



Chapter 3

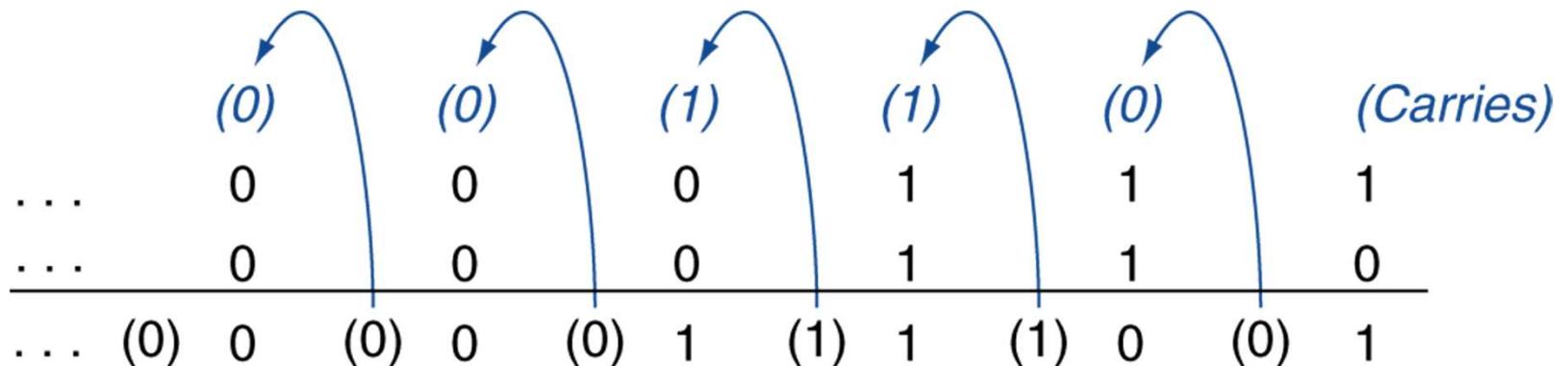
Arithmetic for Computers

Arithmetic for Computers

- Operations on integers
 - Addition and subtraction
 - Multiplication and division
 - Dealing with overflow
- Floating-point real numbers
 - Representation and operations

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 Subtraction

- Add negation of second operand
- Example: $7 - 6 = 7 + (-6)$

+7:	0000 0000 ... 0000 0111
-6:	1111 1111 ... 1111 1010
<hr/>	
+1:	0000 0000 ... 0000 0001
- 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

Dealing with Overflow

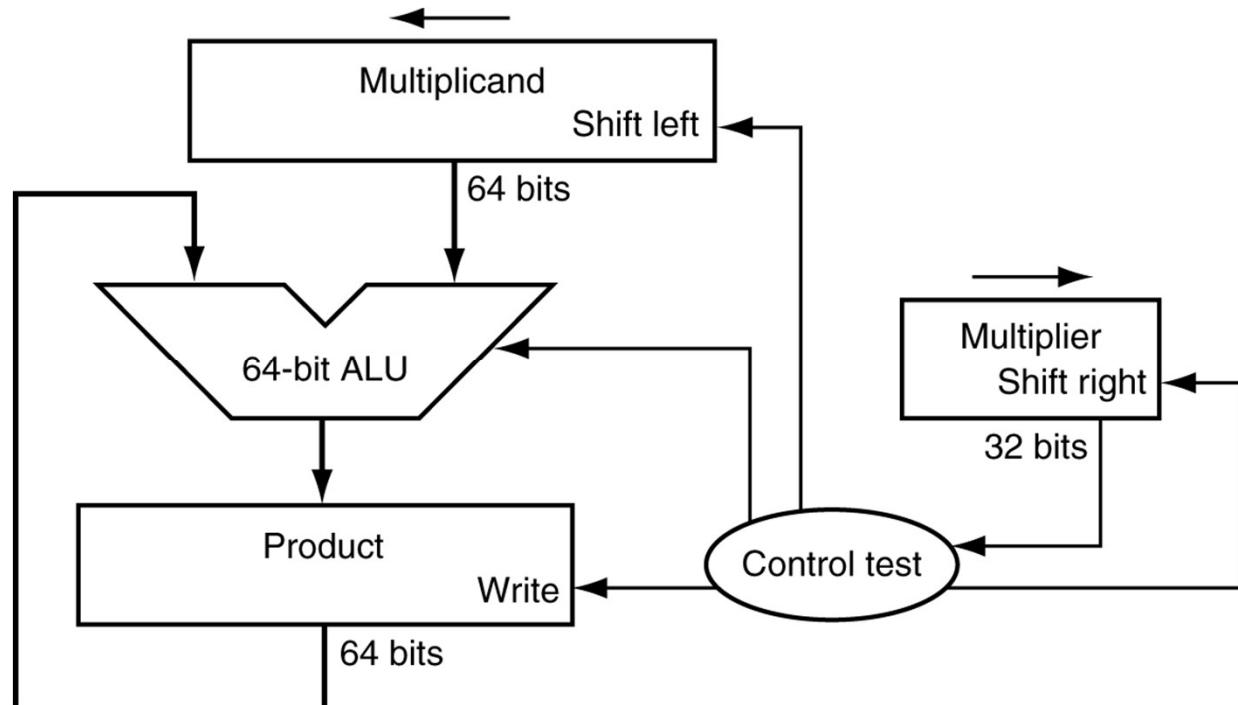
- Some languages (e.g., C) ignore overflow
 - Use MIPS addu, addui, subu instructions
- Other languages (e.g., Ada, Fortran) require raising an exception
 - Use MIPS add, addi, sub instructions
 - On overflow, invoke exception handler
 - Save PC in exception program counter (EPC) register
 - Jump to predefined handler address
 - mfc0 (move from system control reg) instruction can retrieve EPC value, to return after corrective action

Multiplication

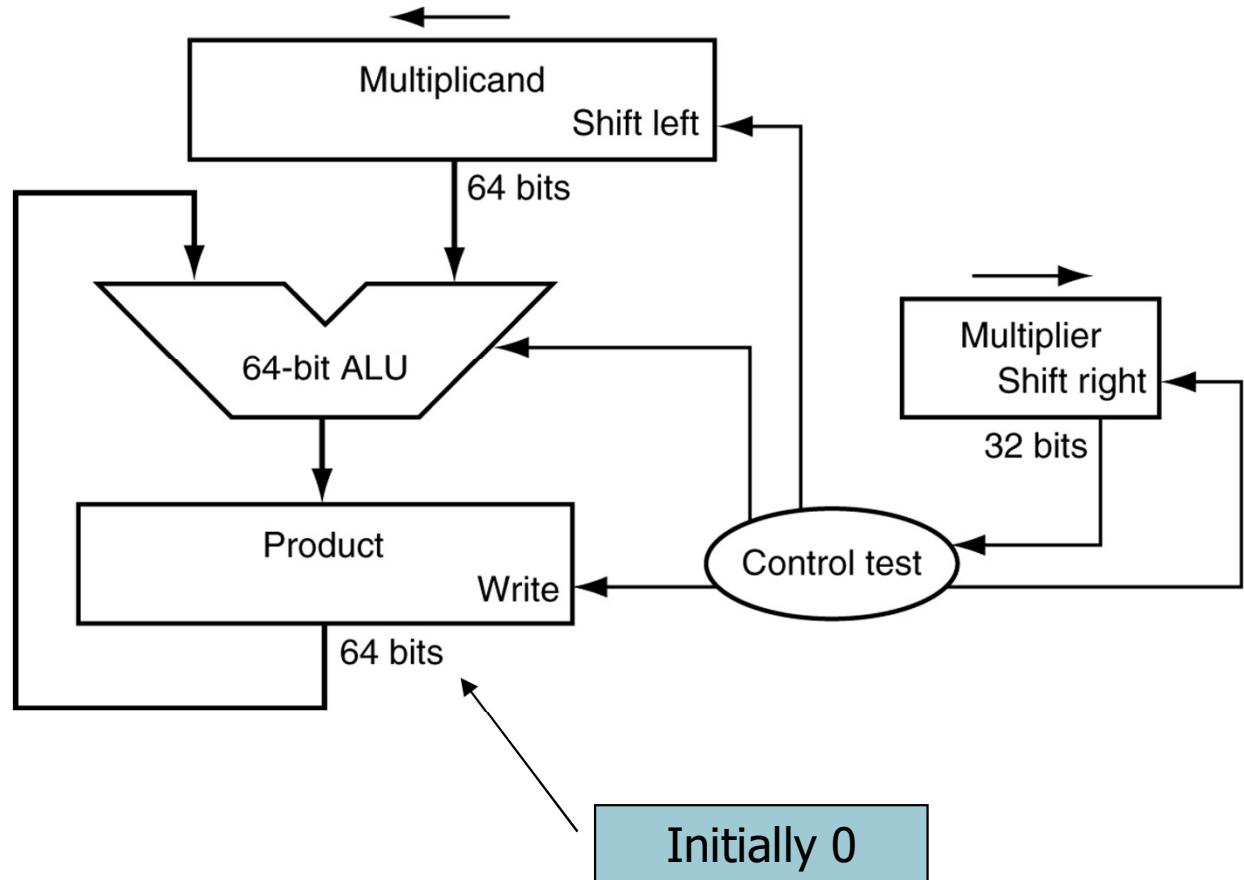
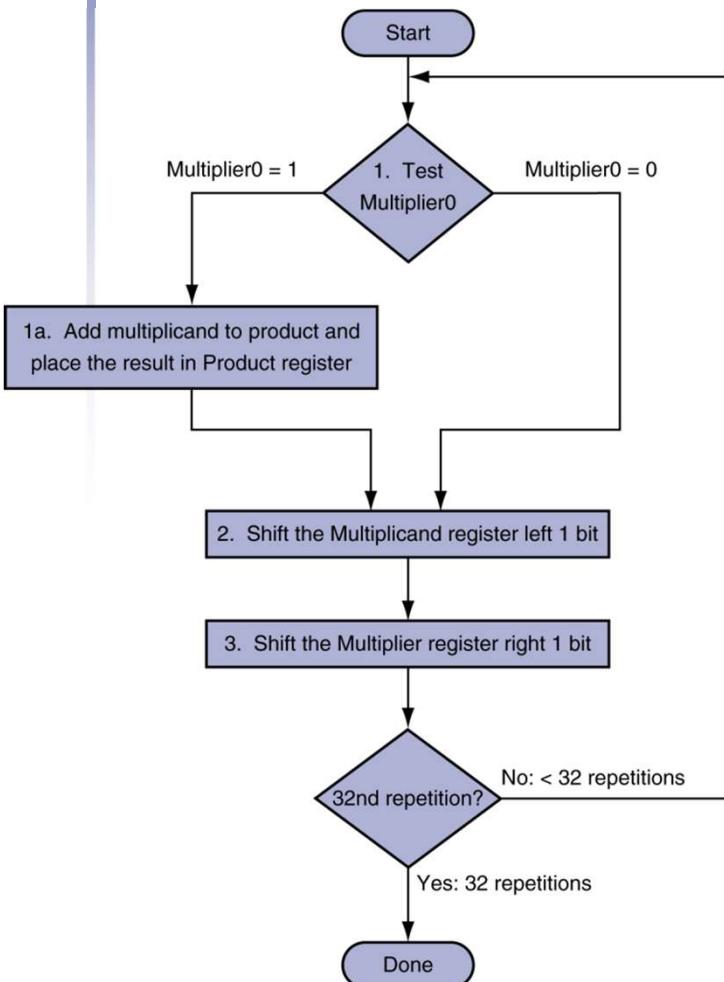
- Start with long-multiplication approach

