

TE2003B

SoC Design: Computer organisation & architecture

Computer Arithmetic

Isaac Pérez Andrade

ITESM Guadalajara
School of Engineering & Science
Department of Computer Science
February – June 2021



References

The following material has been adopted and adapted from

Patterson, D. A., Hennessy, J. L., *Computer Organization and design: The hardware/software interface – ARM edition*, Morgan Kaufmann, 2017.

S. L. Harris and D. M. Harris, *Digital design and computer architecture - ARM edition*, Morgan Kaufmann, 2016.

J. Yiu, *The definitive guide to ARM Cortex-M0 and Cortex-M0+ processors*, Second edition, Elsevier, 2015.

Arithmetic for Computers

Operations on integers

- Addition and subtraction
- Multiplication and division
- Dealing with overflow

Floating-point real numbers

- Representation and operations

Integer operations

Two's complement review

Assume two's complement format

- Q: What's the range (minimum and maximum values that can be represented) of an N -bit two's complement number?

A: $[-(2^{(N-1)}), 2^{(N-1)} - 1]$

- For example, an 8-bit two's complement number may represent values in the range

$$[-2^{8-1}, 2^{8-1} - 1] = [-2^7, 2^7 - 1] = [-128, 127]$$

Overflow & underflow

- Q: What is **overflow**?

A: A condition when the result of a calculation **exceeds** the **maximum** value that can be represented in a numeric format.

- Q: What is **underflow**?

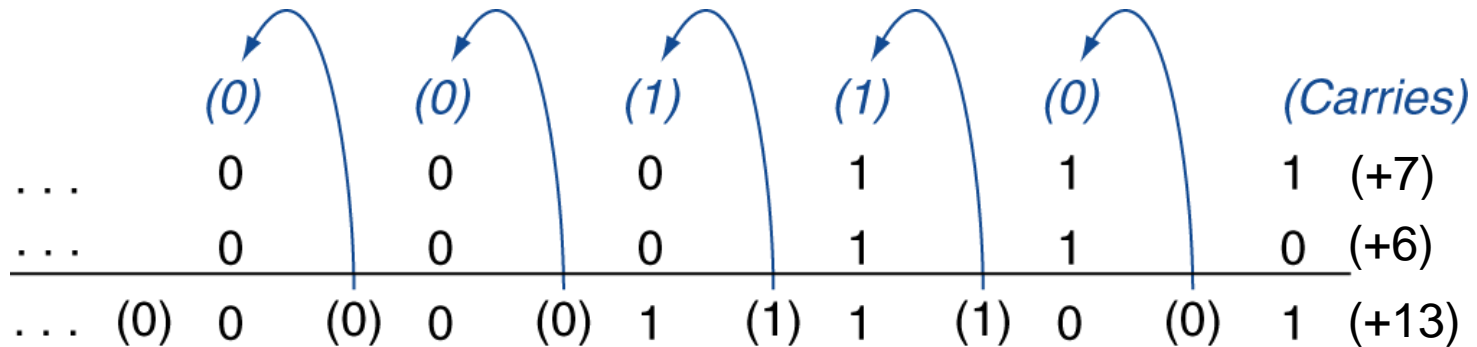
A: A condition when the result of a calculation is **smaller** than the **minimum** value that can be represented in a numeric format.

- Sometimes, the term overflow is used for describing both conditions.

Addition

Integer addition

Example: $7 + 6$



Overflow if result out of range

- Adding +ve and -ve operands, no overflow
- Adding two +ve operands
Overflow if result sign is 1
- Adding two -ve operands
Overflow if result sign is 0

Integer addition

- Example: Adding two 4-bit two's complement numbers

$$5 + 1$$

$$\begin{array}{r} +5: \quad 0101 \\ +1: \quad 0001 \\ \hline +6: \quad 0110 \end{array}$$

$$-2 + 5$$

$$\begin{array}{r} -2: \quad 1110 \\ +5: \quad 0101 \\ \hline +3: \quad 0011 \end{array}$$

$$3 + 6$$

$$\begin{array}{r} +3: \quad 0011 \\ +6: \quad 0110 \\ \hline -7: \quad 1001 \end{array}$$


$$-7 + (-1)$$

$$\begin{array}{r} -7: \quad 1001 \\ -1: \quad 1111 \\ \hline -8: \quad 1000 \end{array}$$

$$-3 + (-6)$$

$$\begin{array}{r} -3: \quad 1101 \\ -6: \quad 1010 \\ \hline +7: \quad 0111 \end{array}$$

Overflow: +9 and -9
can not be represented in 4-bit
two's complement.



Subtraction

Integer subtraction

Addition with negation of second operand

Example: $7 - 6 = 7 + (-6)$

$$\begin{array}{r} +7: \quad 0000 \ 0000 \ \dots \ 0000 \ 0111 \\ -6: \quad 1111 \ 1111 \ \dots \ 1111 \ 1010 \\ \hline +1: \quad 0000 \ 0000 \ \dots \ 0000 \ 0001 \end{array}$$

- Overflow if result out of range
 - Subtracting two +ve or two -ve operands, no overflow
 - Subtracting +ve from -ve operand
 - Overflow if result sign is 0
 - Subtracting -ve from +ve operand
 - Overflow if result sign is 1

Integer subtraction

- Example: Subtracting two 4-bit two's complement numbers

$$5 - (+1)$$

$$\begin{array}{r} +5: \quad 0101 \\ -1: \quad 1111 \\ \hline +4: \quad 0100 \end{array}$$

$$-2 - (-5)$$

$$\begin{array}{r} -2: \quad 1110 \\ +5: \quad 0101 \\ \hline +3: \quad 0011 \end{array}$$

$$-3 - (+6)$$

$$\begin{array}{r} -3: \quad 1101 \\ -6: \quad 1010 \\ \hline +7: \quad 0111 \end{array}$$

$$7 - (-1)$$

$$\begin{array}{r} +7: \quad 0111 \\ +1: \quad 0001 \\ \hline -8: \quad 1000 \end{array}$$

Overflow: -9 and +8 can not be represented in 4-bit two's complement.



Addition & subtraction overflow summary

- Overflow conditions for additions and subtraction in two's complement.

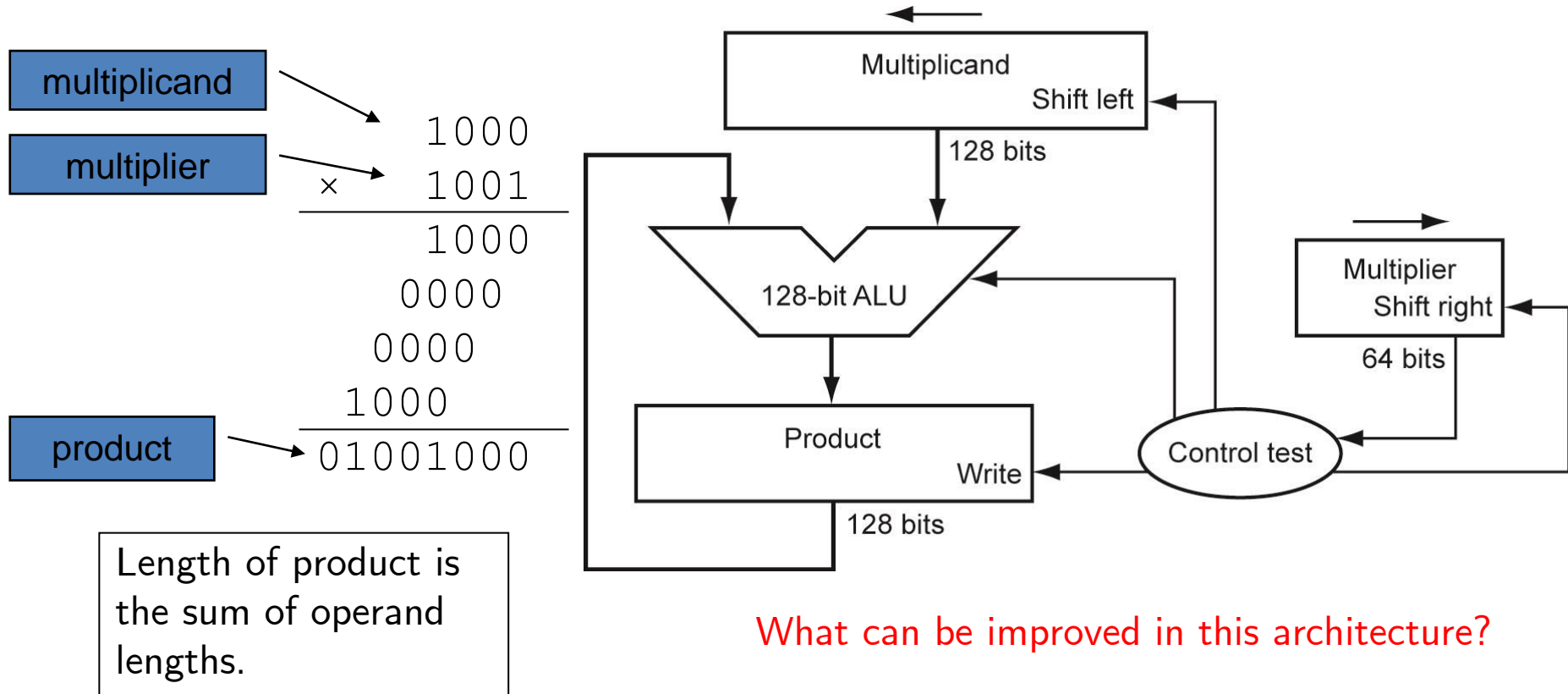
Operation	Operand A	Operand B	Result indicating overflow
$A + B$	≥ 0	≥ 0	< 0
$A + B$	< 0	< 0	≥ 0
$A - B$	≥ 0	< 0	< 0
$A - B$	< 0	≥ 0	≥ 0

Multiplication

Multiplication

- Start with long-multiplication approach

Assume we want to multiply two 64-bit numbers

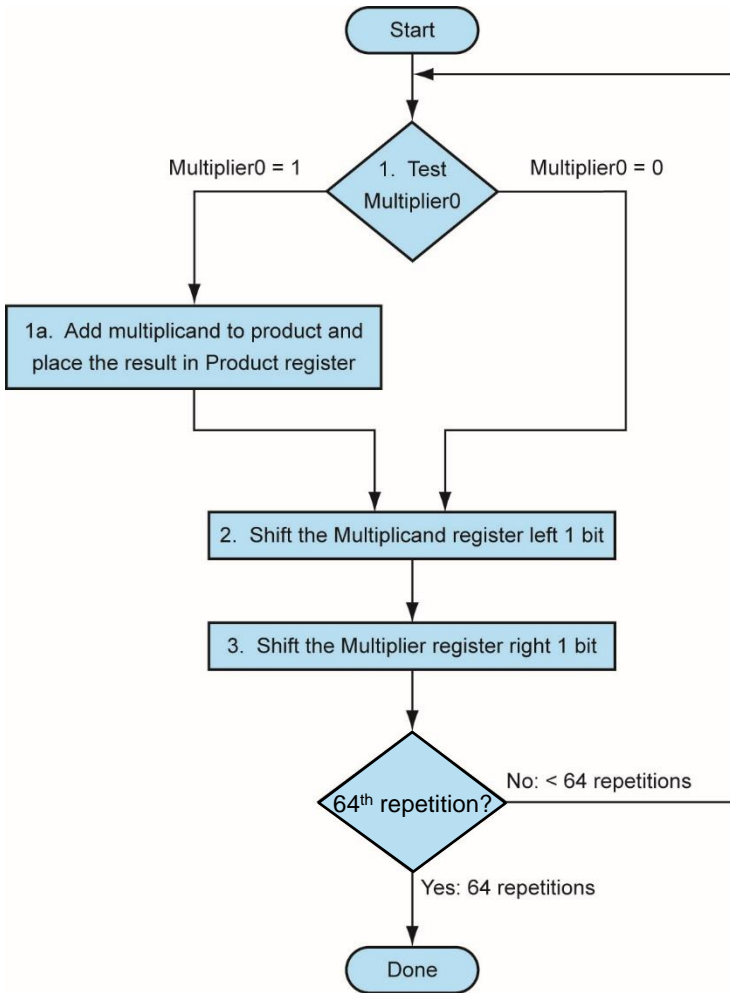


What can be improved in this architecture?

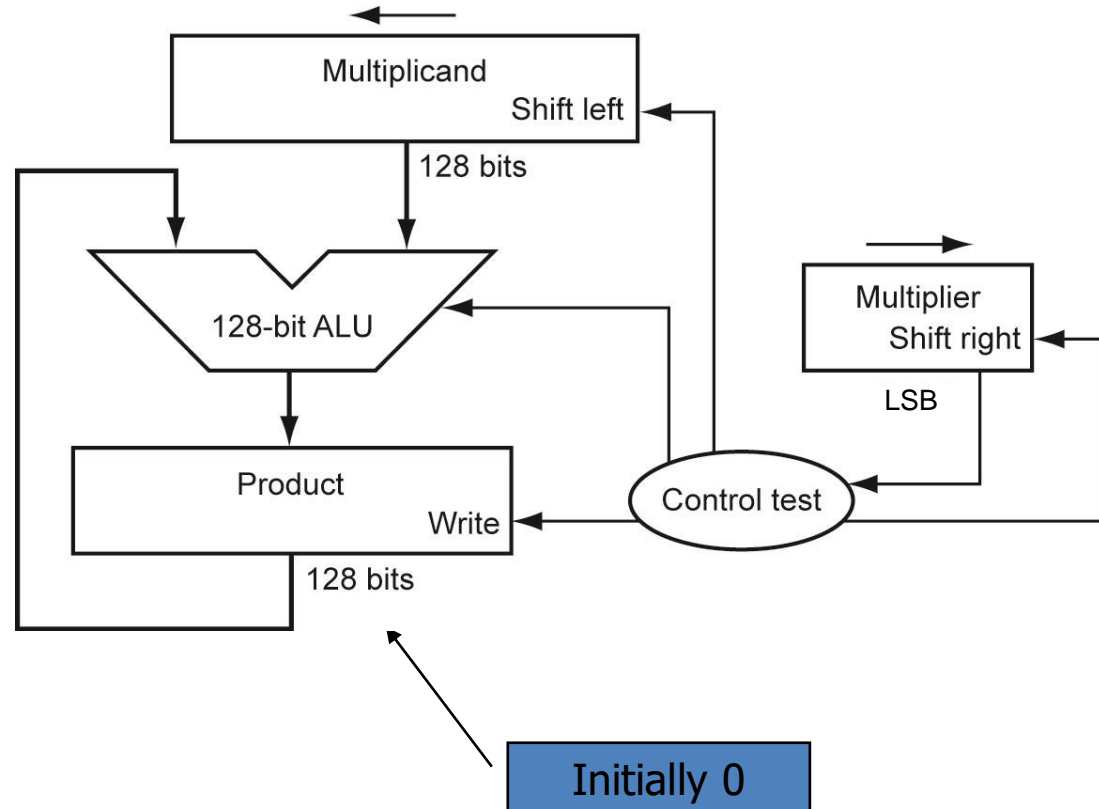
$$A_{M-bits} \times B_{N-bits} = X_{(M+N)-bits}$$

Multiplication hardware

There's one error in the flow chart. Can you spot it?

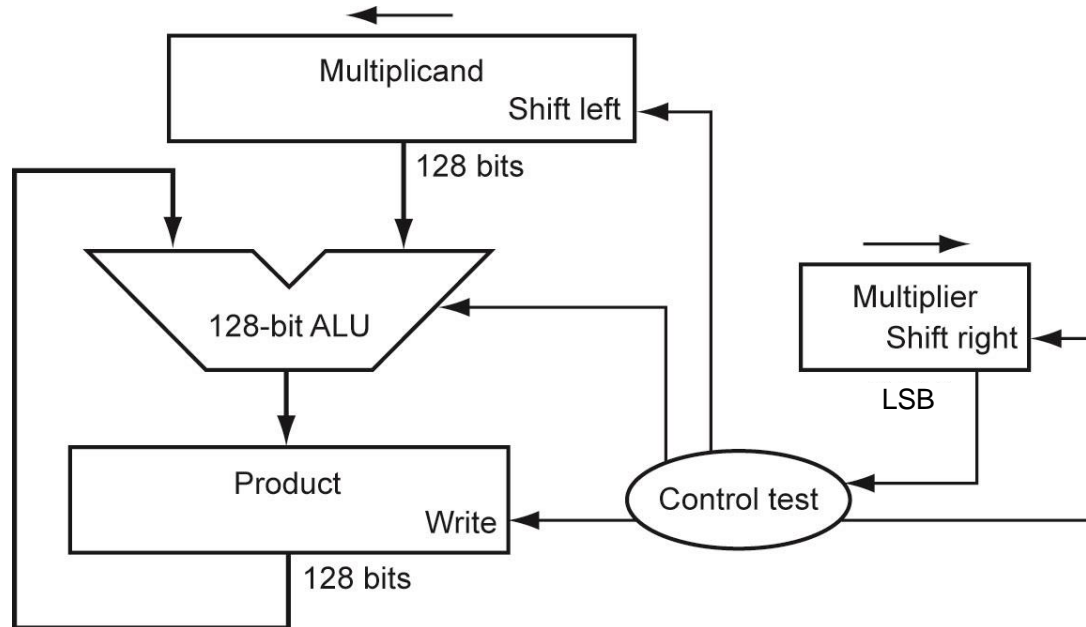


What can be improved in this architecture?



Multiplication hardware

This architecture has a major flaw.
Can you spot it?



