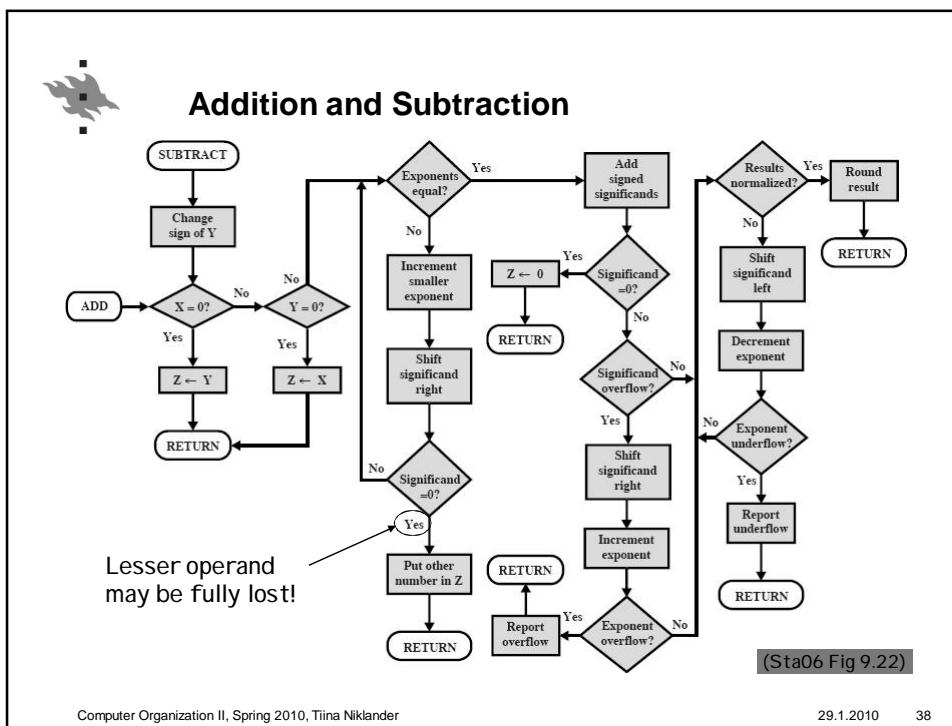
Floating Point Numbers	Arithmetic Operations
$X = X_S \times B^{X_E}$ $Y = Y_S \times B^{Y_E}$	$X + Y = \left(X_S \times B^{X_E - Y_E} + Y_S \right) \times B^{Y_E}$ $X - Y = \left(X_S \times B^{X_E - Y_E} - Y_S \right) \times B^{Y_E} \quad X_E \leq Y_E$ $X \times Y = (X_S \times Y_S) \times B^{X_E + Y_E}$ $\frac{X}{Y} = \left(\frac{X_S}{Y_S} \right) \times B^{X_E - Y_E}$

$X = 0.3 \times 10^2 = 30$
 $Y = 0.2 \times 10^3 = 200$
(Sta06 Table 9.5)

$X + Y = (0.3 \times 10^{2-3} + 0.2) \times 10^3 = 0.23 \times 10^3 = 230$
 $X - Y = (0.3 \times 10^{2-3} - 0.2) \times 10^3 = (-0.17) \times 10^3 = -170$
 $X \times Y = (0.3 \times 0.2) \times 10^{2+3} = 0.06 \times 10^5 = 6000$
 $X \div Y = (0.3 \div 0.2) \times 10^{2-3} = 1.5 \times 10^{-1} = 0.15$