multiplicand	1000
multiplier	x 1001
	<hr/>
	1000
	0000
	0000
product	1000
	<hr/>
	1001000

Length of product is
the sum of operand
lengths

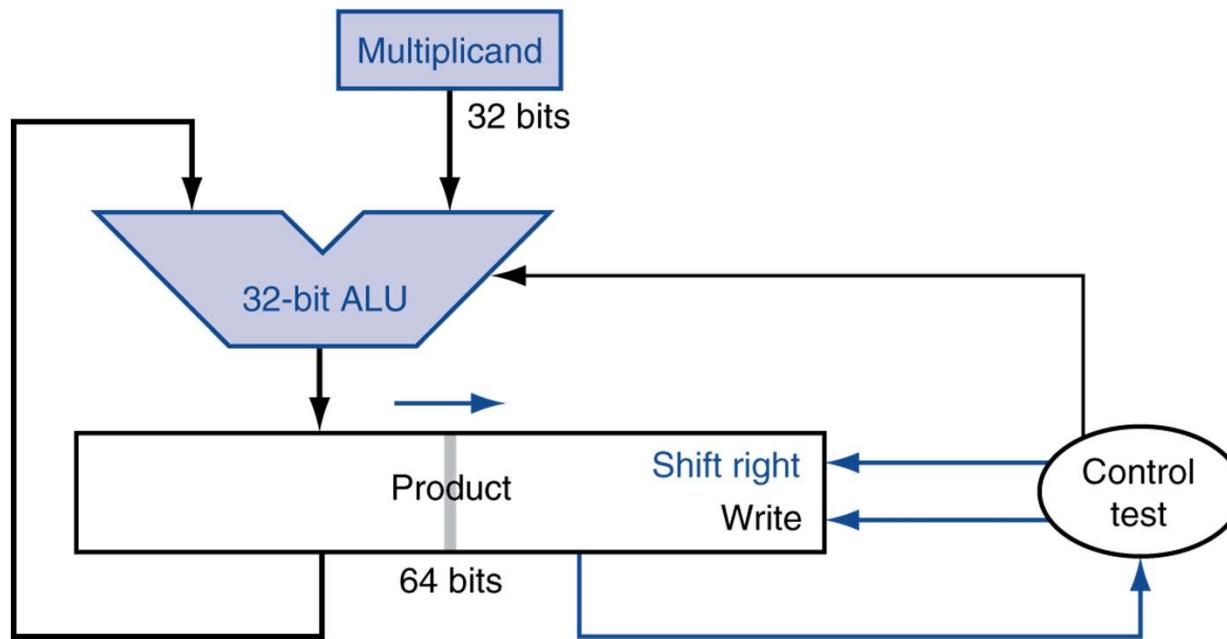


Multiplication Hardware



Optimized Multiplier

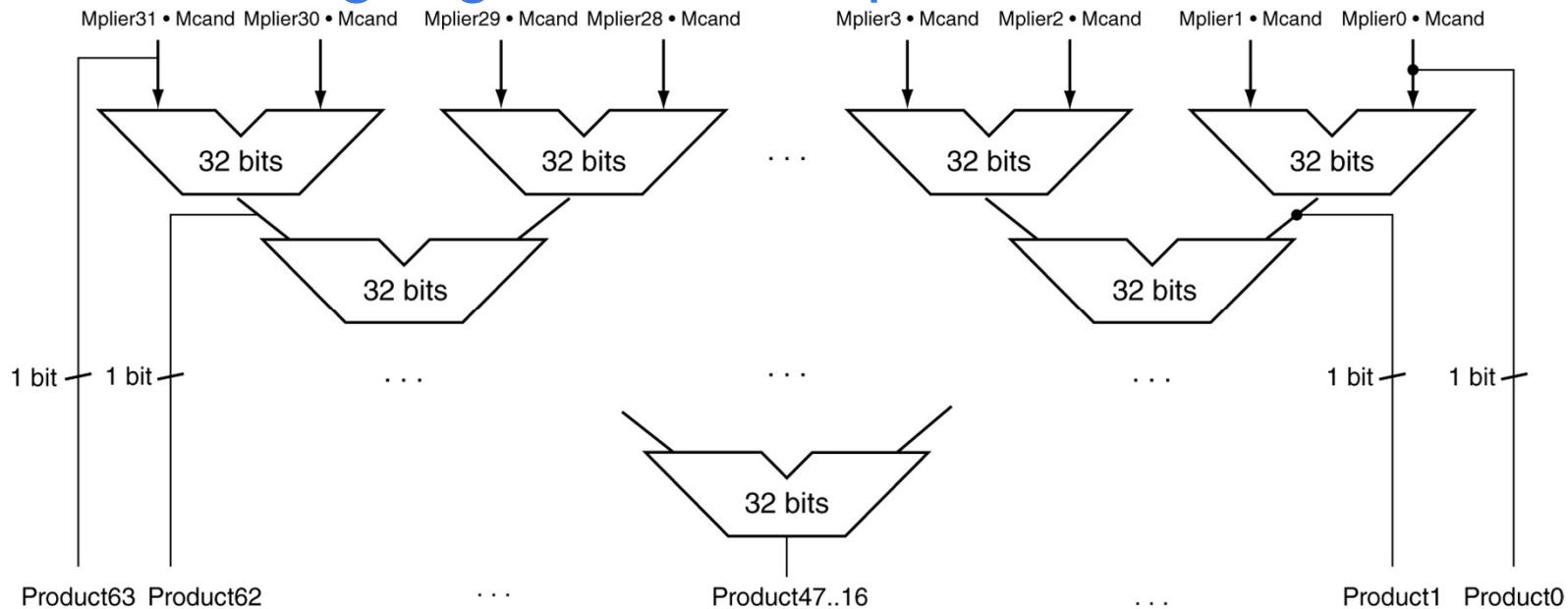
- Perform steps in parallel: add/shift



- One cycle per partial-product addition
 - That's ok, if frequency of multiplications is low

Faster Multiplier

- Uses multiple adders
 - Cost/performance tradeoff
 - Following Figure is incomplete