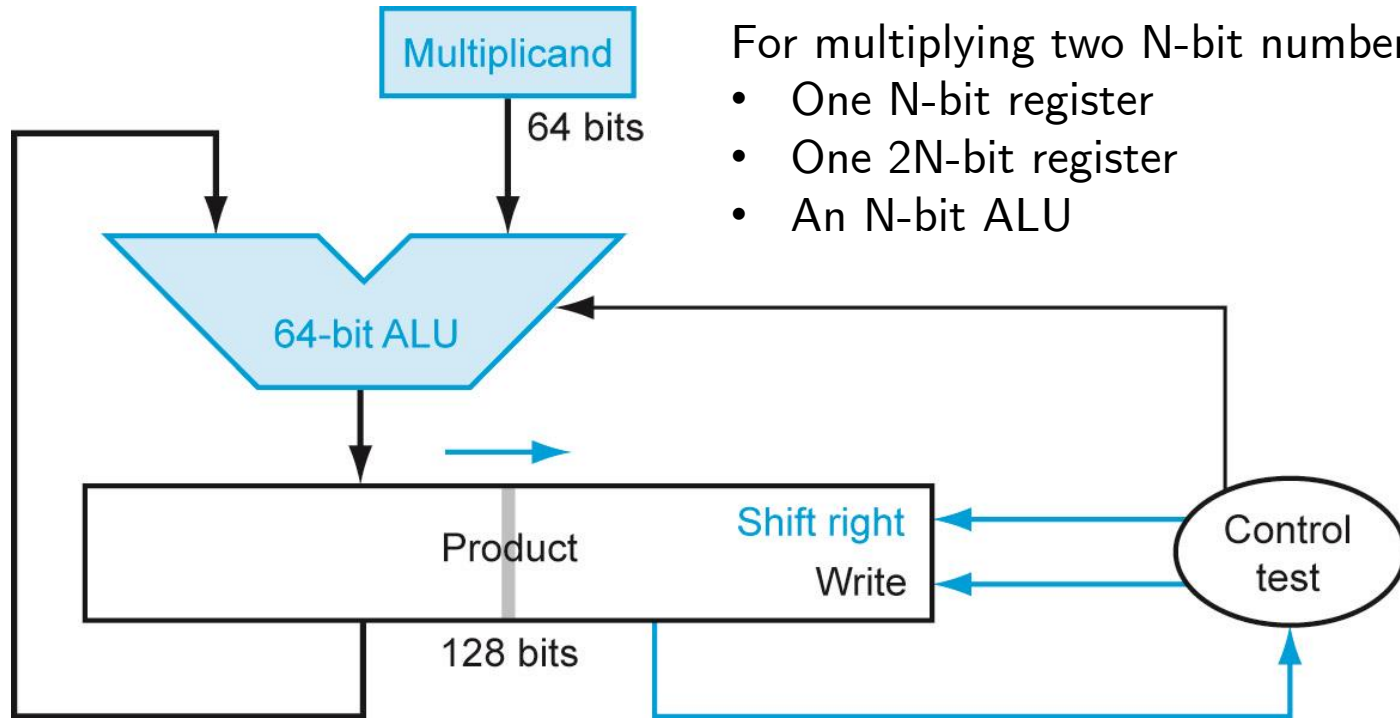
For multiplying two N -bit number, it requires:

- Two $2N$ -bit registers
- One N -bit register
- A $2N$ -bit ALU

This is a waste of resources!

Optimised multiplier

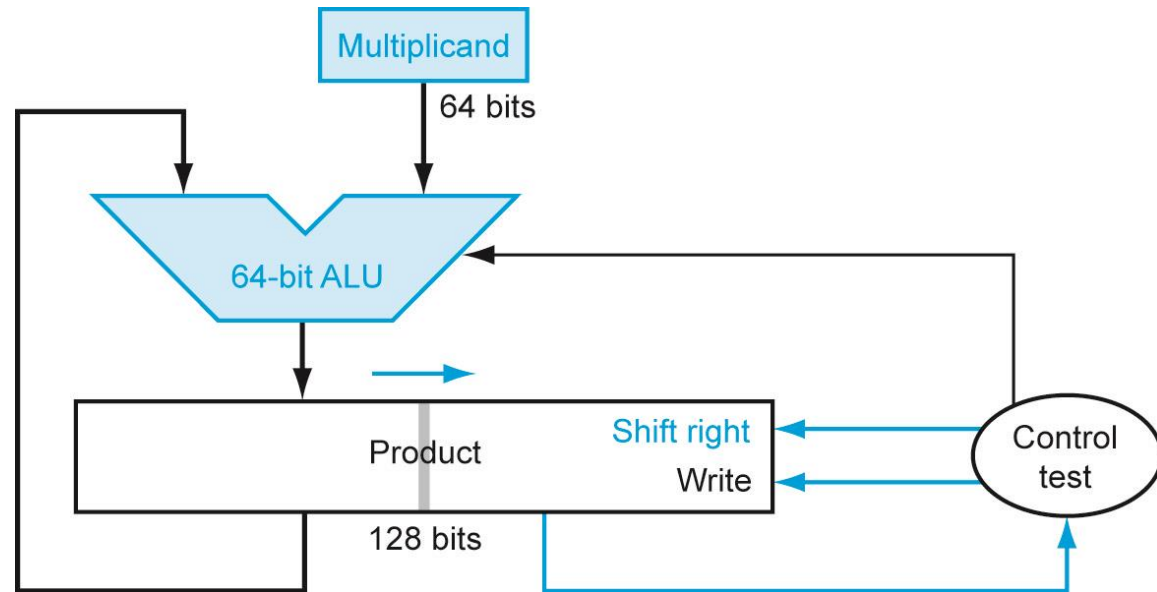
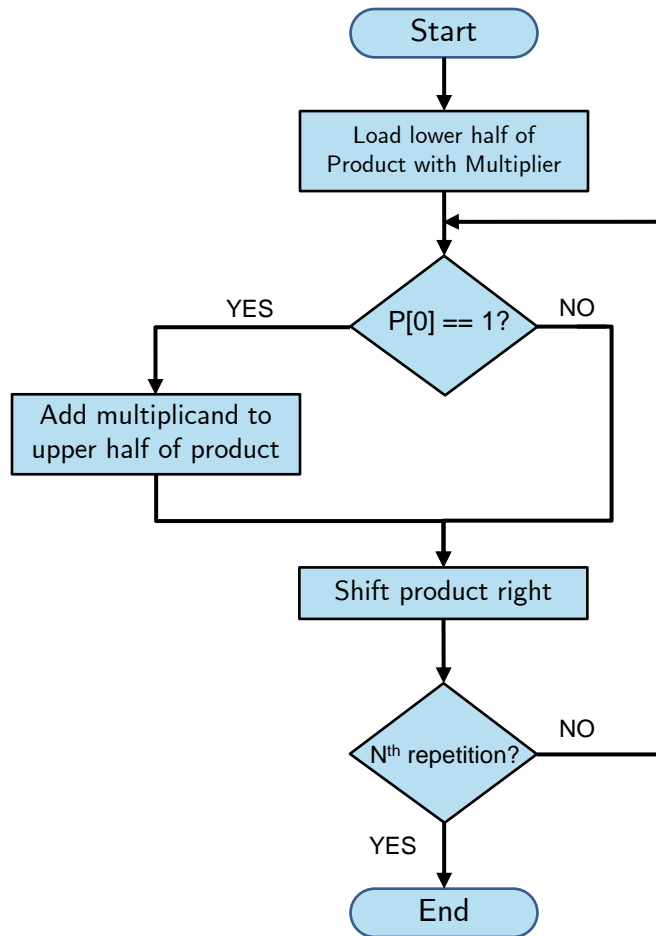
- Perform steps in parallel: add/shift



One cycle per partial-product addition

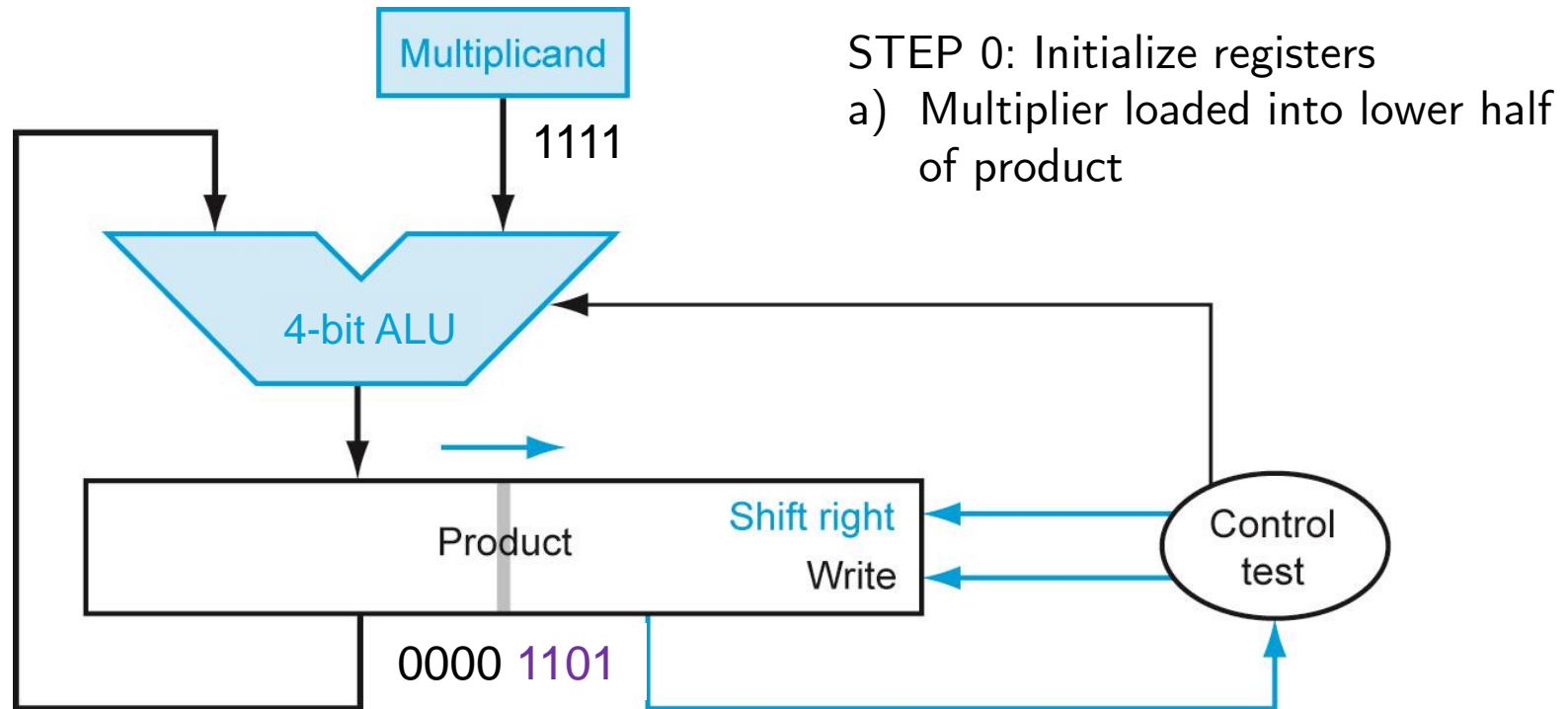
- That's ok, if frequency of multiplications is low