Computer Organization II, Spring 2010, Tiina Niklander

29.1.2010 37



Computer Organization II, Spring 2010, Tiina Niklander

29.1.2010 38



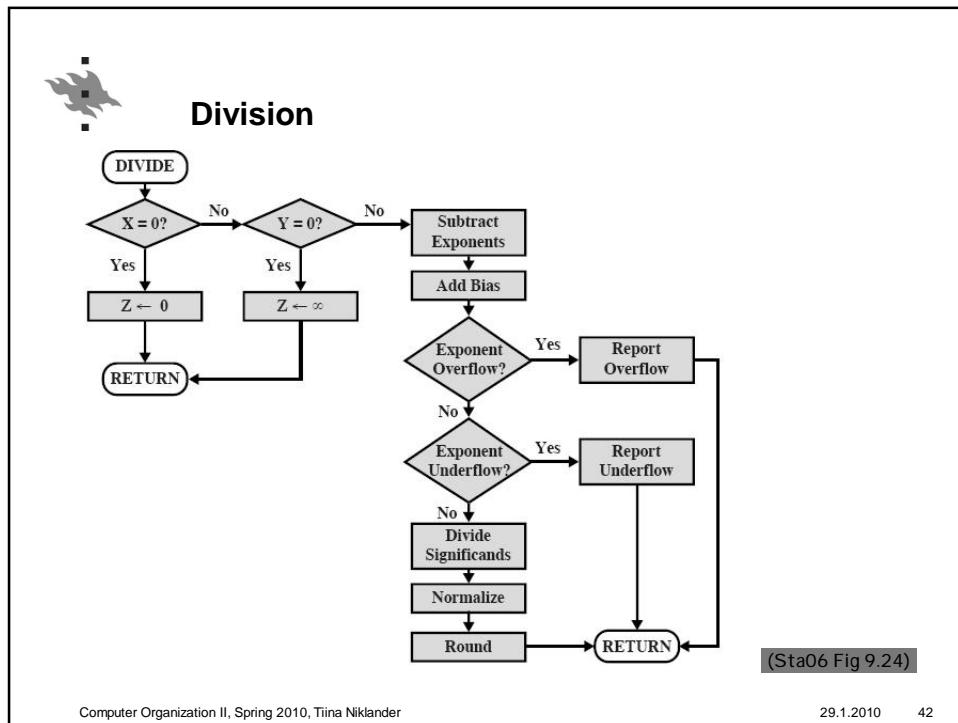
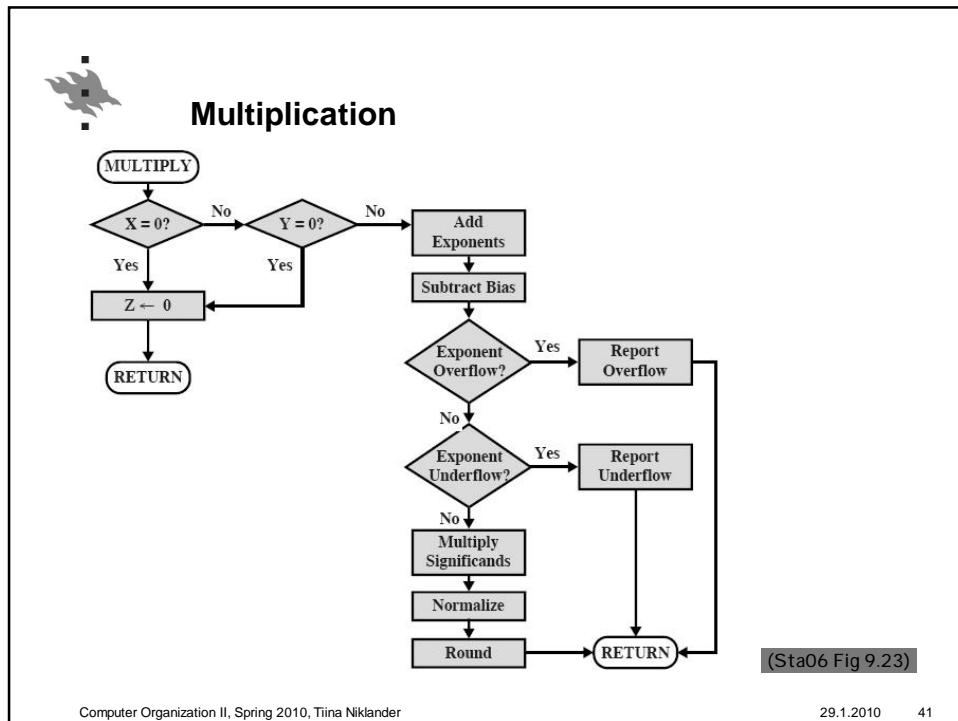
Special cases

- Exponent overflow (*eksponentin yliviuto*)
 - Very large number (above max) Programmable option
 - Value ∞ or $-\infty$, alternatively cause exception
- Exponent underflow (*eksponentin aliviuto*)
 - Very small number (below min) Programmable option
 - Value 0 (or cause exception)
- Significand overflow (*mantissan yliviuto*)
 - Normalise! Fix it!
- Significand underflow (*mantissan aliviuto*)
 - Denormalizing may lose the significand accuracy
 - All significant bits lost? Ooops, lost data!



Rounding (pyöristys)

- Example
 - Value has four decimals
 - Present it using only 3 decimals
- Normal rounding rule
 - round to nearest value
 - Always towards ∞ (*ylöspäin*)
 - Always towards $-\infty$ (*alaspäin*)
 - Always towards 0
- For example, Intel Itanium supports all of these alternatives
 - 3.1234, -4.5678
 - 3.123, -4.568
 - 3.124, -4.567
 - 3.123, -4.568
 - 3.123, -4.567





Review Questions / Kertauskysymyksiä

- Why we use twos complement?
- How does twos complement “expand” to a large number of bits (8b → 16 b)?
- Format of single-precision floating point number?
- When does underflow happen?

- Miksi käytetään 2:n komplementtimuotoa?
- Miten 2:n komplementtiesitys laajenee “suurempaan tilaan” (esim. 8b esitys → 16 b:n esitys)?
- Millainen on yksinkertaisen tarkkuuden liukuluvun esitysmuoto?
- Milloin tulee liukuluvun alivuoto?