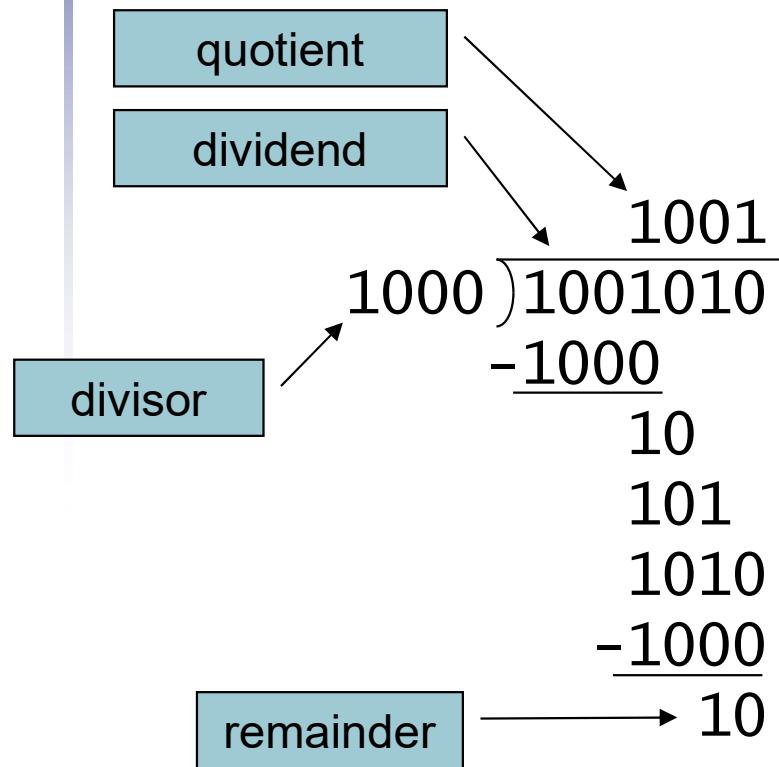


- Can be pipelined
 - Several multiplication performed in parallel

MIPS Multiplication

- Two 32-bit registers for product
 - HI: most-significant 32 bits
 - LO: least-significant 32-bits
- Instructions
 - `mult rs, rt / multu rs, rt`
 - 64-bit product in HI/LO
 - `mfhi rd / mflo rd`
 - Move from HI/LO to rd
 - Can test HI value to see if product overflows 32 bits
 - `mul rd, rs, rt`
 - Least-significant 32 bits of product → rd

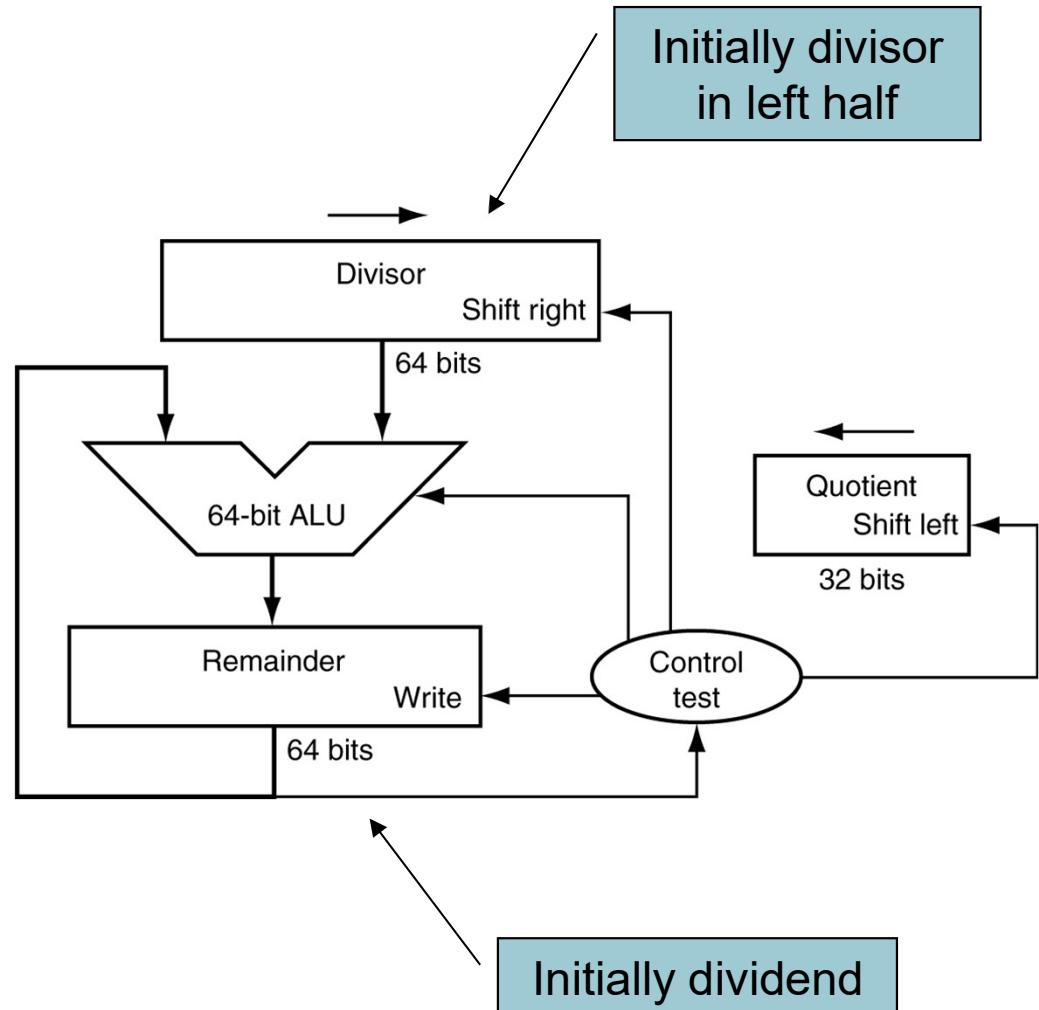
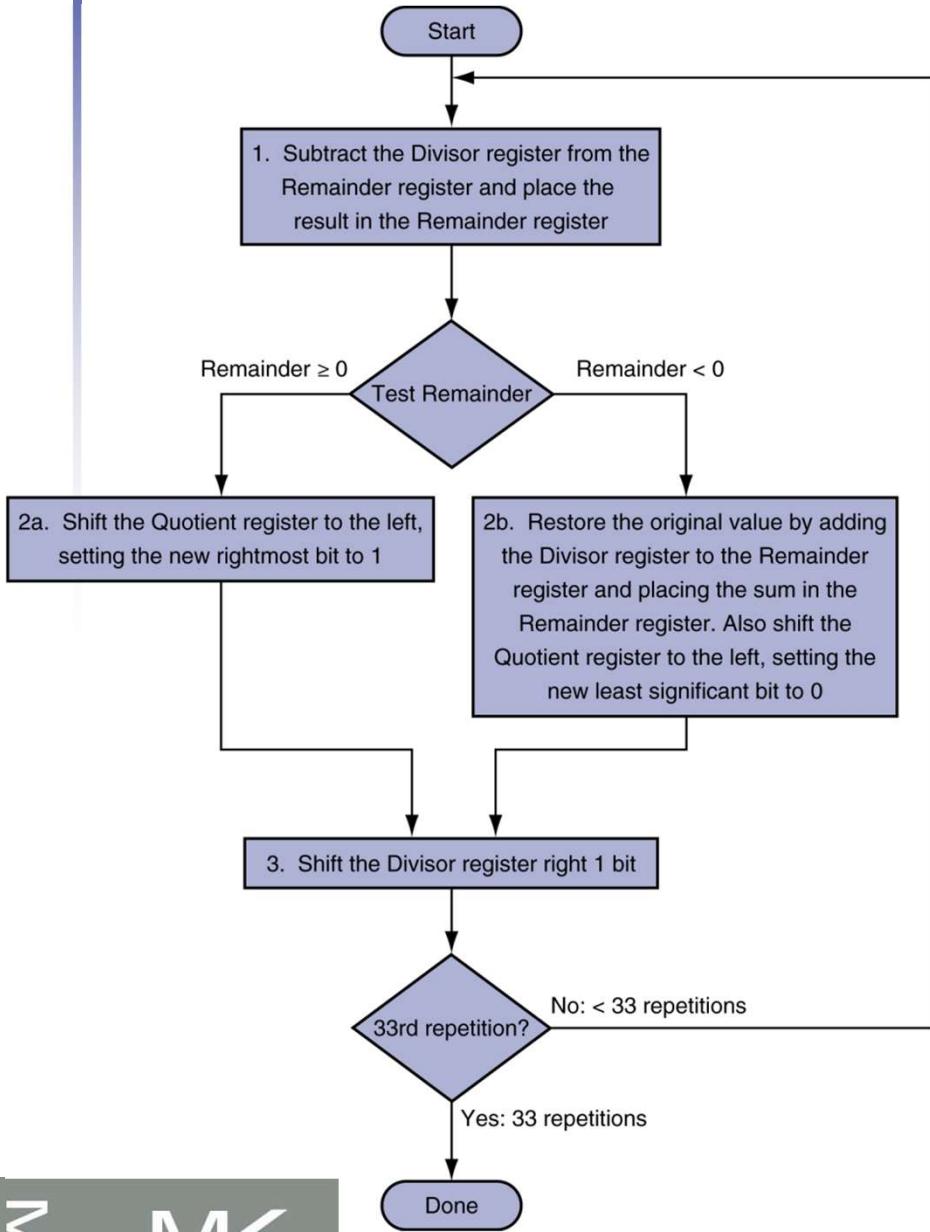
Division