Optimised multiplier



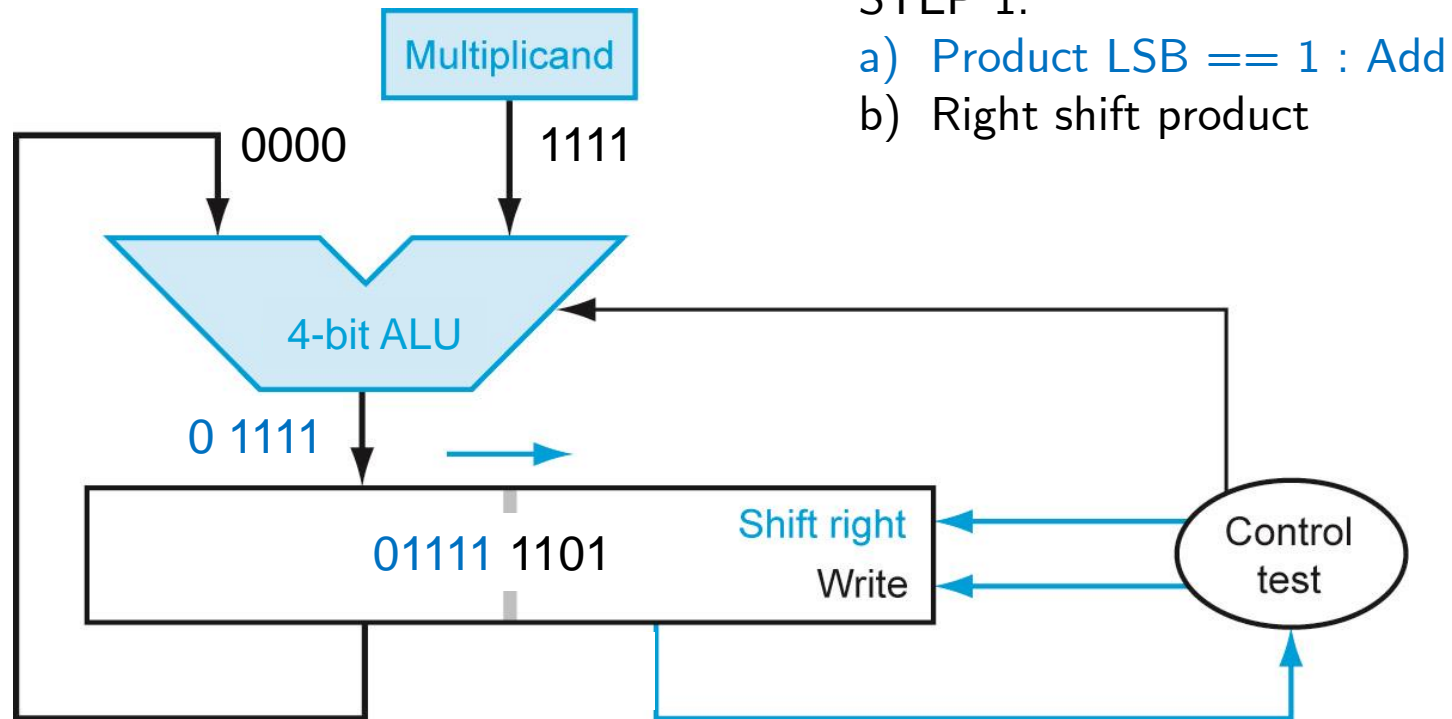
Optimised multiplier - example

- Multiplicand = 15_{10} : 1111
- Multiplier = 13_{10} : 1101



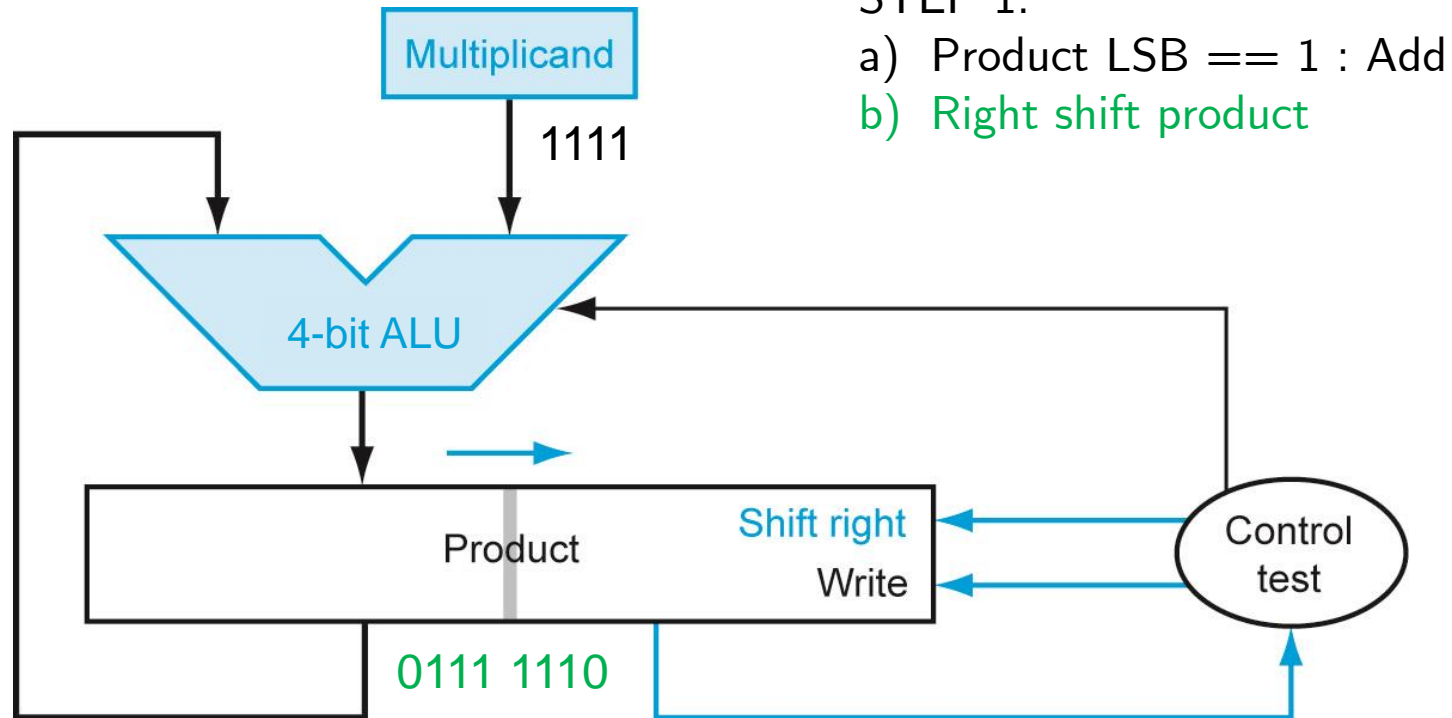
Optimised multiplier - example

- Multiplicand = 15_{10} : 1111
- Multiplier = 13_{10} : 1101



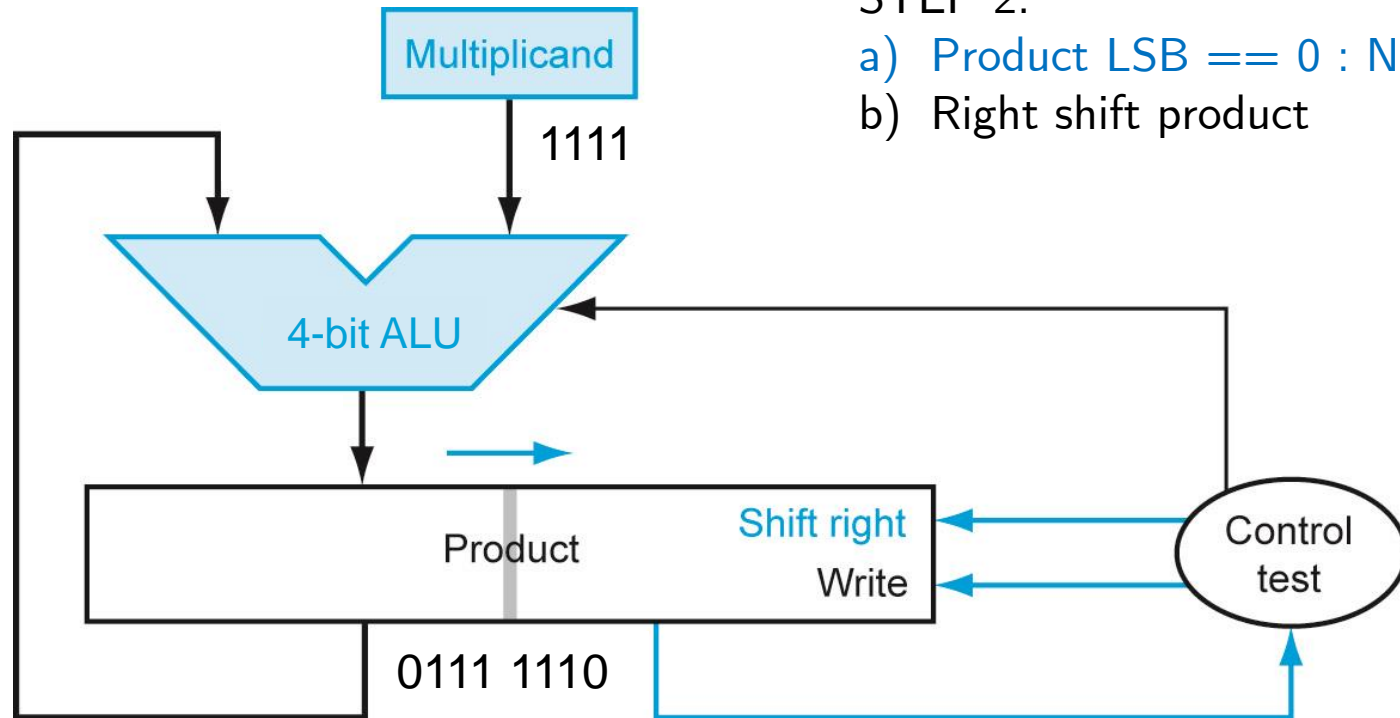
Optimised multiplier - example

- Multiplicand = 15_{10} : 1111
- Multiplier = 13_{10} : 1101



Optimised multiplier - example

- Multiplicand = 15_{10} : 1111
- Multiplier = 13_{10} : 1101

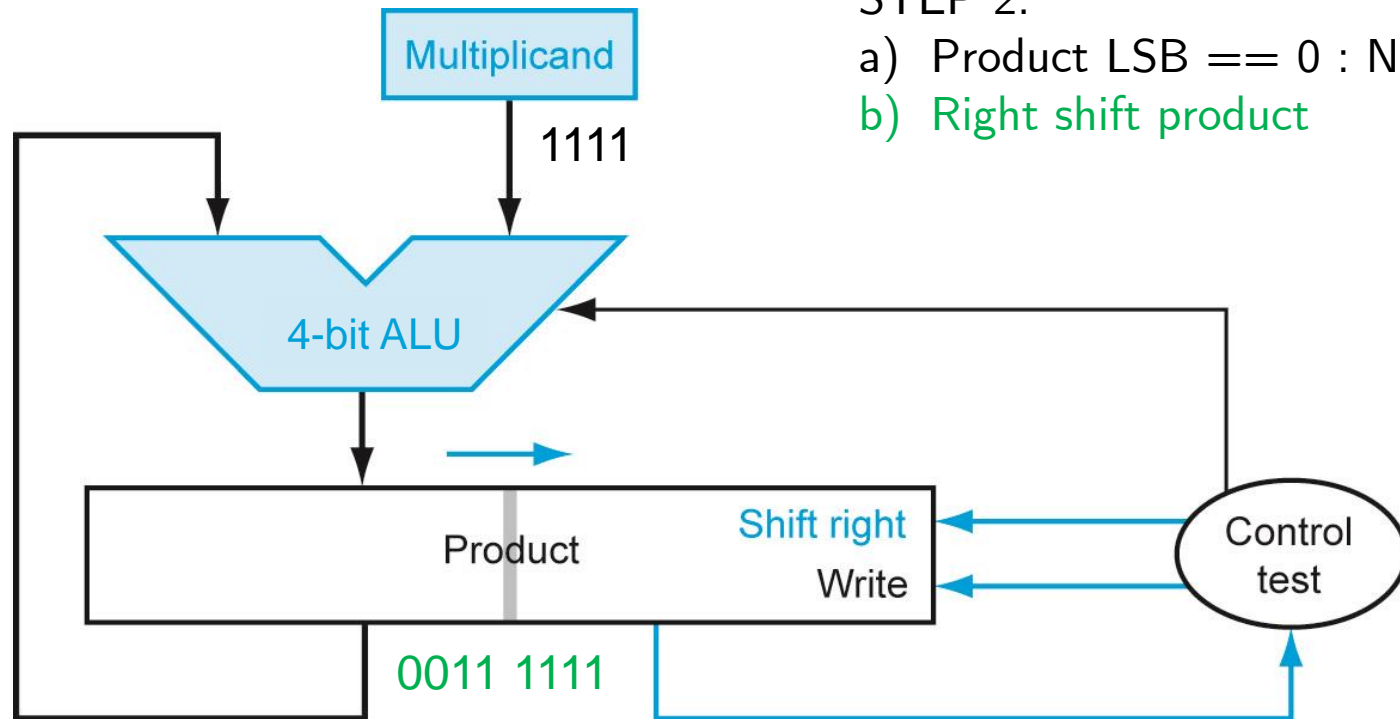


STEP 2:

- a) Product LSB == 0 : No Add
- b) Right shift product

Optimised multiplier - example

- Multiplicand = 15_{10} : 1111
- Multiplier = 13_{10} : 1101



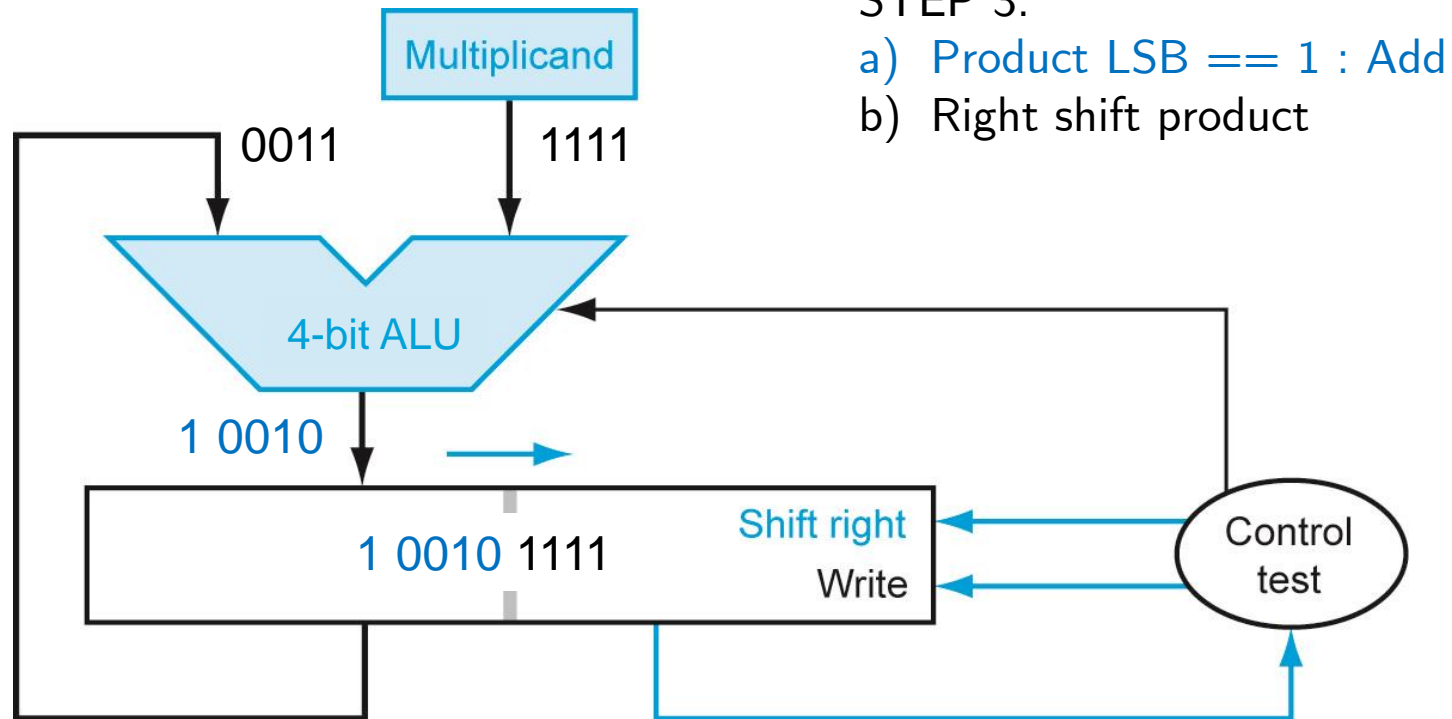
STEP 2:

a) Product LSB == 0 : No Add

b) Right shift product

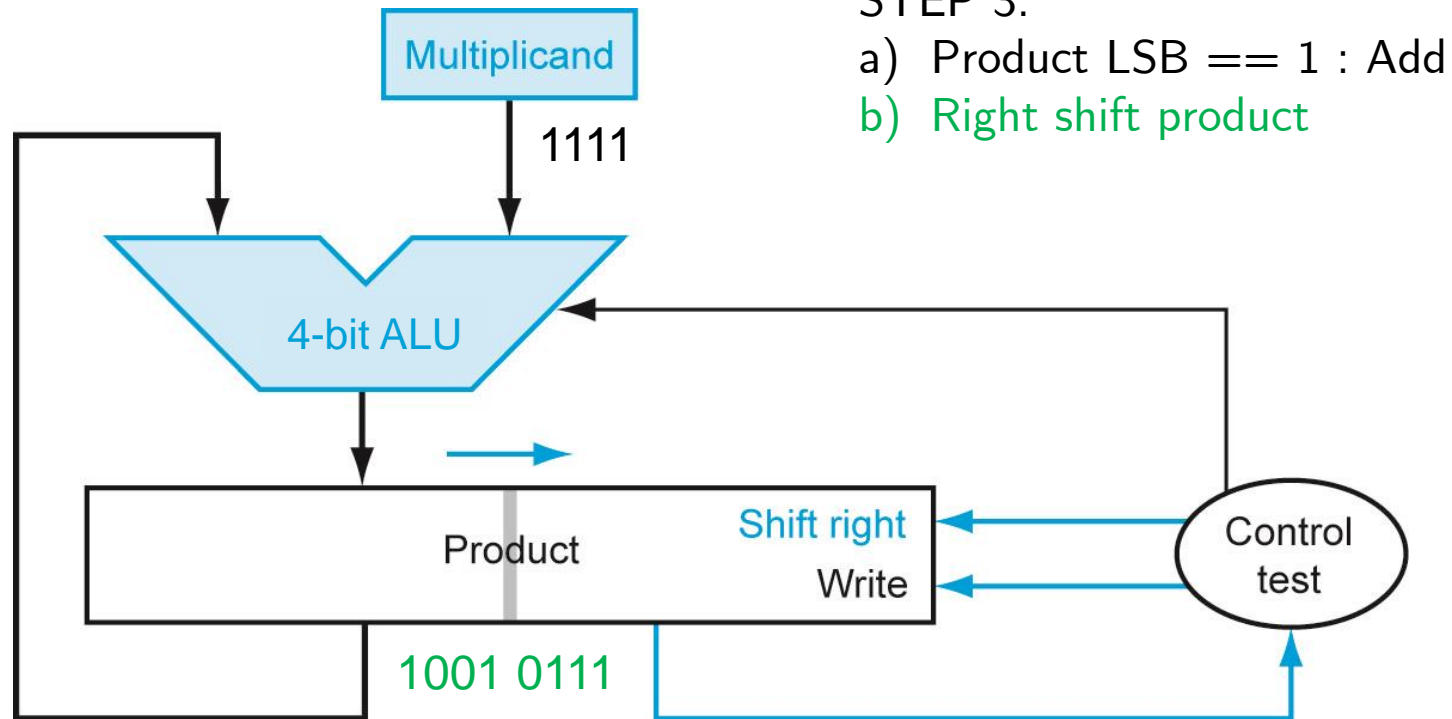
Optimised multiplier - example

- Multiplicand = 15_{10} : 1111
- Multiplier = 13_{10} : 1101



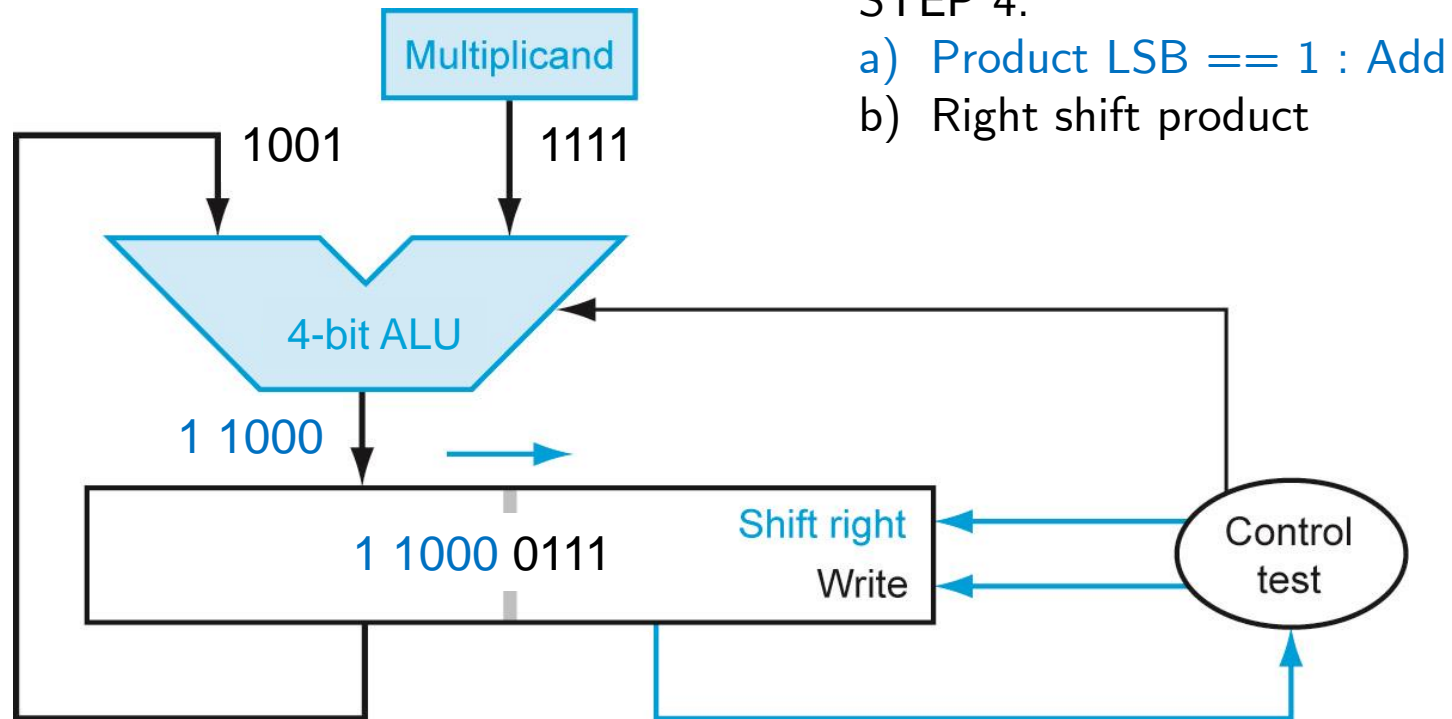
Optimised multiplier - example

- Multiplicand = 15_{10} : 1111
- Multiplier = 13_{10} : 1101



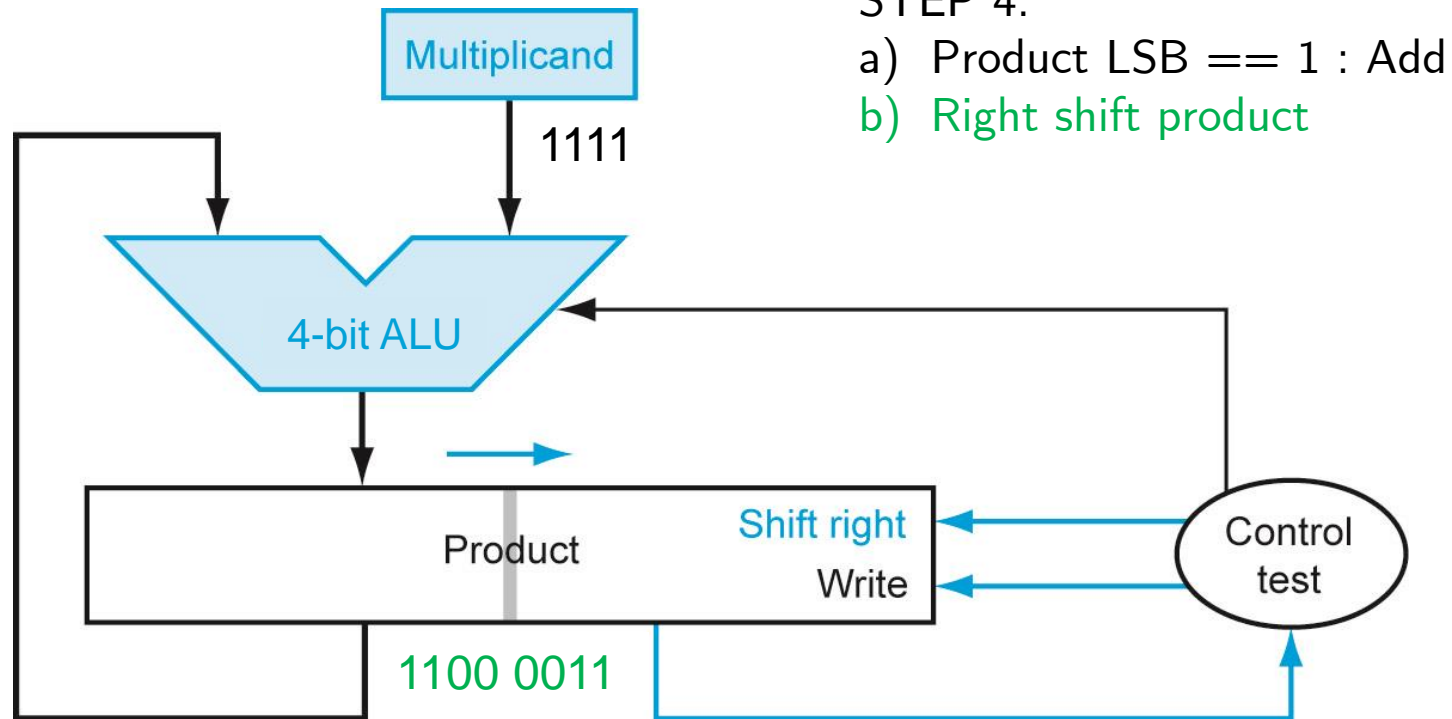
Optimised multiplier - example

- Multiplicand = 15_{10} : 1111
- Multiplier = 13_{10} : 1101



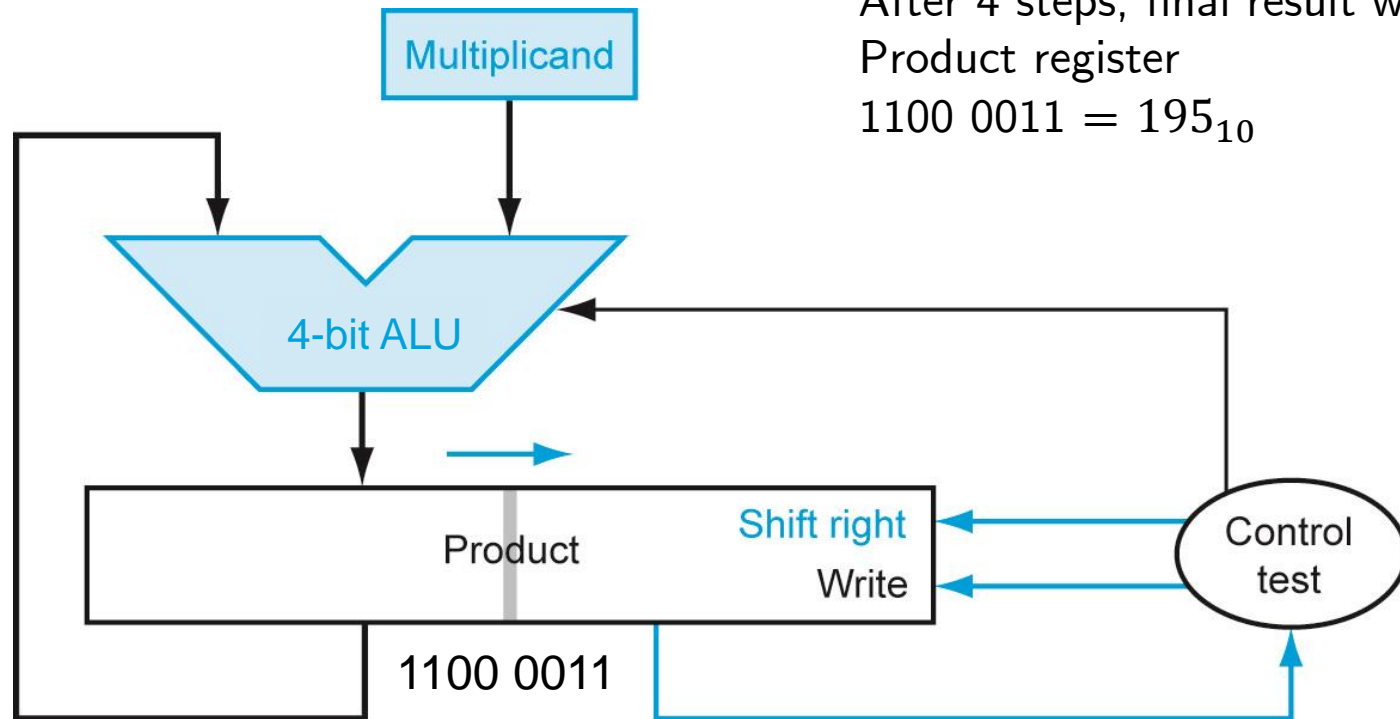
Optimised multiplier - example

- Multiplicand = 15_{10} : 1111
- Multiplier = 13_{10} : 1101



Optimised multiplier - example

- Multiplicand = 15_{10} : 1111
- Multiplier = 13_{10} : 1101

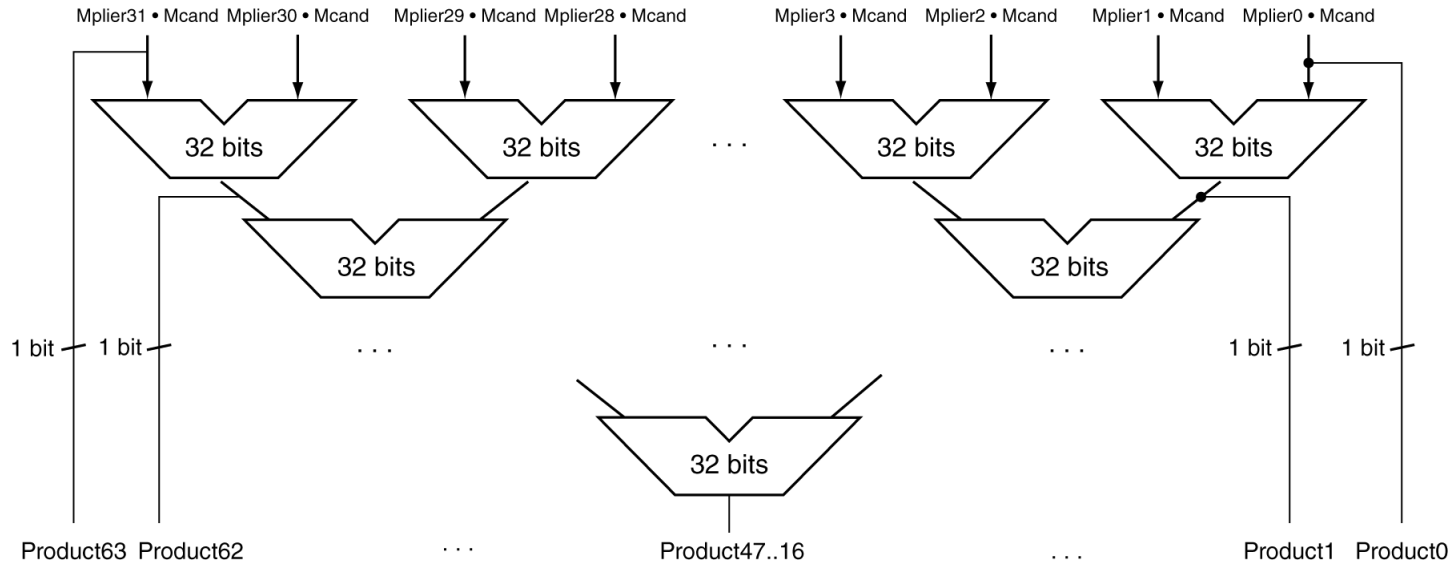


Signed multiplication

- So far, we've only dealt with unsigned operands.
- What happens in signed multiplication?
- For adding two signed N-bit numbers:
 1. Convert both multiplicand and multiplier to positive numbers and keep track of their respective sign.
 2. Apply multiplication algorithm N-1 times.
 3. Negate product if signs are not the same.
- Alternatively, previous algorithm works for signed numbers as long as shifts are performed using sign extension.

Faster Multiplier

- Uses multiple adders
- Cost/performance tradeoff



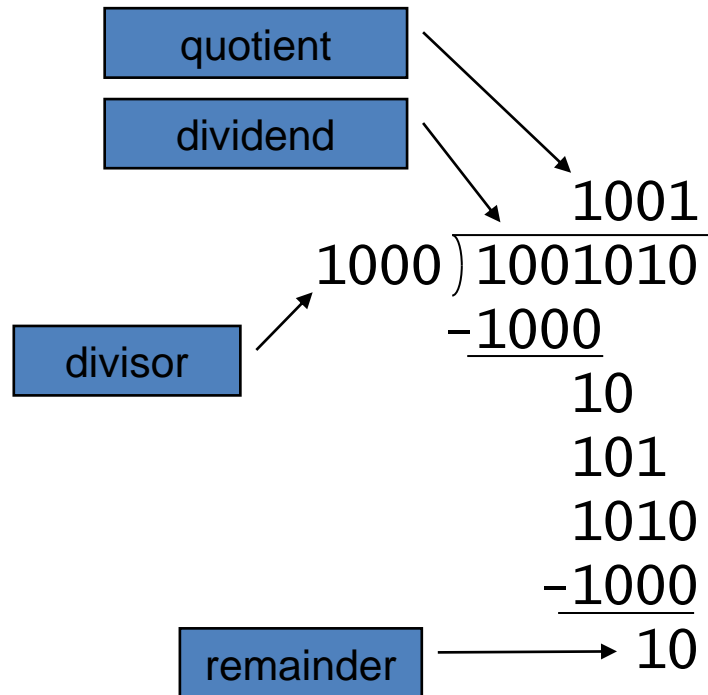
Use $(N-1)$ N -bit adders in parallel, instead of a single N -bit adder $(N-1)$ times.

Can be pipelined

- Several multiplication performed in parallel

Division

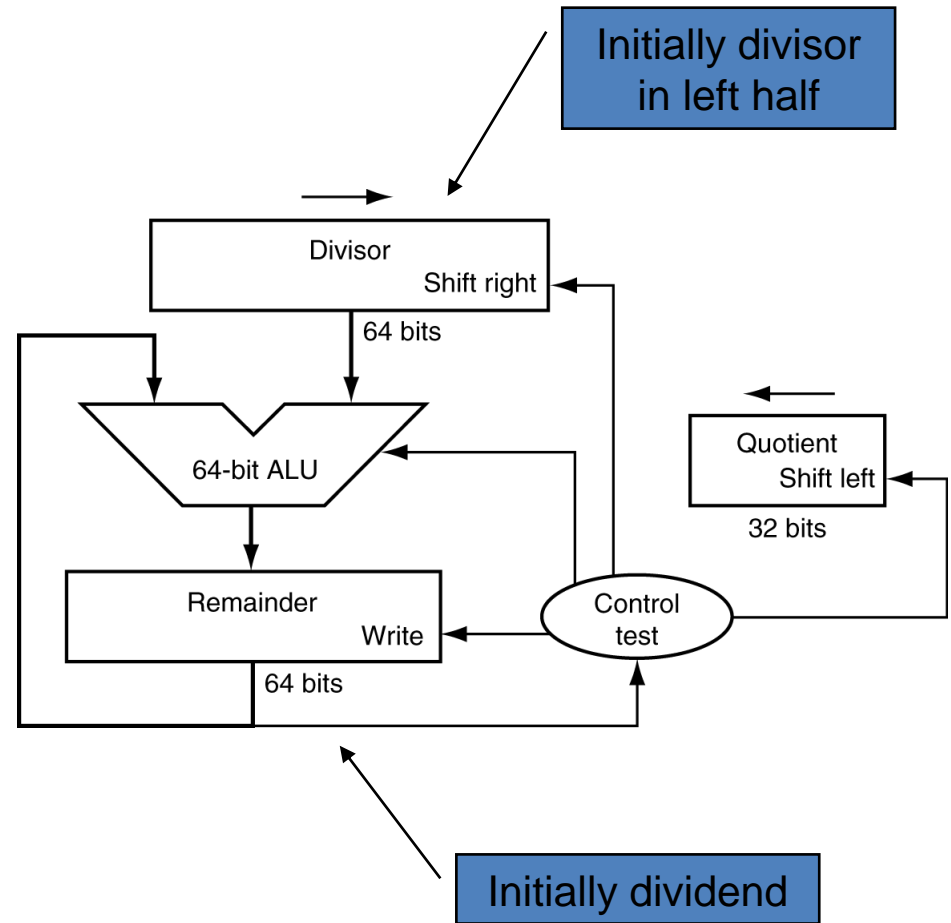
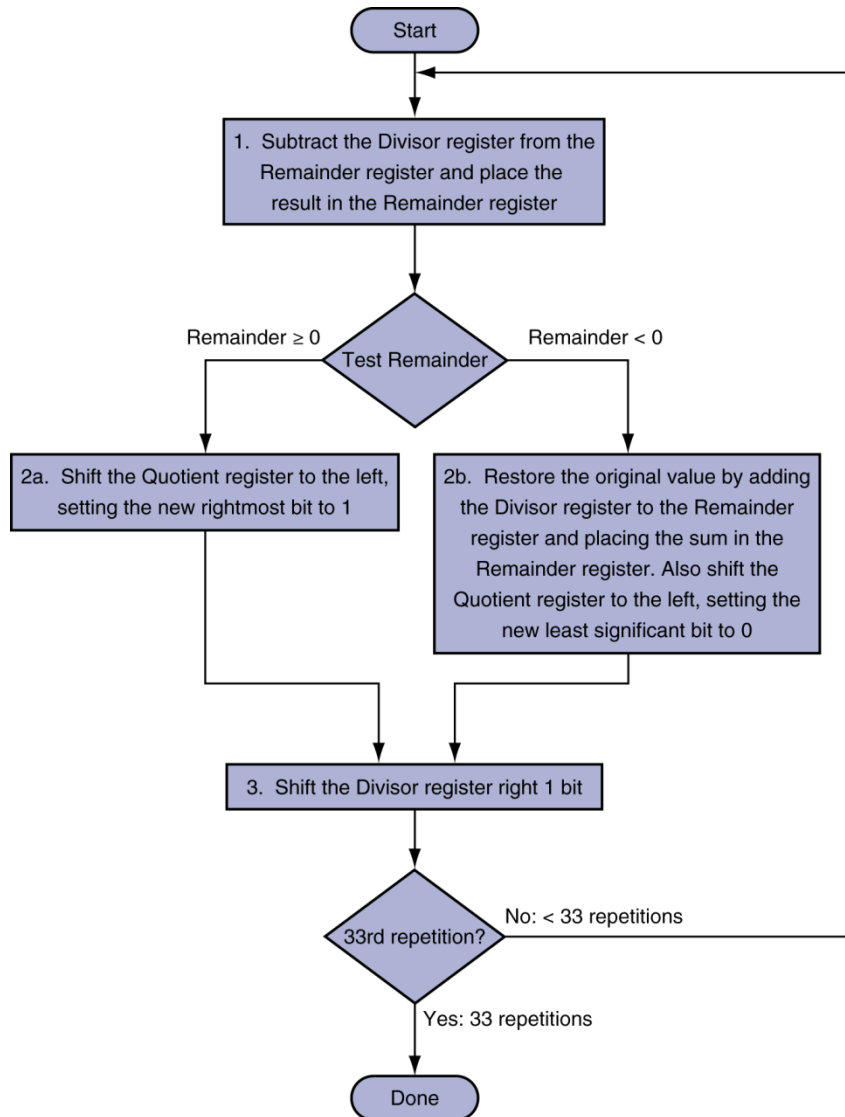
Division



n-bit operands yield *n*-bit quotient and remainder

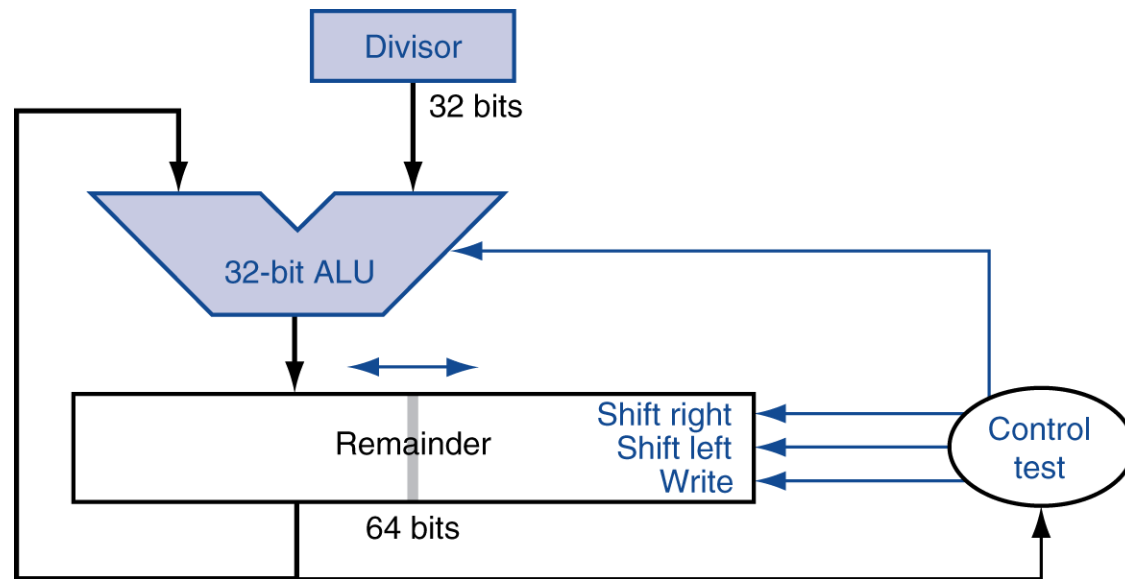
- Check for 0 divisor
- Long division approach
 - If divisor \leq dividend bits
 - 1 bit in quotient, subtract
 - Otherwise
 - 0 bit in quotient, bring down next dividend bit
- Restoring division
 - Do the subtract, and if remainder goes < 0 , add divisor back
- Signed division
 - Divide using absolute values
 - Adjust sign of quotient and remainder as required

Division hardware



Optimized divider

- One cycle per partial-remainder subtraction
- Looks a lot like a multiplier!
 - Same hardware can be used for both



Faster Division

- Can't use parallel hardware as in multiplier
 - Subtraction is conditional on sign of remainder
- Faster dividers generate multiple quotient bits per step
 - Still require multiple steps

Fixed-point representation

Fixed-point introduction

- Real (fractional) numbers may not be represented with integer numbers.

```
integer a,b;
```

```
a = 1.5;
```

```
b = a + a; // b = ?
```

- Fixed-point representation allows real number representation with limited precision.
 - Q_{m.n} representation
 - $m \rightarrow$ number of bits for representing integer part.
 - $n \rightarrow$ number of bits for representing fractional part.
 - Range $[-(2^{m-1}), 2^{m-1} - 2^{-n}]$
 - Resolution is 2^{-n}

Fixed-point

- Fixed-point representation is suitable for embedded applications requiring limited degree of fractional precision.
 - DOOM (1993 videogame) originally used a Q16.16 format for all non-integer operations
https://doomwiki.org/wiki/Fixed_point
- What about high-precision applications?

Fixed-point limitations

- Example:

- Consider Avogadro's number: 6.022×10^{23}

- How many bits would you need to represent Avogadro's number?

$$\lceil \log_2(6.022 \times 10^{23}) \rceil = 79$$

- What about a very small number such as Planck's constant: $6.62607004 \times 10^{-34} \text{ J} \cdot \text{s}$

- How many bits (fractional fixed-point) would you need to represent Planck's constant?

$$\lceil \log_2(6.62607004 \times 10^{-34}) \rceil = 110$$

- We would need at least $79 + 110 = 189$ bits for representing both numbers.

- Not feasible, waste of resources.
- What if need even smaller or larger numbers?

Floating-point representation

Floating-point

- Scientific notation

- A single digit to the left of the decimal point

- $+1.12345 \times 10^{-7}$

- -123.456×10^9

- **Sign**
- **Mantissa/significant**
- **Exponent**

normalized

not normalized

- Normalized scientific notation

- Absolute value of integer part m is in the range
 $[1, 10) \rightarrow 1 \leq m \leq 9$

- Binary numbers may also be represented in scientific notation

$$1.0_2 \times 2^{-1} = 0.1_2 \quad \leftarrow 0.5_{10}$$

$$0.1_2 \times 2^0 = 0.1_2$$

normalized

not normalized

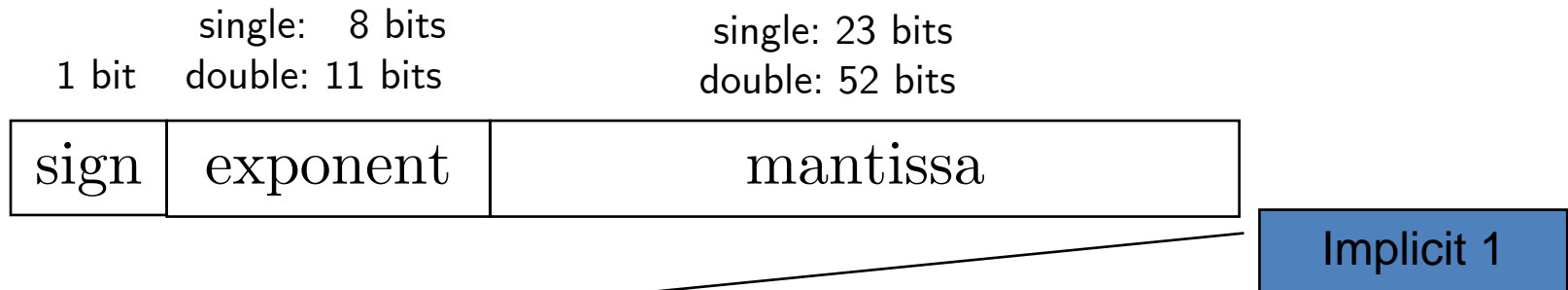
Floating-point

- As the name suggest, binary point is not fixed.
- Representation for non-integral numbers
 - Including very small and very large numbers
 - $(-1)^{\text{sign}} 1.\text{mantissa} \times 2^{(\text{exponent}-\text{bias})}$
 - For simplicity, we'll show the exponent in decimal.
- Programming languages refer to this representation as **float** and **double** types.

Floating-point standard

- Defined by IEEE Std 754-1985
- Developed in response to divergence of representations
 - Portability issues for scientific code
- Now almost universally adopted
- Two representations
 - Single precision (32-bit)
 - Double precision (64-bit)
- Simplifies exchange of data, arithmetic and increases accuracy.

IEEE Floating-point format



$$x = (-1)^{\text{sign}} 1.\text{mantissa} \times 2^{(\text{exponent} - \text{bias})}$$

- **sign**: (0 \Rightarrow non-negative, 1 \Rightarrow negative)
- Significand: **1.mantissa**
- Normalized significand: $1.0 \leq |\text{significand}| < 2.0$
- **exponent** = **actual exponent** + bias
 - Ensures exponent is unsigned
 - Single: bias = 127
 - Double: bias = 1203

IEEE Floating-point format

1 bit	single: 8 bits double: 11 bits	single: 23 bits double: 52 bits
sign	exponent	mantissa

$$x = (-1)^{\text{sign}} 1.\text{mantissa} \times 2^{(\text{exponent} - \text{bias})}$$

- First example: 0.5 to single-precision floating-point
 $+0.5_{10} = +0.1_2 = +1.0_2 \times 2^{-1}$

For simplification, we are using decimal notation for the exponent

- From this, we can gather the following information:
 - Sign: 0
 - Actual exponent: -1
 - Mantissa: 0 (normalized)
- Everything together:
 - Adjusted exponent = $-1 + 127 = 126_{10}$ or 01111110_2
 - Normalized floating-point:
 $(-1)^0 1.0 \times 2^{(126 - 127)}$

0 01111110 000000000000000000000000

IEEE Floating-point format

1 bit	single: 8 bits double: 11 bits	single: 23 bits double: 52 bits
sign	exponent	mantissa

$$x = (-1)^{\text{sign}} 1.\text{mantissa} \times 2^{(\text{exponent} - \text{bias})}$$

- First example: 0.5 to single-precision floating-point

$$+0.5_{10} = (-1)^{\overline{0}} \underline{\overline{1.0}} \times 2^{(126-127)}$$

0 01111110 000000000000000000000000

This 1 is not
actually required
here

- What about 0.5 in double-precision?

$$0.5_{10} = (-1)^{\text{green}} \text{red}1.\text{purple}0 \times 2^{(\text{blue}1022-1023)}$$

```
0 0111111110 000000000000000000000000000000000000000000  
000000
```

Floating-point example

- Represent -0.75
 - $-0.75 = (-1)^1 \times 1.1_2 \times 2^{-1}$
 - $\text{sign} = 1$
 - $\text{mantissa} = 1000\dots00_2$
 - $\text{exponent} = -1 + \text{bias}$
 - Single: $-1 + 127 = 126 = 01111110_2$
 - Double: $-1 + 1023 = 1022 = 01111111110_2$
- Single: $10111111101000\dots00$
- Double: $101111111111101000\dots00$

Floating-point example

single: 8 bits double: 11 bits		single: 23 bits double: 52 bits
1 bit		
sign	exponent	mantissa

- What number is represented by the single-precision float

11000000101000000000000000000000

- sign = 1
- mantissa = 01000...00₂
- exponent = 10000001₂ = 129
- $x = (-1)^1 \times (1.01_2) \times 2^{(129 - 127)}$
 $= (-1) \times 1.25 \times 2^2$
 $= -5.0$

Remember that this 1 is always implied!

Single precision range

Exponents 00000000 and 11111111 reserved

Smallest value

- exponent: 00000001
 \Rightarrow actual exponent = $1 - 127 = -126$
- Fraction: 000...00 \Rightarrow significand = 1.0
- $\pm 1.0 \times 2^{-126} \approx \pm 1.2 \times 10^{-38}$

Largest value

- exponent: 11111110
 \Rightarrow actual exponent = $254 - 127 = +127$
- Fraction: 111...11 \Rightarrow significand ≈ 2.0
- $\pm 2.0 \times 2^{+127} \approx \pm 3.4 \times 10^{+38}$

Double precision range

Exponents 0000000000 and 1111111111 reserved

Smallest value

- Exponent: 0000000001
 \Rightarrow actual exponent = $1 - 1023 = -1022$
- Fraction: 000...00 \Rightarrow significand = 1.0
- $\pm 1.0 \times 2^{-1022} \approx \pm 2.2 \times 10^{-308}$

Largest value

- exponent: 1111111110
 \Rightarrow actual exponent = $2046 - 1023 = +1023$
- Fraction: 111...11 \Rightarrow significand ≈ 2.0
- $\pm 2.0 \times 2^{+1023} \approx \pm 1.8 \times 10^{+308}$

Floating-point precision

- Relative precision
 - Single: approx. $2^{-23} \rightarrow \approx 7$ decimal digits
 - Double: approx. $2^{-52} \rightarrow \approx 16$ decimal digits

Floating-point special representation

- Denormal numbers

- In normalised numbers, significand have an implicit leading 1

$$x = (-1)^{\text{sign}} 1.\text{mantissa} \times 2^{(\text{exponent}-\text{bias})}$$

- Denormal numbers have a leading 0 in the significand.
- Biased exponent is 0.
- These numbers allow to represent numbers smaller than the smaller normalised number, as well as special representation such as $\pm\infty$ and NaN ($0 \div 0$).

Denormal numbers

$$x = (-1)^{\text{sign}} \times (0.\text{mantissa}) \times 2^{0-\text{bias}}$$

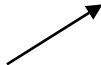
- Smaller than normal numbers.

- Zero

- $\text{sign} = 0, 1$
- $\text{biased exponent} = 0$
- $\text{mantissa} = 0$

$$x = (-1)^{\text{sign}} \times (0 + 0) \times 2^{-\text{bias}} = \pm 0.0$$

Two representations
of 0.0!



Denormalized numbers

Smallest
denormalized
value

- Single precision:
 $\pm 2^{-23} \times 2^{-126} \approx 1.4 \times 10^{-45}$
- Double precision:
 $\pm 2^{-52} \times 2^{-1022} \approx 4.9 \times 10^{-324}$

Largest
denormalized
value

- Single precision:
 $\pm(1 - 2^{-23}) \times 2^{-126} \approx \pm 1.17 \times 10^{-38}$
- Double precision:
 $\pm(1 - 2^{-52}) \times 2^{-1022} \approx \pm 2.22 \times 10^{-308}$

Infinites and NaNs

- Exponent = 111...1, mantissa = 000...0
 - $\pm\infty$
 - Can be used in subsequent calculations, avoiding need for overflow check
- Exponent = 111...1, mantissa \neq 000...0
 - Not-a-Number (NaN)
 - Indicates illegal or undefined result
 - e.g., $0.0 / 0.0$
 - Can be used in subsequent calculations

Floating-point special formats summary

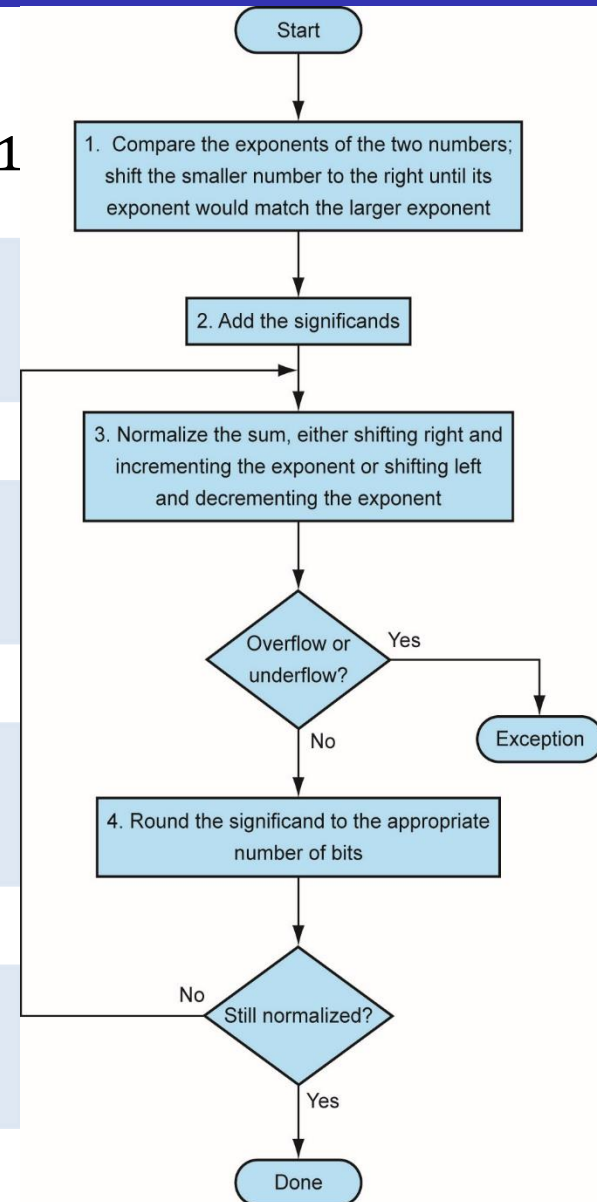
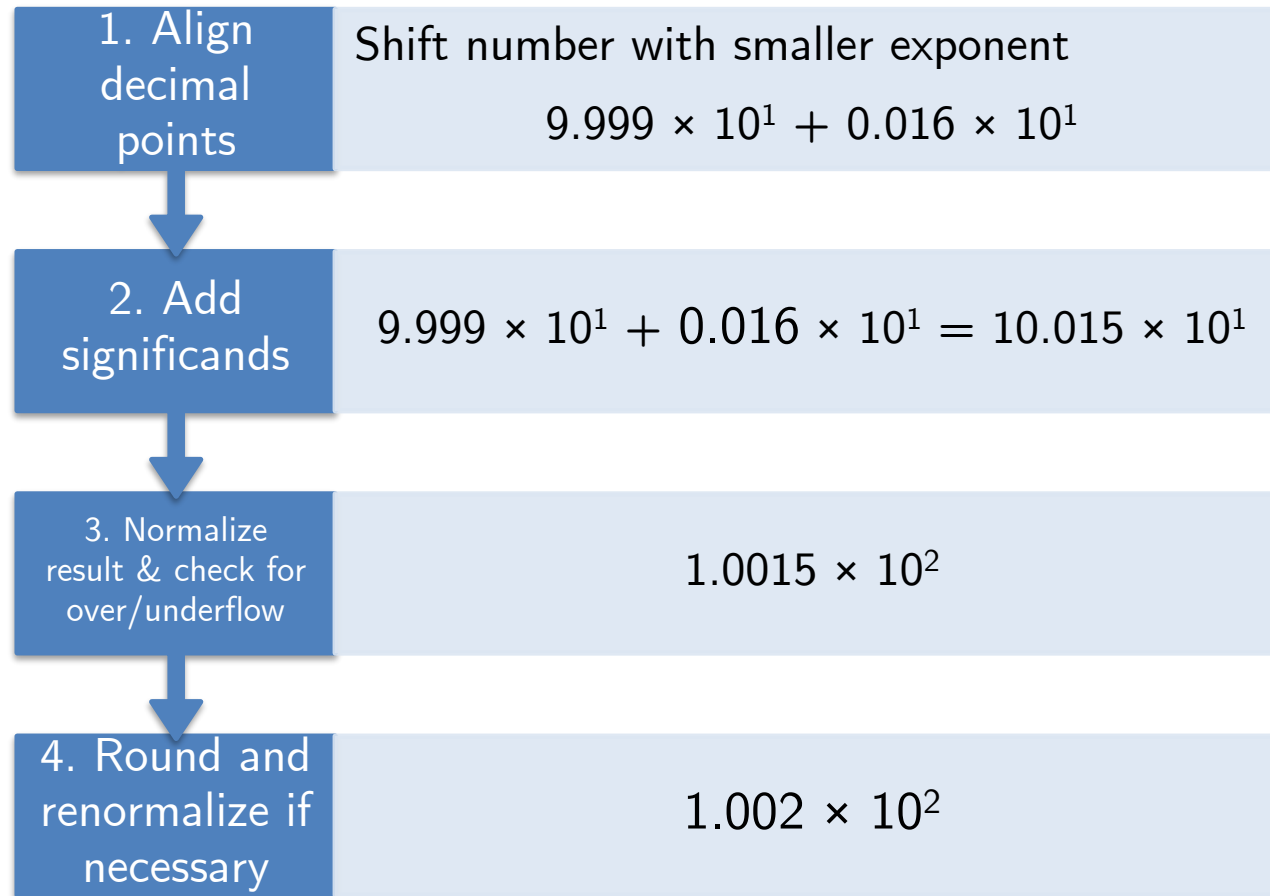
Single precision		Double precision		Object represented
Exponent	Fraction	Exponent	Fraction	
0	0	0	0	0
0	Nonzero	0	Nonzero	\pm denormalized number
1–254	Anything	1–2046	Anything	\pm floating-point number
255	0	2047	0	\pm infinity
255	Nonzero	2047	Nonzero	NaN (Not a Number)