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

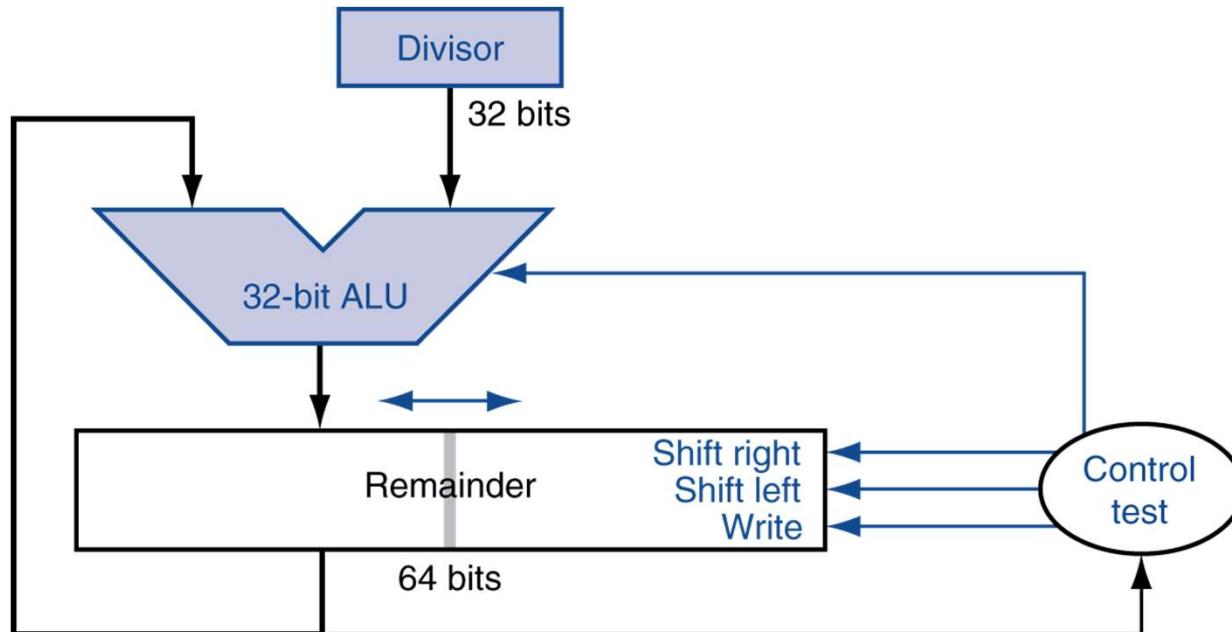
Division Hardware



Division Example

Iteration	Step	Quotient	Divisor	Reminder
0	Initial values	0000	0010 0000	0000 0111
1	1: Rem = Rem – Div	0000	0010 0000	①110 0111
	2b: Rem < 0 \Rightarrow +Div, sll Q, Q0 = 0	0000	0010 0000	0000 0111
	3: Shift Div right	0000	0001 0000	0000 0111
2	1: Rem = Rem – Div	0000	0001 0000	①111 0111
	2b: Rem < 0 \Rightarrow +Div, sll Q, Q0 = 0	0000	0001 0000	0000 0111
	3: Shift Div right	0000	0000 1000	0000 0111
3	1: Rem = Rem – Div	0000	0000 1000	①111 1111
	2b: Rem < 0 \Rightarrow +Div, sll Q, Q0 = 0	0000	0000 1000	0000 0111
	3: Shift Div right	0000	0000 0100	0000 0111
4	1: Rem = Rem – Div	0000	0000 0100	②000 0011
	2a: Rem \geq 0 \Rightarrow sll Q, Q0 = 1	0001	0000 0100	0000 0011
	3: Shift Div right	0001	0000 0010	0000 0011
5	1: Rem = Rem – Div	0001	0000 0010	②000 0001
	2a: Rem \geq 0 \Rightarrow sll Q, Q0 = 1	0011	0000 0010	0000 0001
	3: Shift Div right	0011	0000 0001	0000 0001

Optimized Divider



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

Optimized Divider

Iteration	Step	Divisor)	Remainder (R/ Q)	
0	Init	0010	0000	0111
1	Rem = Rem - Sub	0010	Negative!	0111
	Restore Rem Left Shift and Set 0	0010	0000 0000	0111 1110
2	Rem = Rem – Sub	0010	Negative!	1110
	Restore Rem Left Shift and Set 0	0010	0000 0001	1110 1100
3	Rem = Rem – Sub	0010	Negative!	1100
	Restore Rem Left Shift and Set 0	0010	0001 0011	1100 1000
4	Rem = Rem – Sub	0010	0001	1000
	Left Shift and Set 1	0010	0011	0001
5	Rem = Rem – Sub	0010	0001	0001
	Left Shift and Set 1 (special shift)	0010	0001	0011

Optimized Multiplier

Iteration	Step	Multiplicand	Result / Multiplier	
0	Init	0111	0000	0101
1	LSB =1, Add	0111	0111	0101
	RS	0111	0011	1010
2	LSB =0, Pass	0111	0011	1010
	RS	0111	0001	1101
3	LSB =1, Add	0111	1000	1101
	RS	0111	0100	0110
4	LSB =0, Pass	0111	0100	0110
	RS	0111	0010	0011

MIPS Division

- Use HI/LO registers for result
 - HI: 32-bit remainder
 - LO: 32-bit quotient
- Instructions
 - `div rs, rt` / `divu rs, rt`
 - No overflow or divide-by-0 checking
 - Software must perform checks if required
 - Use `mfhi`, `mflo` to access result
 - Q in LO, R in HI

Floating Point

10.00
↑
decimal point

scientific notation \rightarrow base 45 power 4
fifteen

- Representation for non-integral numbers
 - Including very small and very large numbers

- Like scientific notation

■ -2.34×10^{56}

normalized

■ $+0.002 \times 10^{-4}$

not normalized

■ $+987.02 \times 10^9$

non-zero

- In binary

■ $\pm 1.xxxxxxx_2 \times 2^{yyyy}$

$6.673 \times 10^{-31} \rightarrow g, e, c$ (fraction)
 $\rightarrow law$
(standard ন)

kg \rightarrow standard

'standard is not the law'

standardization \rightarrow authority limit /
baseline set করা
ভাব

authority \rightarrow IEEE ...

- Types float and double in C

\rightarrow something, decimal point
যান্তের digit forcefully 1

scientific notation \rightarrow decimal point এর স্থানে ফাট থাকবে

যদি decimal point এর স্থানে একটি digit থাকে

(non-zero)
Normalized standard form

$\pm 1.(\dots)_2 \times 2^{\square}$ এটিও বাইনারীতে থাকবে
১) বাইনারীতে ন্য নিলেও হচ্ছে
বাইনারীতে থাকবে

exponent term যদিও hardware কে binary তে থাকে

অধিক decimal পাখিয়ে কাট বুবো,

binary representation নিখিতে বলা হলে \rightarrow convert to

binary .

$$(1^0 1)_2 \rightarrow 1 \times 2^2 + 0 \times 2^1 + 1 \times 2^0$$

১
 \rightarrow এই ২ মূলে থাকে ন

$$\pm 1 \cdot (\dots)_2 \times 2^{\boxed{}}$$

এই part মুলো
যাবাই লাগবে

1. fraction \rightarrow implicit 1 / floating point এর পরের part
 2. exponent
 3. sign

২  এই power এর উপর depend করে point এর
যথেত্ত্ব জায়াম নিতে আরি।

- float করছে

এই floating point কো হয়।

যদি 16 bit architecture ইয় -
 MIPS \rightarrow 32 bit architecture
 \hookrightarrow ALU 32 bit - operand একটি 32 bit এর হতে হবে
 \rightarrow register 32 bit

$\frac{+}{-} 1. (\dots)$

sign representation

1 bit — লাগবে $(32-1) = 31$ bit বাকি

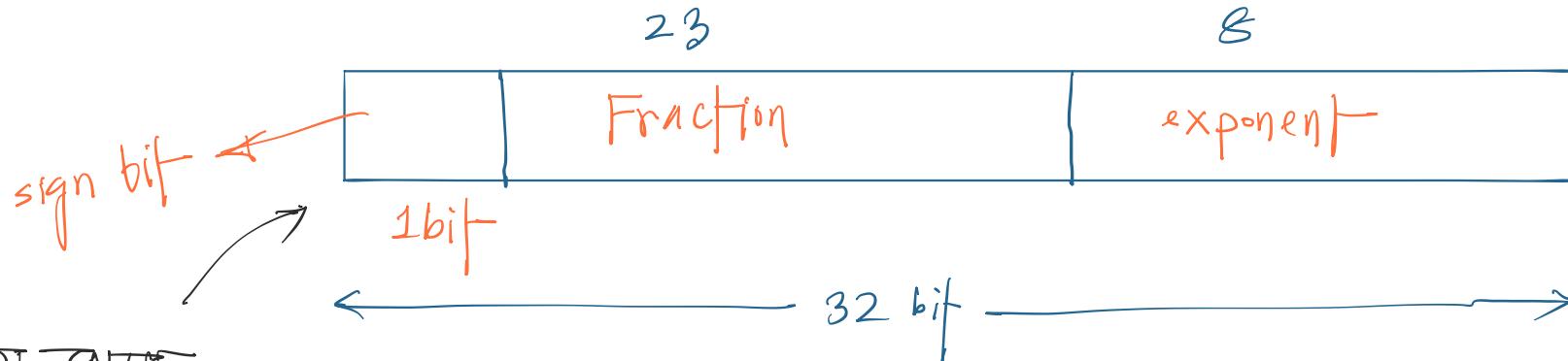
fraction, exponent এ কষ্টিক্ষ দিবা fixed করা যায়?

↪ trade off এরা নাগে। (Good design needs good compromises)

exponent বাড়ালো \rightarrow অনেক বড় একটি number রাখতে পারবে 2^{10}
 $= 1024$
fraction " \rightarrow precision বাড়বে

{ 8 bit exponent
23 bit fraction }

represent কোনো



পুরো বাজে

extra facilities

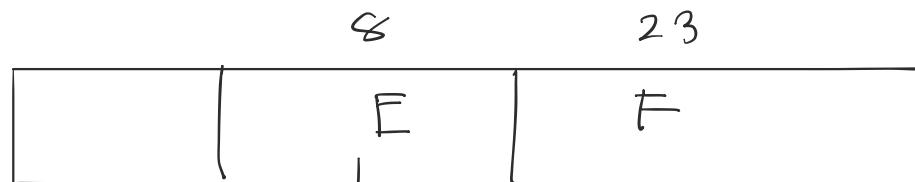
দাখিল

→ fraction মানের রাখলে compare করতে পারবেনা, suppose
fraction -এর but exponent কর্তৃ। But comparison operation
অনেক সহজ। So, comparison faster করতে



comparison → আগে sign bit mismatch check

বিলো টেলে exponent check (high probability যে এখনকার mismatch
করবে)
বিলো " " fraction check



↓

binary to decimal conversion

$$F \times 2^{-1} + F_2 \times 2^{-2}$$

problem এখন — যদি ০ আসে

—আমরা বিহু combination fix করে নির্মাণ

exponent = 0

Fraction full part = 0

(less occurring case)

একটি scenario (0) । আর add

হবে না।

↪ make the common case

first (। add কর্তৃত common),
পরই rare case.

exponent can also be positive or negative.

↳ exponent এবং অগ্রে bit sign bit বাস্থাবো।

-1 → 1 1 1 1 1 1 1

+1 → 0 0 0 0 0 0 1

bit by bit compare করা হাজায় উপায় নেই। তখন -1 বড় আভাবে।

∴ 2's complement করা যাবে না। (exponent comparison এ ঘোট যাবে)

soln: bias calculate :

$$\underbrace{2^{n-1} - 1}$$

$$8\text{ bit এর জন্য} = 2^{8-1} - 1$$

$$= 127$$

127 যাগ করে রাখুলো scale up করবো। negative জন্ম
positive হবে।

$$-1 \rightarrow 127 + (-1) = 126$$

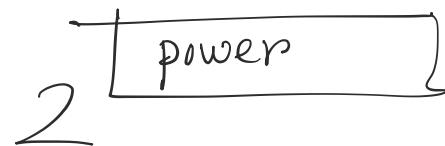
b
binary TO represent

$$1 \rightarrow 128 \xrightarrow{\quad}$$

এখন bit by bit compare.

\therefore এই value সেটা always একই থাকবে।

\therefore exponent এর এই value টি পাবো যেটা $\rightarrow 127$ ইটা power



decimal G
 { precision
 range

কি হবে?

binary 64

precision 23 bit + implicit 1
 $= 24 \text{ bit}$

range $\rightarrow -127 \text{ to } 127$ $-(2^7 - 1) \text{ to } (2^7 - 1)$

decimal 6

$2^{127} = 10^{\square} \rightarrow 38. \dots \text{ (range)}$

$$2^{23} = 10^3$$

6. \curvearrowright precision

float $\rightarrow 4 \text{ byte}$

double $\rightarrow 8 \text{ byte (64 bit)}$

sign	fraction	exponent
1	+ 52	+ 11

sign	fraction	exponent
1	+ 52	+ 11

precision and range both in binary and decimal ?

binary { precision = 52
range = -1023 to 1023

decimal: range $2^{1023} = 10^x$
 $x = 307.95$

precision $2^{52} = 10^x$

$x = 15.65$

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)

IEEE Floating-Point Format

single: 8 bits

double: 11 bits

single: 23 bits

double: 52 bits

S	Exponent	Fraction
---	----------	----------

$$x = (-1)^S \times (1 + \text{Fraction}) \times 2^{(\text{Exponent} - \text{Bias})}$$

- S: sign bit (0 \Rightarrow non-negative, 1 \Rightarrow negative)
- Normalize significand: $1.0 \leq |\text{significand}| < 2.0$
 - Always has a leading pre-binary-point 1 bit, so no need to represent it explicitly (hidden bit)
 - Significand is Fraction with the “1.” restored
- Exponent: excess representation: actual exponent + Bias
 - Ensures exponent is unsigned
 - Single: Bias = 127; Double: Bias = 1203

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
 - $\pm 1.0 \times 2^{-126} \approx \pm 1.2 \times 10^{-38}$
range
- Largest value
 - exponent: 11111110
 \Rightarrow actual exponent = $254 - 127 = +127$
 - Fraction: 111...11 \Rightarrow significand ≈ 2.0
 - $\pm 2.0 \times 2^{+127} \approx \pm 3.4 \times 10^{+38}$



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)

Double-Precision Range

- Exponents 0000...00 and 1111...11 reserved
- Smallest value
 - Exponent: 00000000001
 \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: 11111111110
 \Rightarrow actual exponent = $2046 - 1023 = +1023$
 - Fraction: 111...11 \Rightarrow significand ≈ 2.0
 - $\pm 2.0 \times 2^{+1023} \approx \pm 1.8 \times 10^{+308}$

$$\frac{1 + \text{fraction}}{0}$$

Floating-Point Precision

■ Relative precision

- all fraction bits are significant
- Single: approx 2^{-23}
 - A float has 23 bits of mantissa, and 2^{23} is 8,388,608. 23 bits let you store all 6 digit numbers or lower, and most of the 7 digit numbers. This means that floating point numbers have between 6 and 7 digits of precision, regardless of exponent.
 - Equivalent to $23 \times \log_{10}2 \approx 23 \times 0.3 \approx 6$ decimal digits of precision
- Double: approx 2^{-52}
 - Equivalent to $52 \times \log_{10}2 \approx 52 \times 0.3 \approx 16$ decimal digits of precision

$$-(0.75)_{10}$$

• এখন কাজ বাইনারীতে করওয়া। (2 bit floating)

→ এটকন করবে ?

— এট bit number তোধৰ

$$= -(0.11)_2$$

• ২য় কাজ normalized form করওয়া।

$$= -1.1 \times 2^{-1} \rightarrow \text{bias term add}$$

$$-1 + 127 = 126$$

$$= -1.1 \times 2^{126}$$

→ বাইনারীত represent কৰ
লাগবে,

exp : 01111110

1011111101...

Fraction: 10000...