Floating-point addition

Floating-point addition

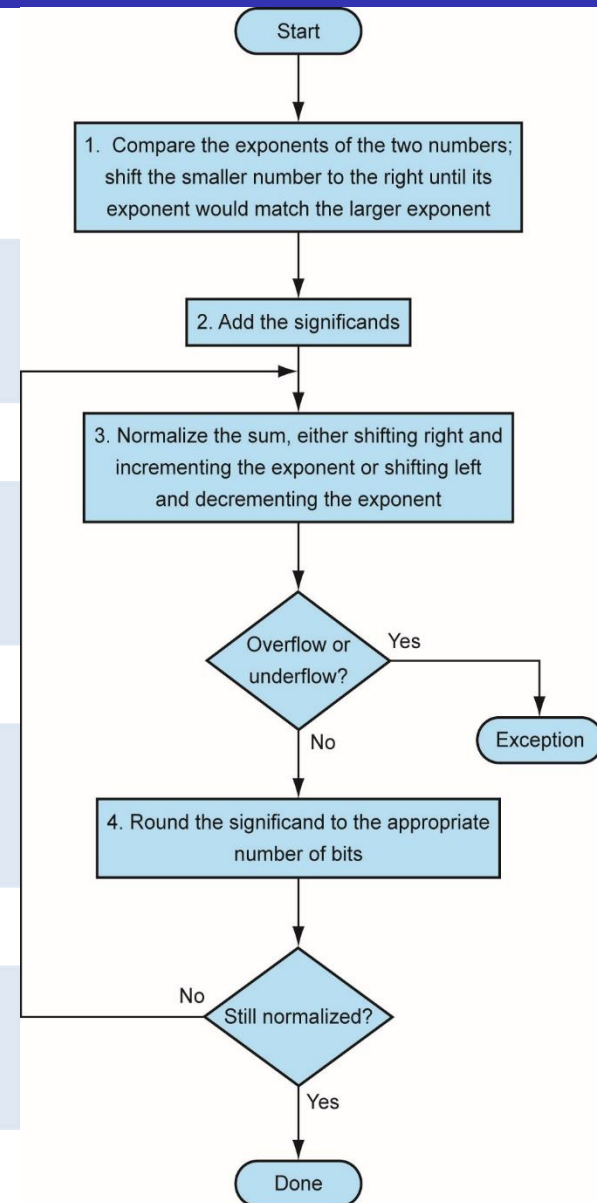
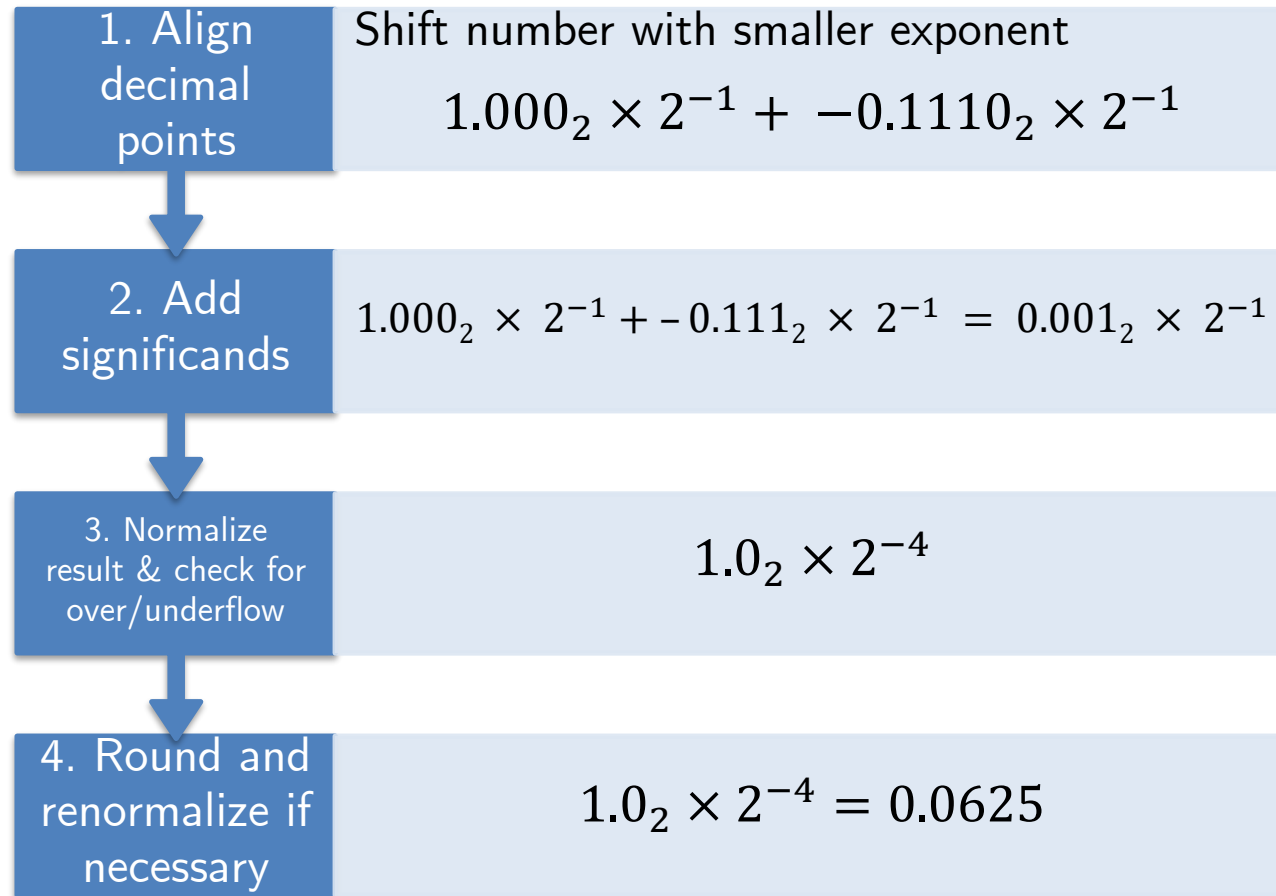
- Consider a decimal example

$$9.999 \times 10^1 + 1.610 \times 10^{-1}$$



Floating-point addition

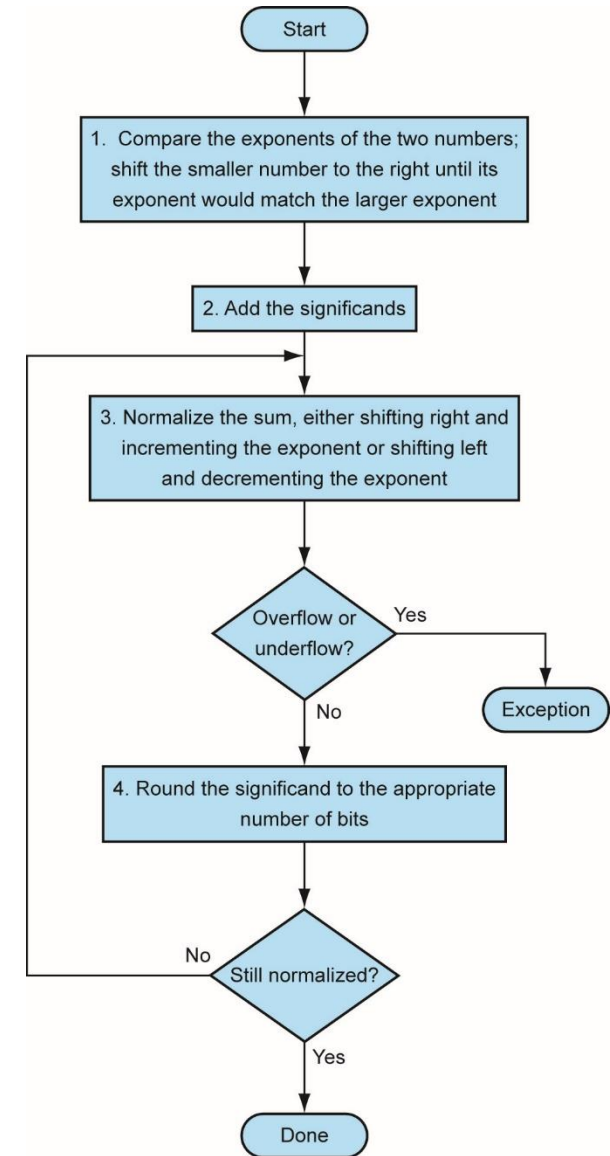
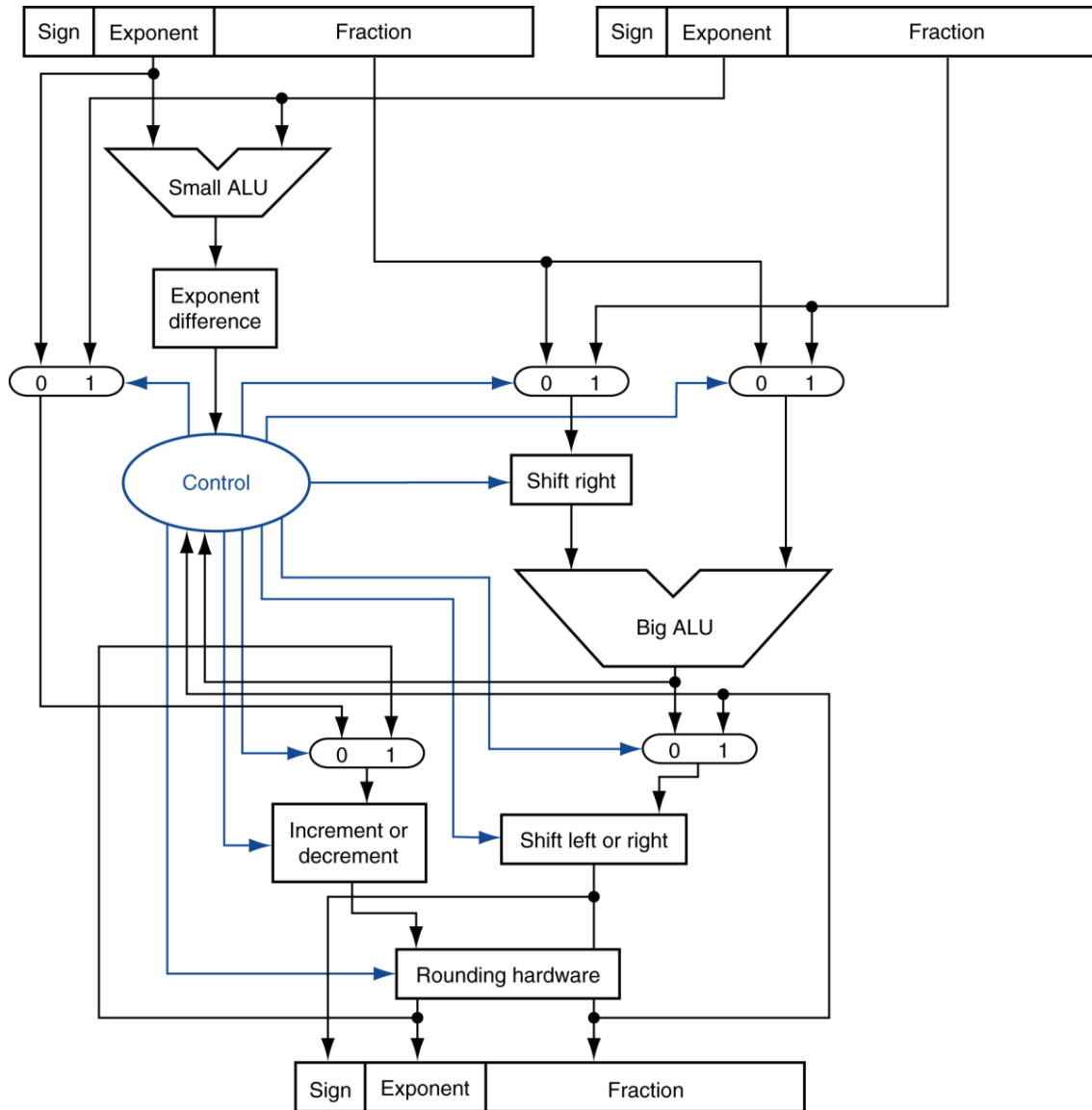
- Consider a floating-point example
 $1.000_2 \times 2^{-1} + -1.110_2 \times 2^{-2}$ or $(0.5 + -0.4375)$



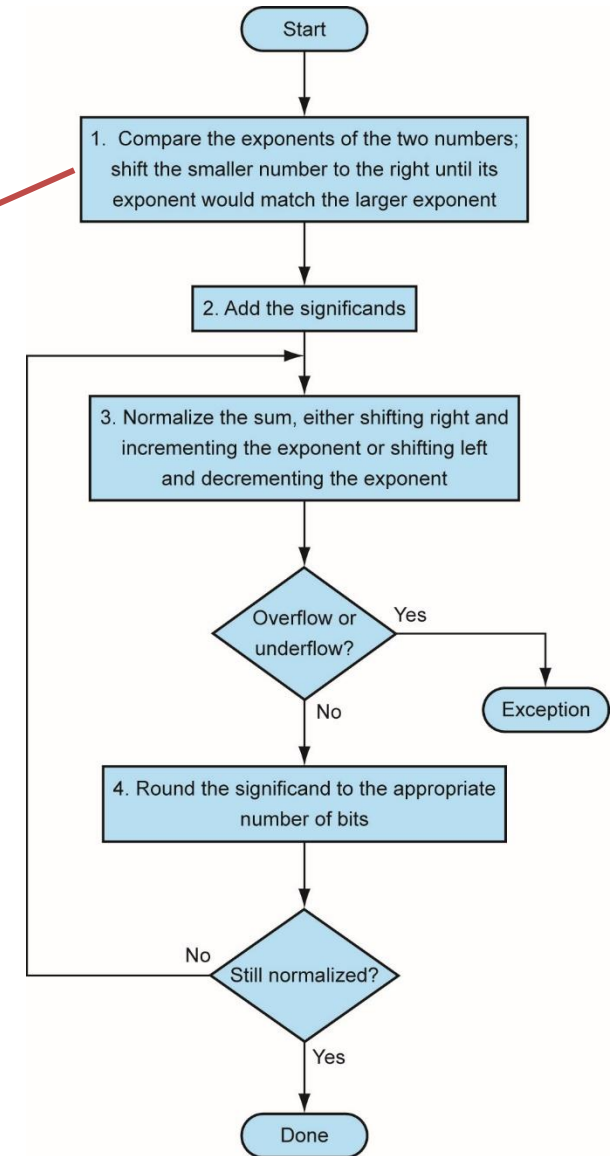
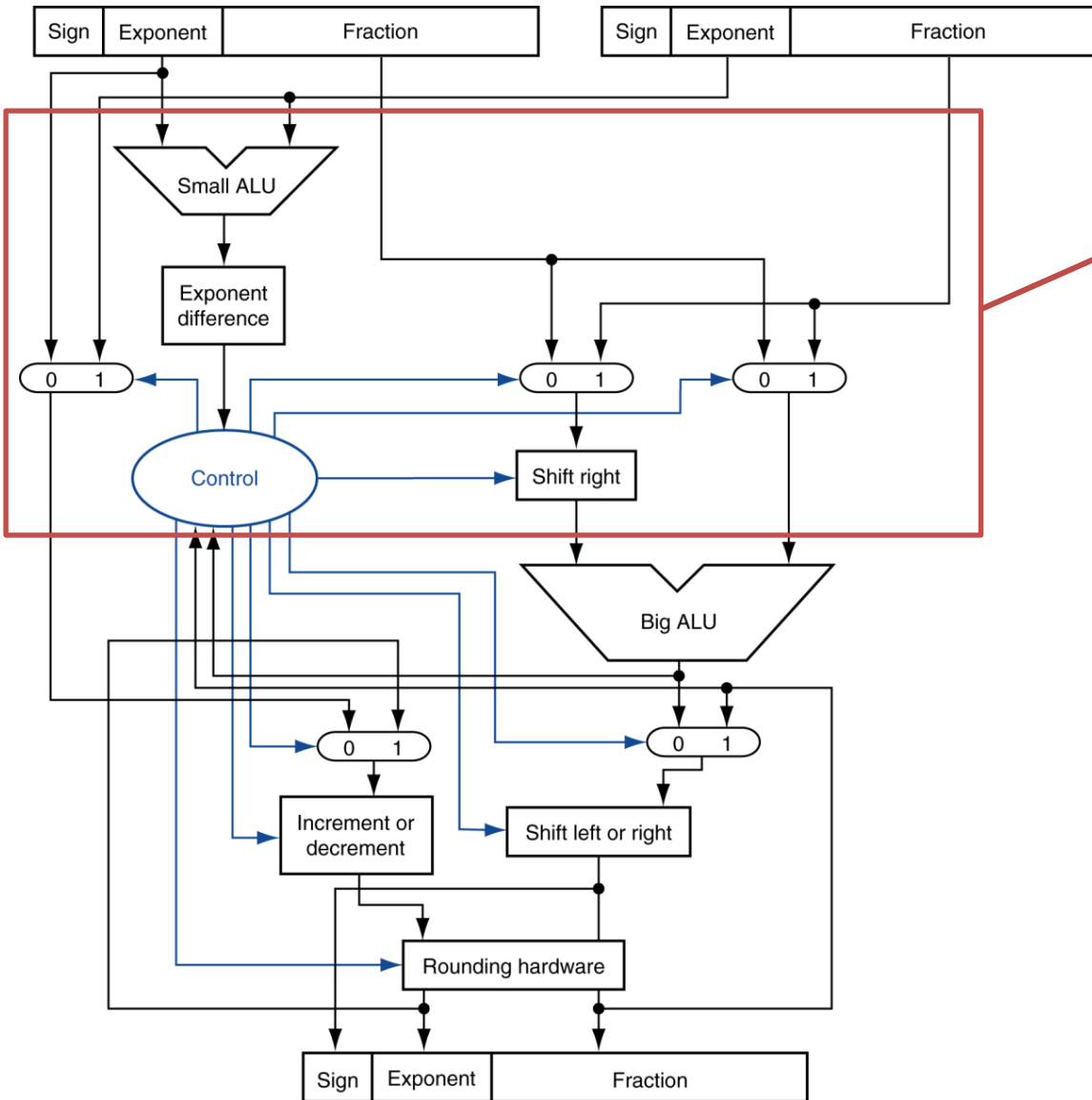
Floating-point adder hardware

- Much more complex than integer adder.
- Doing it in one clock cycle would take too long.
 - Much longer than integer operations.
 - Slower clock would penalize all instructions.
- Floating-point adder usually takes several cycles
 - Can be pipelined