MSB থেকে শুরু হয়।

#

$$1.0 \times 2^0$$

$$\begin{aligned} \exp &= 127 \\ &= 111111 \end{aligned}$$

$$1 \times 2^{127}$$

$$\text{fraction} = 0.000\dots$$

$$\exp = 111111\dots$$

$$S = 0$$

0 এর floating point এ express করতে পারিনা।
implicit 1 আছে।

1 এর করতে পারি। কারণ exponent এ 127 আয়বে।

0 না।

Floating-Point Example

- Represent -0.75
 - $-0.75 = (-1)^1 \times 1.1_2 \times 2^{-1}$
 - $S = 1$
 - Fraction = $1000\dots00_2$
 - Exponent = $-1 + \text{Bias}$
 - Single: $-1 + 127 = 126 = 01111110_2$
 - Double: $-1 + 1023 = 1022 = 01111111110_2$
- Single: $101111101000\dots00$
- Double: $1011111111101000\dots00$

Floating-Point Example

- What number is represented by the single-precision float

11000000101000...00

$$\begin{aligned} (.01)_2 &= 2^{-2} \\ &= (0.25)_{10} \end{aligned}$$

- S = 1
- Fraction = 01000...00₂
- Exponent = 10000001₂ = 129
- $x = (-1)^1 \times (1 + 01_2) \times 2^{(129 - 127)}$
 $= (-1) \times 1.25 \times 2^2$
 $= -5.0$

addition:

$$0.75 \times 10^2$$

$0.004 \times 10^{\textcircled{1}}$ exponent টা হোଇ অবান কରୁ ଲାଗିବେ।
fraction ଦ୍ୱାରା right shift.

$$0.0004 \times 10^2 \quad \text{exponent অବାନ ହେଲେ,}$$

$$\begin{array}{r} 0.7500 \times 10^2 \\ + 0.0004 \times 10^2 \\ \hline 0.7504 \times 10^2 \end{array}$$

$$= 7.504 \times 10^1 \quad \text{normalized form ହିତେ ହେବେ,}$$

binary ଟାଙ୍କି same. But implicit 1 କେ କିମ୍ବା left shift/right shift
କରୁ ଲାଗିବେ, exponent ହୁଏଇ ଅବାନ ଏବତେ,

1. 1

$$\text{right shift} = 0.11$$

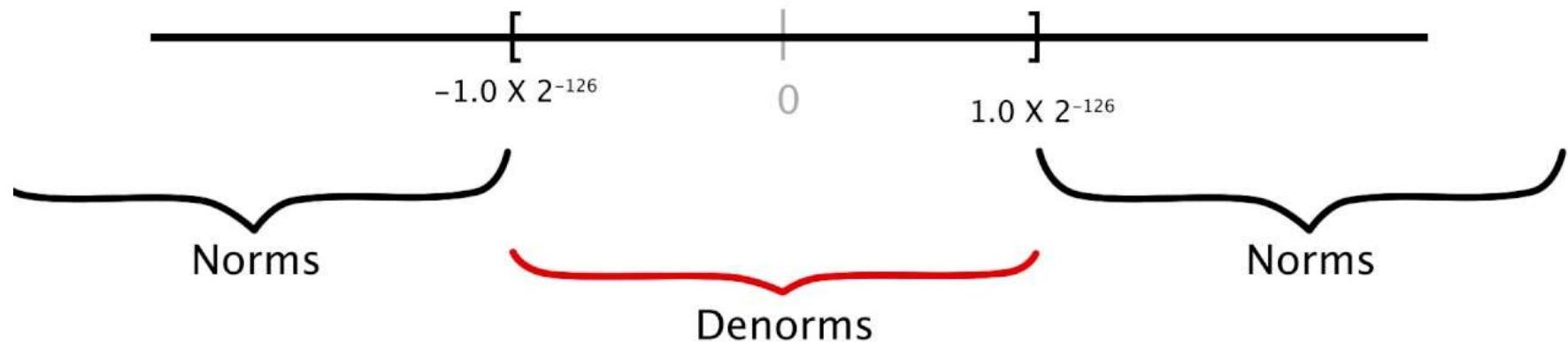
implicit 1 ହିତେ

\therefore implicit value ହେବନ୍ତ 01

Denormal Numbers

Denormalized Numbers $\pm 0.B \times 2^e$

Every denormalized number is smaller in absolute value than every normalized number.



Denormal Numbers

- Exponent = 000...0 ⇒ hidden bit is 0

$$x = (-1)^s \times (0 + \text{Fraction}) \times 2^{-\text{Bias}}$$

- Smaller than normal numbers
 - allow for gradual underflow, with diminishing precision
- Denormal with fraction = 000...0

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

Two representations
of 0.0!

Infinities and NaNs

- Exponent = 111...1, Fraction = 000...0
 - $\pm\infty$
 - Can be used in subsequent calculations, avoiding need for overflow check
- Exponent = 111...1, Fraction \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 Addition

- Consider a 4-digit decimal example
 - $9.999 \times 10^1 + 1.610 \times 10^{-1}$
- 1. Align decimal points
 - Shift number with smaller exponent
 - $9.999 \times 10^1 + 0.016 \times 10^1$
- 2. Add significands
 - $9.999 \times 10^1 + 0.016 \times 10^1 = 10.015 \times 10^1$
- 3. Normalize result & check for over/underflow
 - 1.0015×10^2
- 4. Round and renormalize if necessary
 - 1.002×10^2

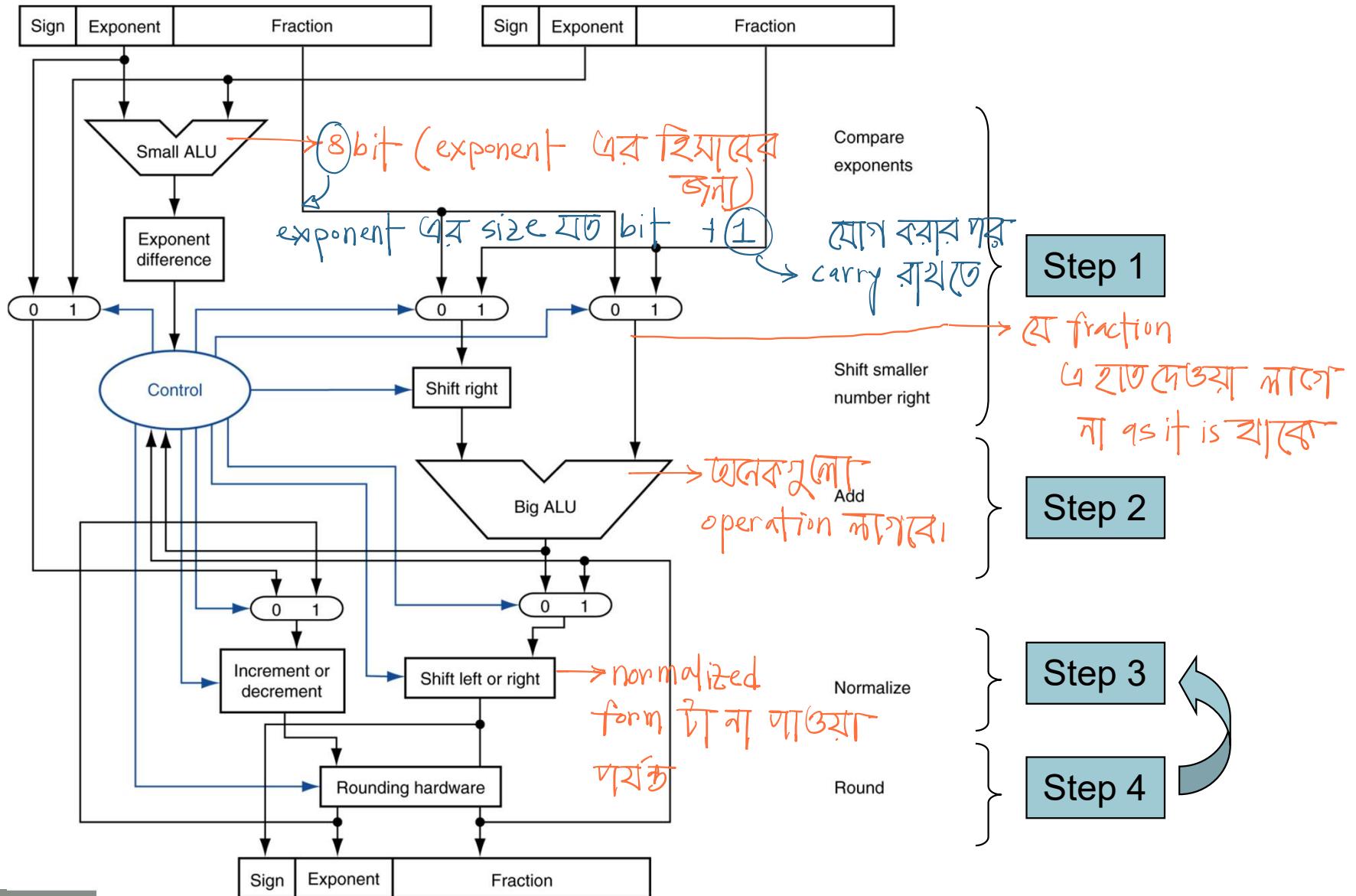
Floating-Point Addition

- Now consider a 4-digit binary example
 - $1.000_2 \times 2^{-1} + -1.110_2 \times 2^{-2}$ ($0.5 + -0.4375$)
- 1. Align binary points
 - Shift number with smaller exponent
 - $1.000_2 \times 2^{-1} + -0.111_2 \times 2^{-1}$
- 2. Add significands
 - $1.000_2 \times 2^{-1} + -0.111_2 \times 2^{-1} = 0.001_2 \times 2^{-1}$
- 3. Normalize result & check for over/underflow
 - $1.000_2 \times 2^{-4}$, with no over/underflow
- 4. Round and renormalize if necessary
 - $1.000_2 \times 2^{-4}$ (no change) = 0.0625

FP 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
- FP adder usually takes several cycles
 - Can be pipelined