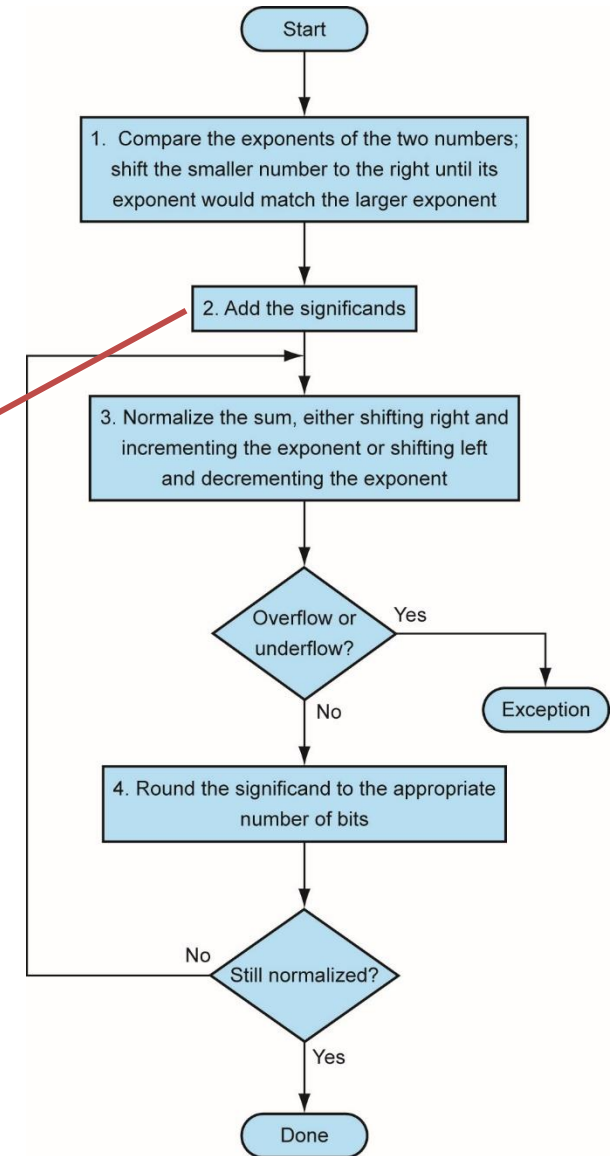
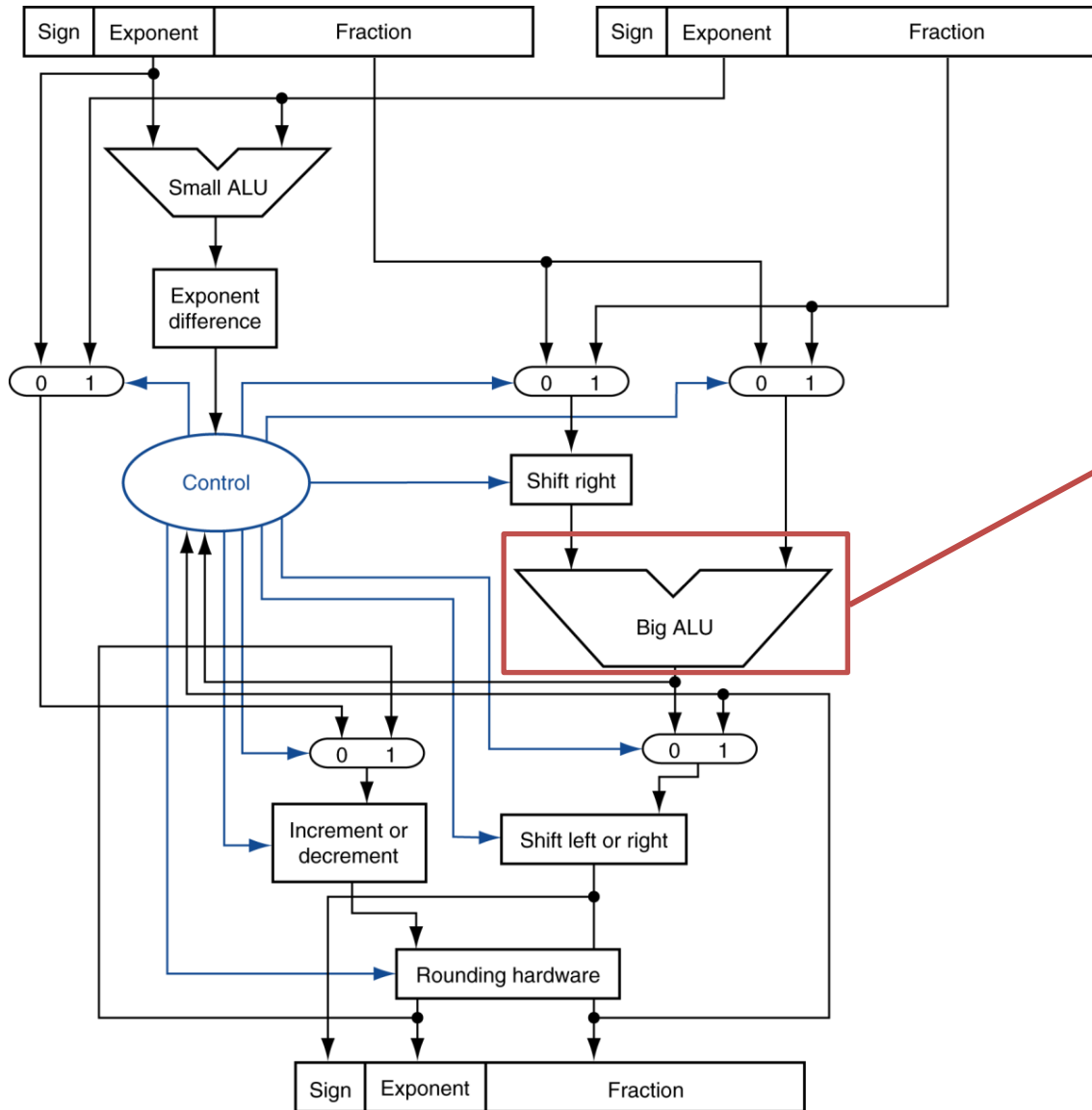
Floating-point adder hardware



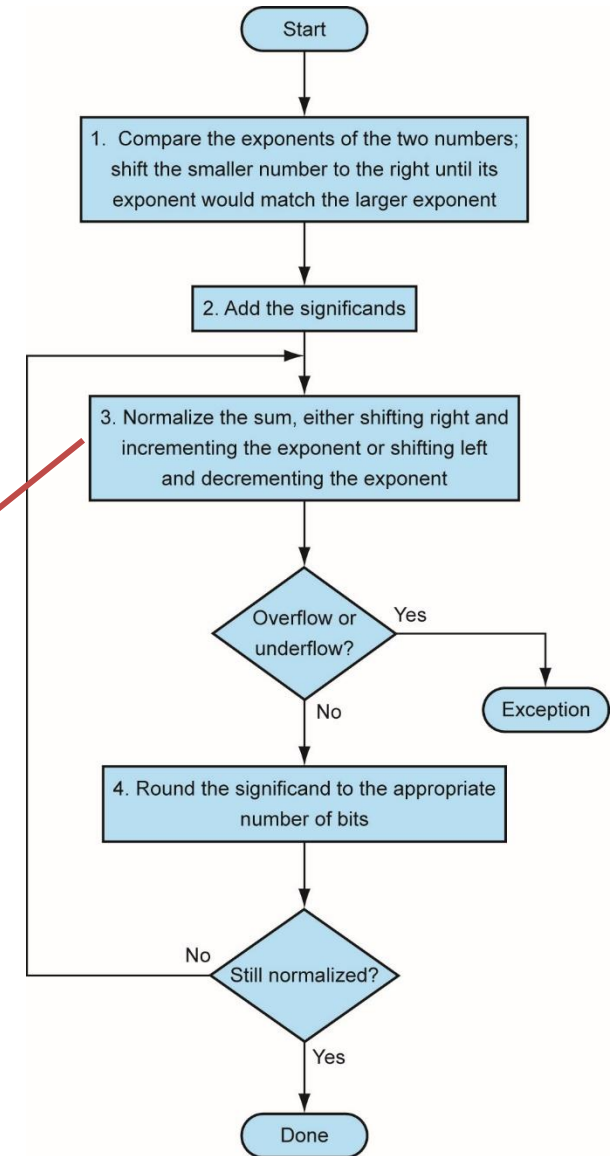
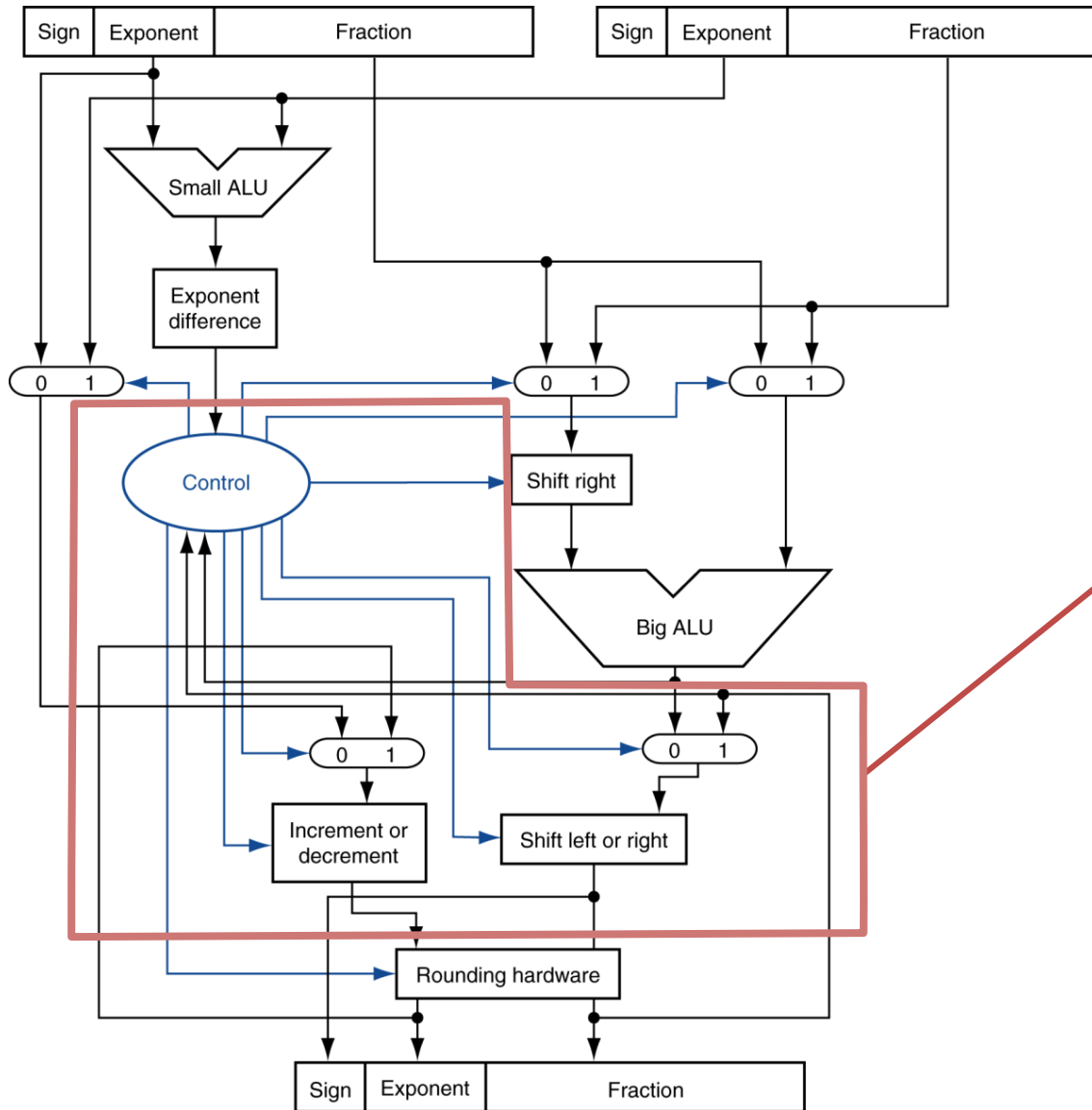
Floating-point adder hardware



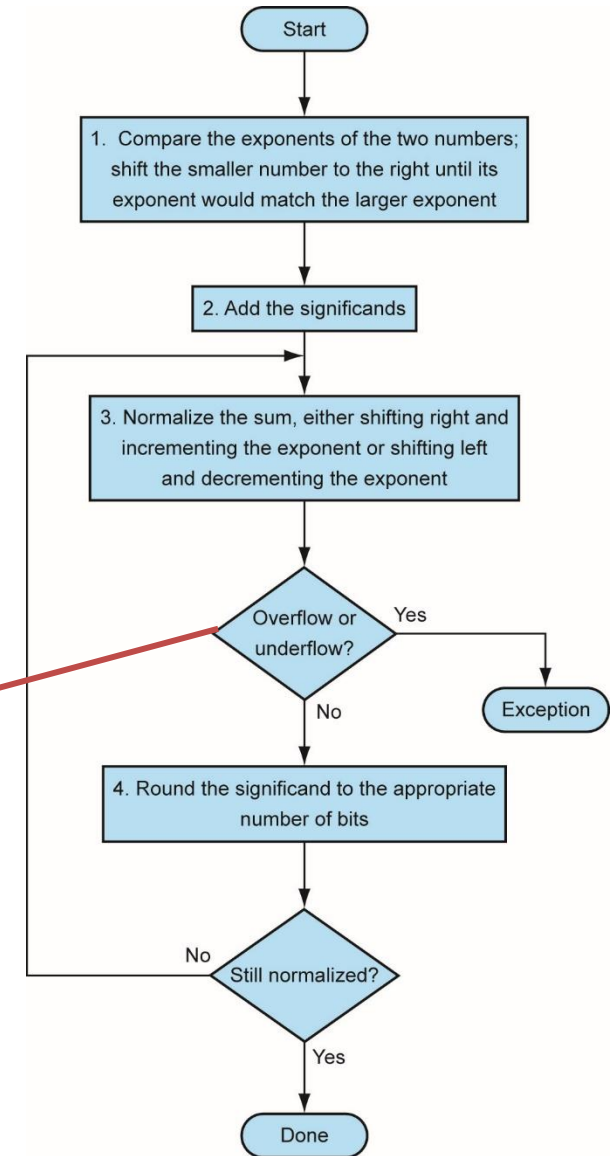
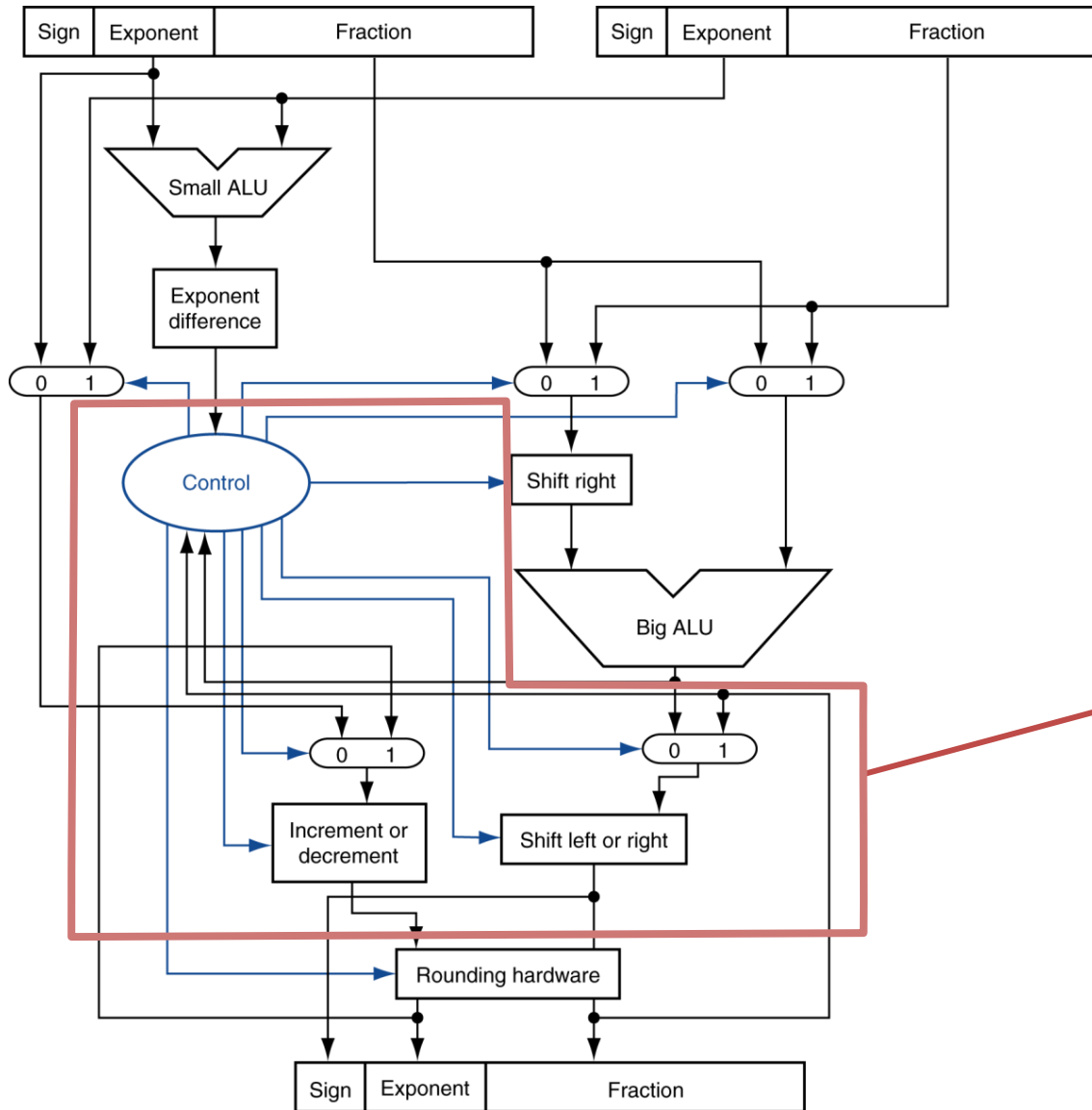
Floating-point adder hardware



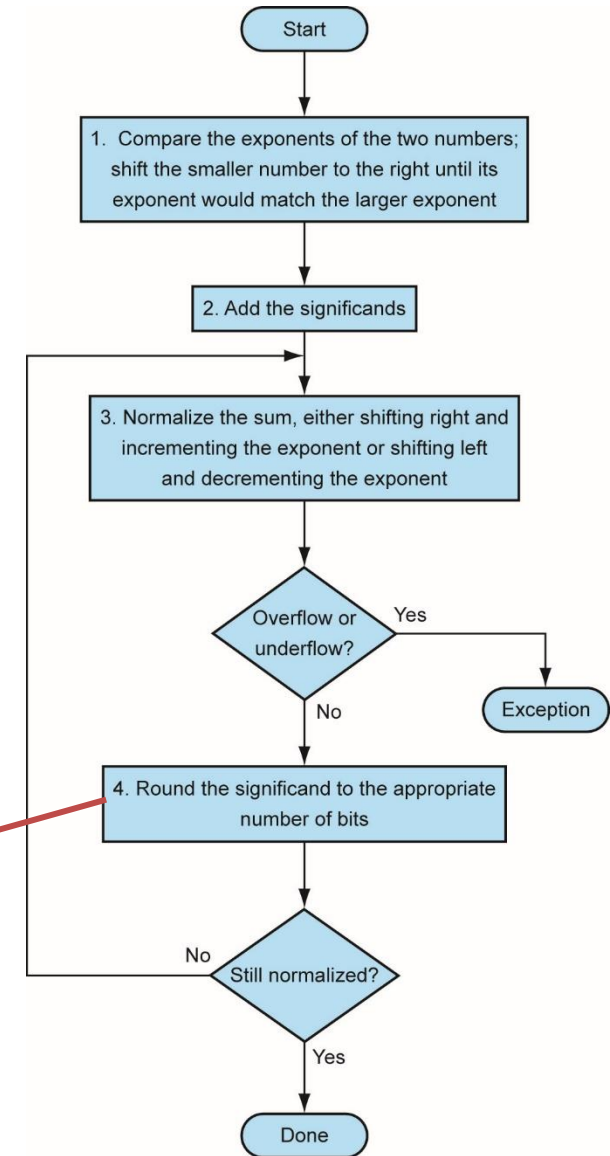
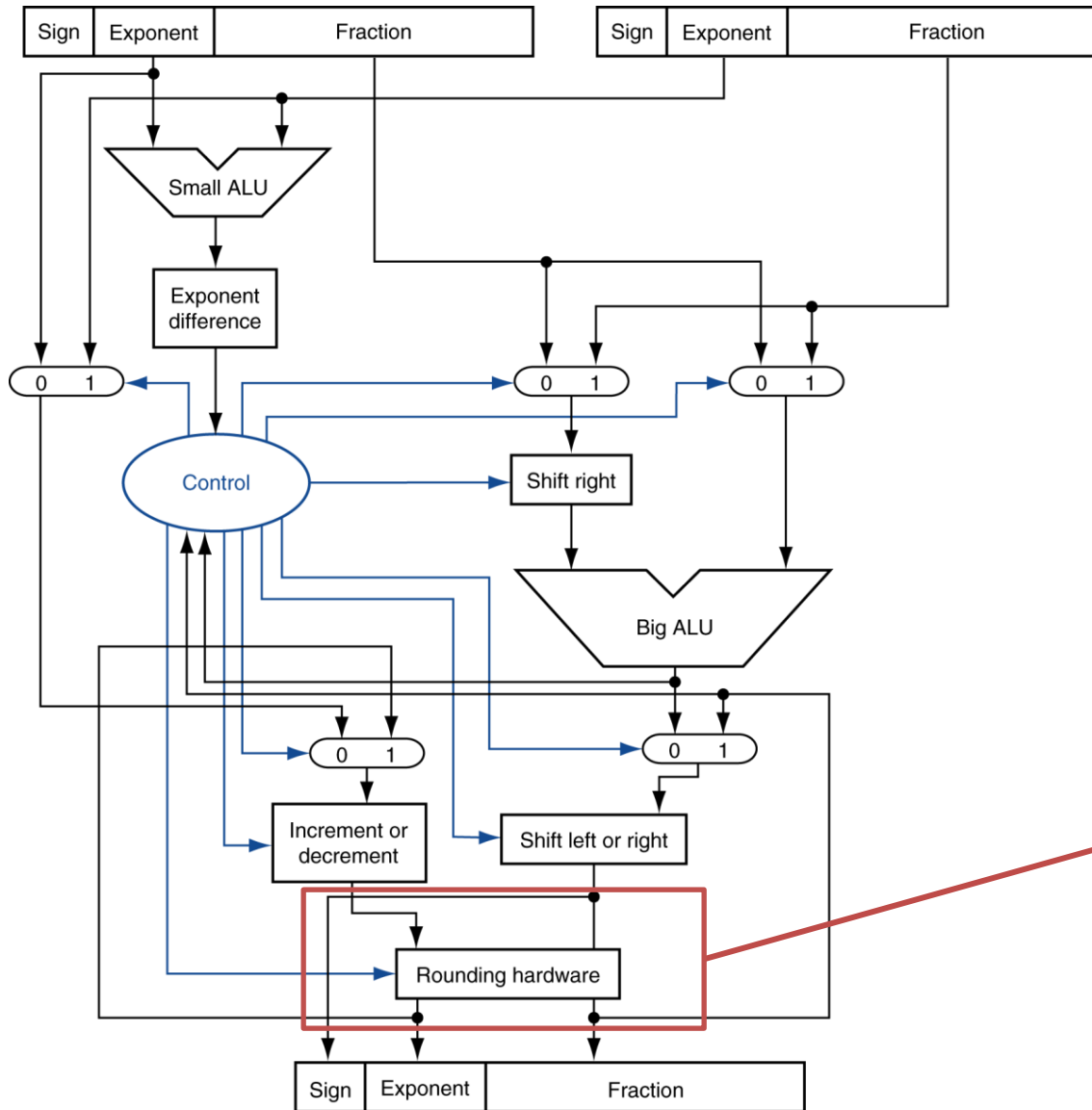
Floating-point adder hardware