FP Adder Hardware



Rounding Bits

- Guard and round digits, and sticky bit
- When computing result, assume there are several extra digits available for shifting and computation. This improves accuracy of computation.
- **Guard digit:** first extra digit/bit to the right of mantissa -- used for rounding addition results
- **Round digit:** second extra digit/bit to the right of mantissa -- used for rounding multiplication results
- **Sticky bit:** third extra digit/bit to the right of mantissa – used for resolving ties such as 0.50...00 vs. 0.50...01

Rounding Bits

- Rounding using Guard and round digits, and sticky bit

G	R	S	Action
0	0	0	Truncate
0	0	1	Truncate
0	1	0	Truncate
0	1	1	Truncate
1	0	0	Round to Even
1	0	1	Round Up
1	1	0	Round Up
1	1	1	Round Up

1.100GRS

1.10001 = 1.53125

1.10010 = 1.56250

1.10011 = 1.59375

Multiplication:

$$\begin{array}{r}
 \text{multiplicand} & 1000 \\
 \text{multiplier} & 1001 \\
 \hline
 & 1000 \\
 & 0000 \\
 & 0000 \\
 & 1000 \\
 \hline
 & 1001000
 \end{array}$$

} Intermediate value
 (store এখন লাগবে memory (র))
 finally add এবার স্মরণযোগ্য memory
 থেকে যান লাগবে 2টি ক্ষেত্রে
 আনতে পারবে।

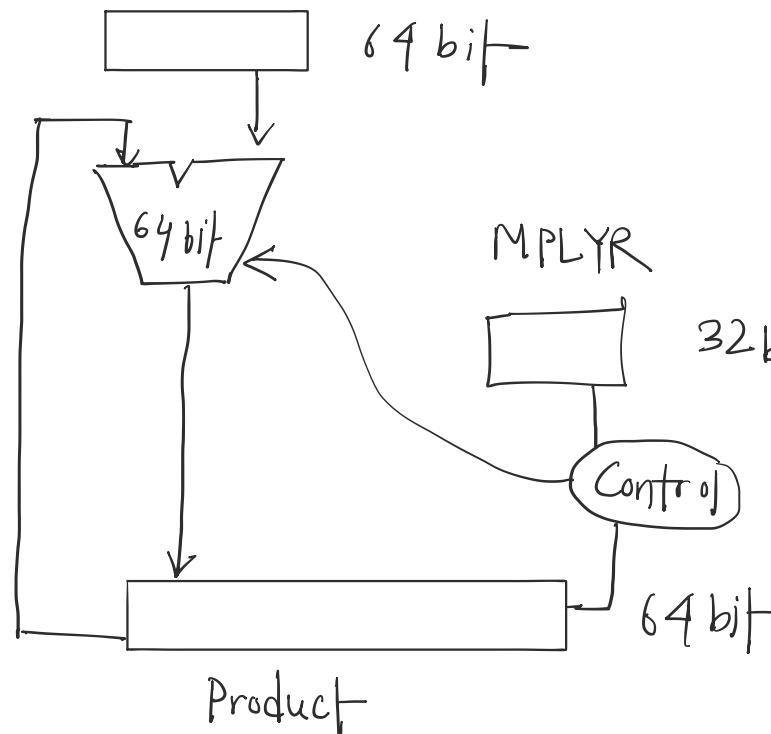
প্রতি memory operation \approx 2টি load, 2টি store

→ memory operation ঘনের costly.

current bit
 multiplier \wedge 1 যদি মুলে add হবে (current multiplicand + product)
 0 " add হবে না।

প্রতি step \wedge multiplier এবার right shift করবে
 multiplicand left shift করব।

MPCND



left shift করছি,
কোণে part বাদ দিতে পারবো না,

32bit : rightmost bit কোনে ঠিকেও সমস্যা
নাই

option 1 multiplier এর right most bit দ্বারে add করবো কিনা product

এ যোগ্য decide \rightarrow implement করে easy . control দ্বারে decide

option 2 ALU কে কাজ করতে দিবো না।

\rightarrow but - যে আগের result ই আছাৰে।

0 0 0 0		1 0 1 1
---------	--	---------

0 0 0 0		0 0 0 0
---------	--	---------

0101

পতিবার add করে, ফিল্ট
multiplier এর value 1
করে addition এর result
store করে।

0 0 0 | 0 1 1 0

left shift

sum করে but 0 থাকায়

0 0 1 0

store করেনা।

right shift

0 0 0 0 | 0 1 1

51

s2

0 0 1 0 1 1 0 0



$$\begin{array}{r}
 00101100 \\
 + 00001011 \\
 \hline
 \text{Sum} = 00110111
 \end{array}$$

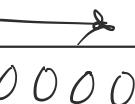
0 0 1 1 0 1 1 1

s3

0 1 0 1 1 0 0 0



0 0 0 0



→ কাজ ক্লিশ না। এবাব

করা নাগরে।

0 0 1 1 0 1 1 1

s4

counter বালৈ

4 steps done.

multiplication

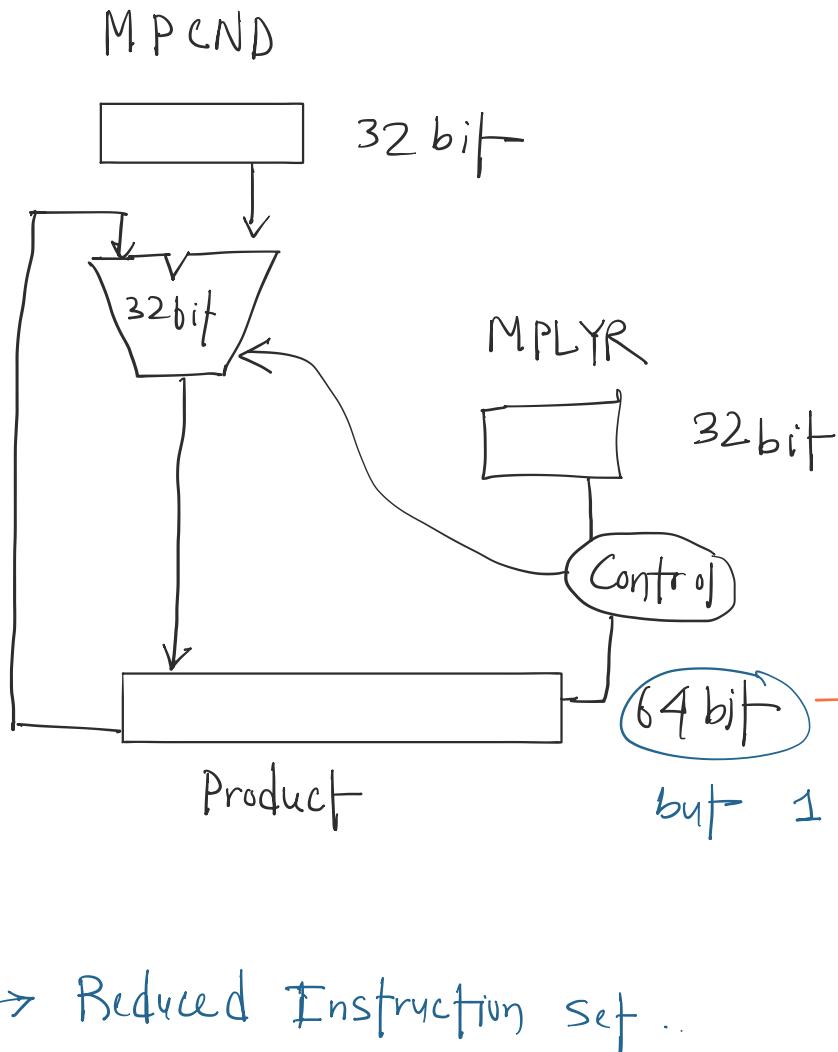
শেষ

1 0 1 1 0 0 0

0 0 0 0

0 0 1 1 0 1 1 1

32 bit processor ടു മാർജ്ജ് കുത്തോടെ 64 bit ALU പാബോനാ।



ഡിസ്ട്രൈബ് കരേ

ലഭിക്കേം, MIPS കു സ്പെഷ്യൽ register, but 1 ടി 64 നാ ഏറ്റു കു 32 bit.

High ← → Low

→ Reduced Instruction set ..

power consumption ← RISC → instruction set reduce (hardware ...) ARM architecture
എൻബൈ കരു

• Apple vision

CISC → ഏഴ് resource ലാ അനേക ഫാക്ട്
complex instruction set .. instruction set അനേക powerful .