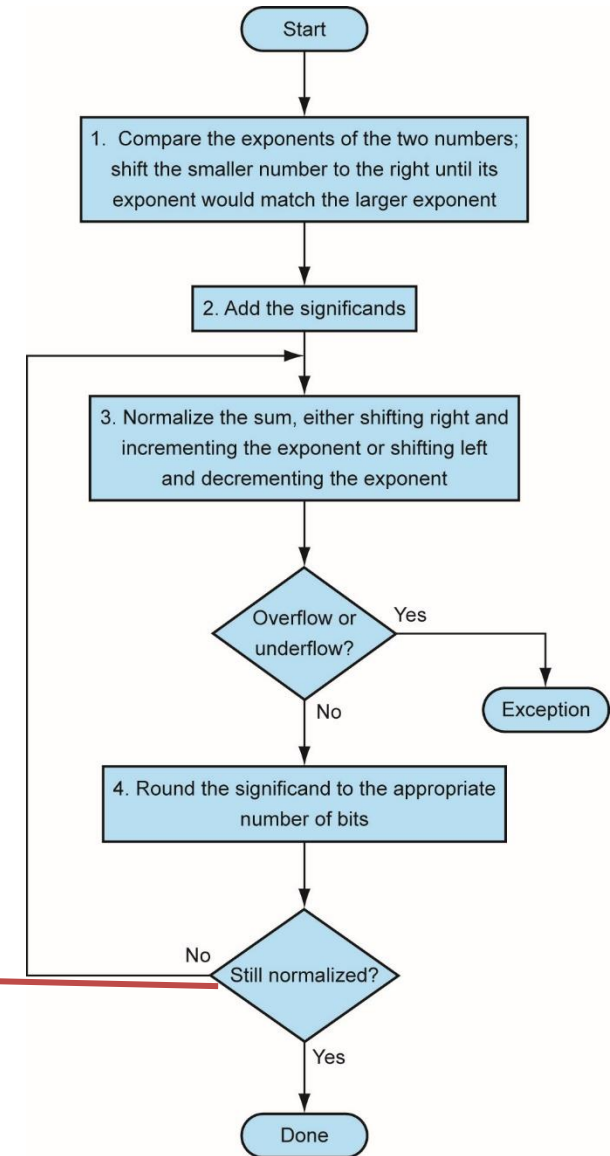
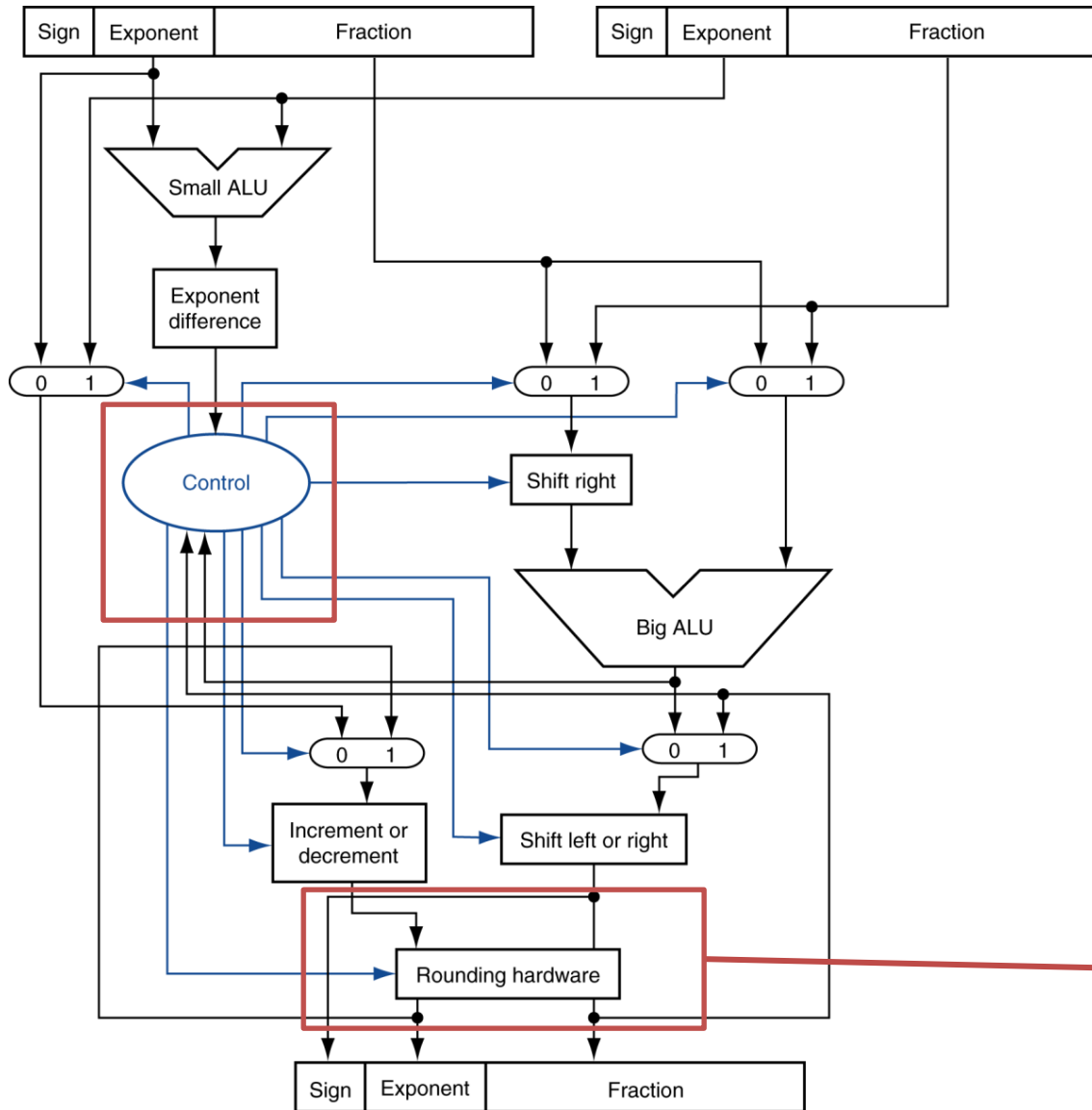
Floating-point adder hardware



Floating-point adder hardware



Floating-point adder hardware



Floating-point multiplication

Floating-point multiplication

Consider a decimal example

$$1.110 \times 10^{10} \times 9.200 \times 10^{-5}$$

1. Add exponents

- For biased exponents, subtract bias from sum
- New exponent = $10 + -5 = 5$

2. Multiply significands

$$\begin{aligned} &1.110 \times 9.200 \\ &= 10.212 \Rightarrow 10.212 \times 10^5 \end{aligned}$$

3. Normalize result & check for over/underflow

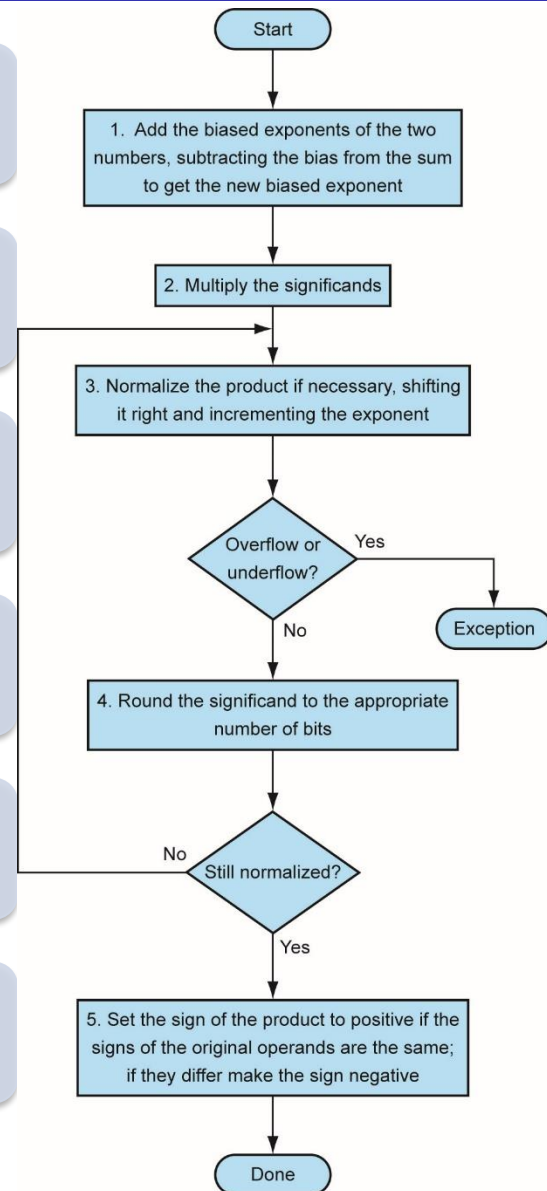
$$1.0212 \times 10^6$$

4. Round and renormalize if necessary

$$1.021 \times 10^6$$

5. Determine sign of result from signs of operands

$$+1.021 \times 10^6$$



Floating-point multiplication

Consider a floating-point example

$$-14.25 \times 3.125 \text{ in decimal, or} \\ -1.11001_2 \times 2^3 \times 1.1001_2 \times 2^1$$

1. Add exponents

- Unbiased: $3 + 1 = 4$
- Biased: $(3 + 127) + (1 + 127) - 127 = 4 + 254 - 127 = 4 + 127$

2. Multiply significands

$$\begin{aligned} &-1.11001_2 \times 1.1001_2 \\ &= 10.110010001_2 \\ &\Rightarrow 10.110010001_2 \times 2^4 \end{aligned}$$

3. Normalize result & check for over/underflow

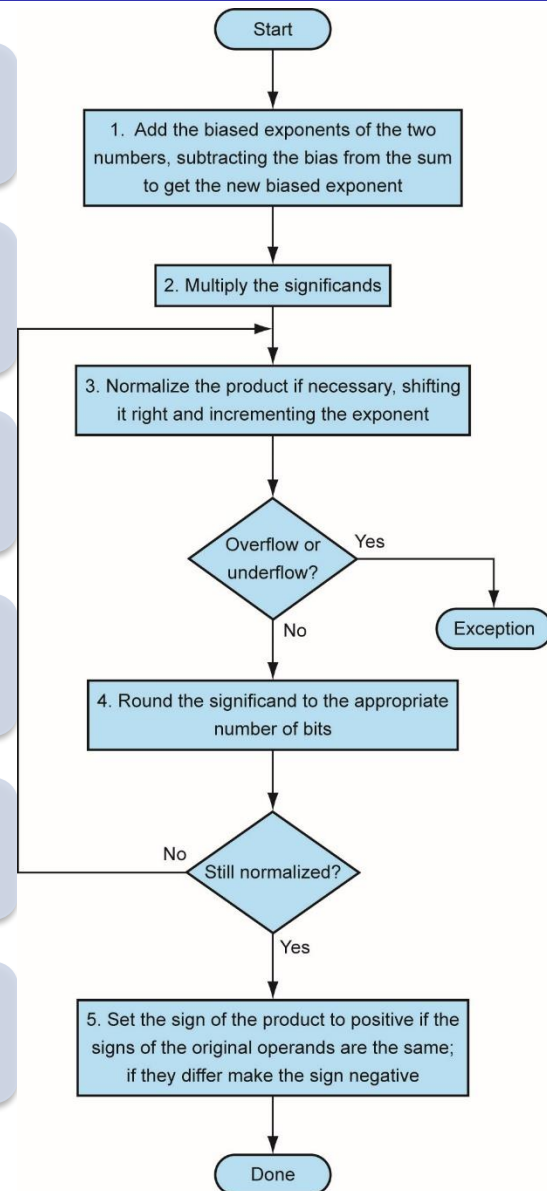
$$1.0110010001_2 \times 2^5 \text{ with no over/underflow}$$

4. Round and renormalize if necessary

$$1.0110010001_2 \times 2^5 \text{ (no change)}$$

5. Determine sign of result from signs of operands

$$-1.0110010001_2 \times 2^5 = -44.53125$$



Floating-point multiplication

- What's the floating-point of the previous result?

$$\begin{aligned} -14.25_{10} \times 3.125_{10} &= -44.53125_{10} \\ (-1.11001_2 \times 2^3) \times (1.1001_2 \times 2^1) &= -1.0110010001 \times 2^5 \end{aligned}$$

- sign = 1
- exponent
 - Single precision: $5 + 127 = 132 \rightarrow 10000100$
 - Double precision: $5 + 1023 = 1028 \rightarrow 10000000100$
- mantissa = 011001000100000000000000
- Floating point representation:

Single precision: 1 10000100 011001000100000000000000

Double precision: 1 10000000100 011001000100 ... 00

Floating-point arithmetic hardware

- Floating-point multiplier is of similar complexity to floating-point adder
 - But uses a multiplier for significands instead of an adder
- Floating-point arithmetic hardware usually does
 - Addition, subtraction, multiplication, division, reciprocal, square-root
 - Floating-point \leftrightarrow integer conversion
- Operations usually takes several cycles
 - Can be pipelined

Floating-point rounding

Floating-point rounding

- Let's try to represent 0.1_{10} in floating-point.

1. Represent 0.1_{10} in binary

1. Integer part is 0

2. Fractional part is:

$$0.1 \times 2 = 0 + 0.2$$

$$0.2 \times 2 = 0 + 0.4$$

$$0.4 \times 2 = 0 + 0.8$$

$$0.8 \times 2 = 1 + 0.6$$

$$0.6 \times 2 = 1 + 0.2$$

$$0.2 \times 2 = 0 + 0.4$$

$$0.4 \times 2 = 0 + 0.8$$

$$0.8 \times 2 = 1 + 0.6$$

$$0.6 \times 2 = 1 + 0.2$$

$$\vdots$$

0.00011

This sequence repeats infinite times!

0.1 is not a machine number, which means, it may not be exactly represented in a computing system.

Floating-point rounding

- Our floating-point representation will have to be as close as possible to 0.1_{10}
 $0.00011001100110011001100110011001 \dots_2$
- Remember that mantissa is 23 and 52 bits for single- and double-precision, respectively.
- IEEE employs rounding.
 - If first extra bit is 1, we add 1 to the rest of the mantissa bits – This is rounding up
 - If first extra bit is 0, we drop all extra bits – This is rounding down.
 - Special case if extra bits are 1000....000
 - Round up if last mantissa bit is 1
 - Round down if last mantissa bit is 0

Floating-point rounding

- Rounded normalized value:

$$1.10011001100110011001100110\mathbf{1}_2 \times 2^{-4}$$

We **rounded up** for this example

$$1.10011001100110011001100110\mathbf{1}_2 \times 2^{123-127} \text{ (single)}$$

- Floating-point representation:

Single: 1 0111011 10011001100110011001101

- Which represents the value of

$$0.10000000\mathbf{1490116119384765625}$$

- Similarly, double precision represents 0.1 as

$$0.100000000000000000\mathbf{55511151231257827021181583404541015625}$$

Floating-point rounding

- In your favourite programming language try the following code using float or double data types.

$0.1 + 0.1 + 0.1 == 0.3$

Is the result TRUE or FALSE?

- Problems with accuracy
 - Several failures (in some cases with fatal consequences) have been reported due to numerical errors.

<http://ta.twi.tudelft.nl/users/vuik/wi211/disasters.html>

Concluding Remarks

- Bits have no inherent meaning
 - Interpretation depends on the instructions applied
- Computer representations of numbers
 - Finite range and precision
 - Need to account for this in programs

Concluding Remarks

- ISAs support arithmetic operations
 - Signed and unsigned integers.
 - Floating-point approximation to reals.
- Bounded range and precision
 - Operations can overflow and underflow