RISC



hardware বাড়ায়।

power usage কমায়।

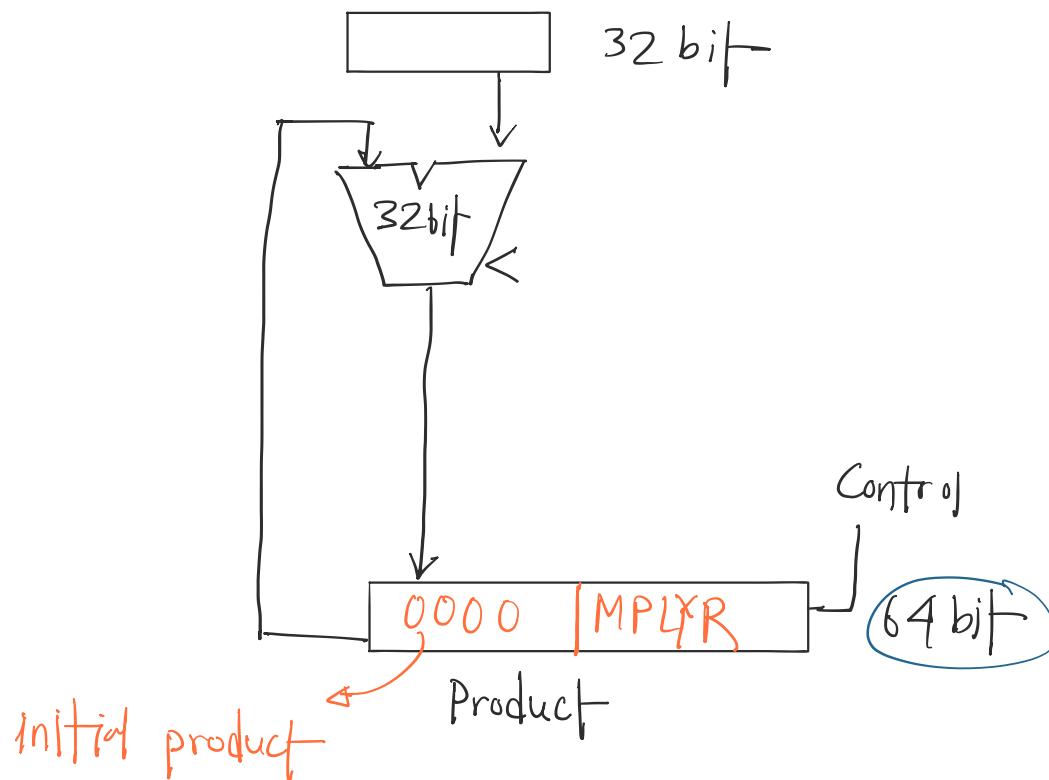
পরিস্থিতি আববে

C/C++ multiplication/division হাববে

- 8 bit হাবতে পারে

→ MPCND left shift করতে পারবেনা। তাইলে 32 bit-exceed করে, so hi product-কে right shift করবে।

MPCND



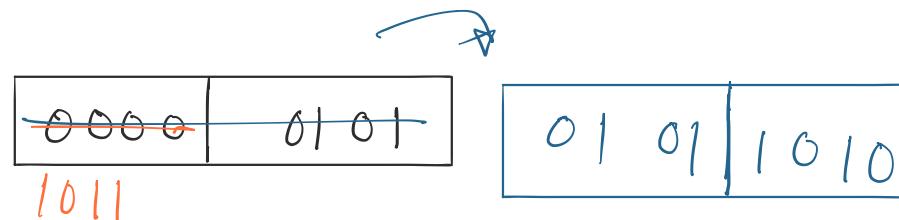
প্রতিবার multiplicand এর সাথে product এর left half add.

↳ update করবে কি না multiplier দ্বারা decide

→ update করবে পুরু লেফ্ট হল্ফ.

1 0 11

(S2)



update

then

right shift

1011

52

0010 | 1101

product ও পিণ্ডেনা and right shift

53

~~1101 | 1101~~ right shift

0110 | 1110

add
update
right shift

1011

54

0011 | 0111

product

যতক্ষণ multiplier
exhaust রহেনা
চলুন করবো।

Numerical example এইটি যথাযাই এই hardware এর মাধ্যমে আচ্ছে হো।
পর্যায় example না।

Floating-Point Multiplication

- Consider a 4-digit 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
 - $1.110 \times 9.200 = 10.212 \Rightarrow 10.212 \times 10^5$
- 3. Normalize result & check for over/underflow
 - 1.0212×10^6
- 4. Round and renormalize if necessary
 - 1.021×10^6
- 5. Determine sign of result from signs of operands
 - $+1.021 \times 10^6$



Floating-Point Multiplication

- Now consider a 4-digit binary example
 - $1.000_2 \times 2^{-1} \times -1.110_2 \times 2^{-2}$ (0.5×-0.4375)
- 1. Add exponents
 - Unbiased: $-1 + -2 = -3$
 - Biased: $(-1 + 127) + (-2 + 127) = -3 + 254 - 127 = -3 + 127$
- 2. Multiply significands
 - $1.000_2 \times 1.110_2 = 1.110_2 \Rightarrow 1.110_2 \times 2^{-3}$
- 3. Normalize result & check for over/underflow
 - $1.110_2 \times 2^{-3}$ (no change) with no over/underflow
- 4. Round and renormalize if necessary
 - $1.110_2 \times 2^{-3}$ (no change)
- 5. Determine sign: $+ve \times -ve \Rightarrow -ve$
 - $-1.110_2 \times 2^{-3} = -0.21875$

FP Arithmetic Hardware

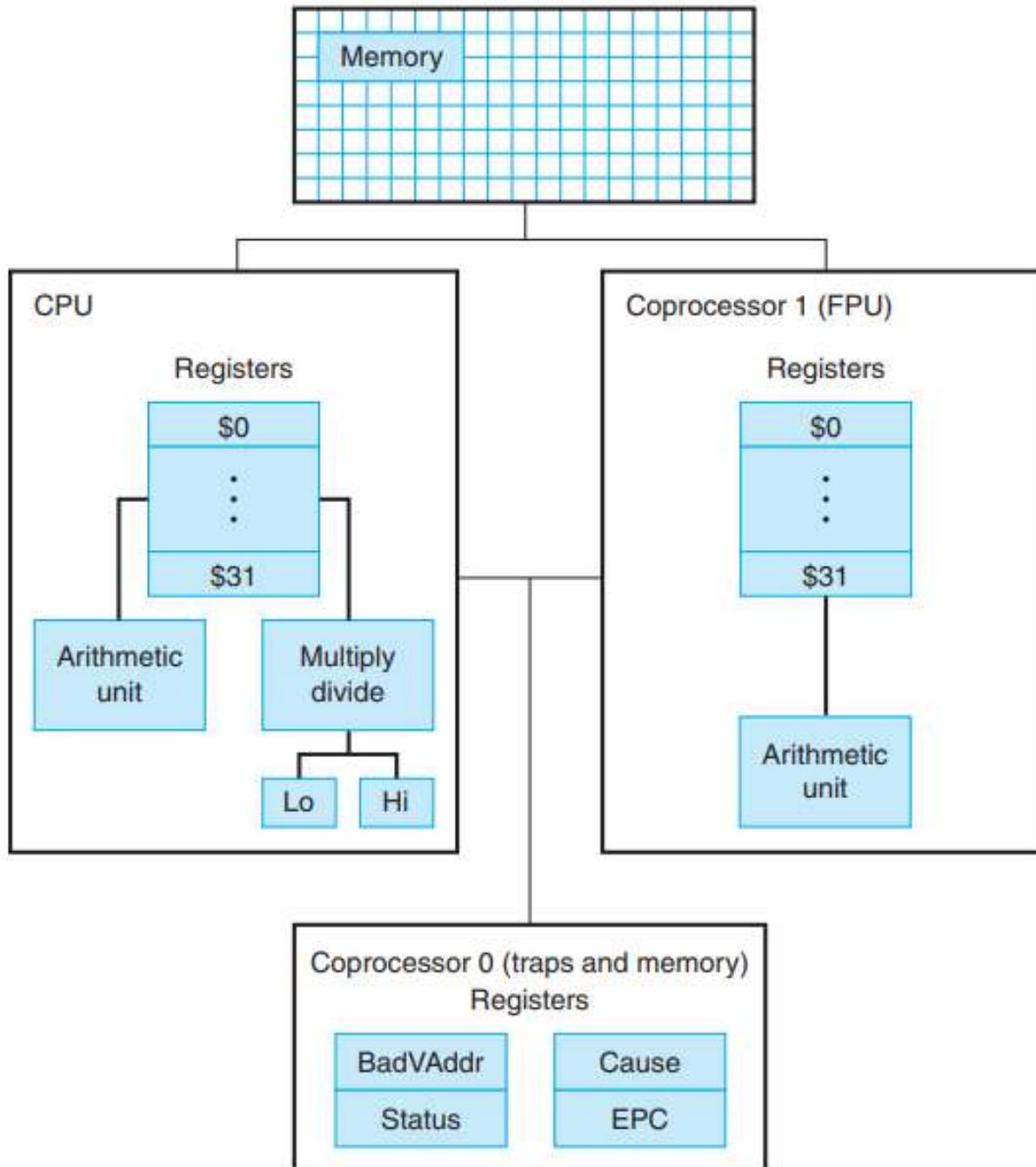
- FP multiplier is of similar complexity to FP adder
 - But uses a multiplier for significands instead of an adder
- FP arithmetic hardware usually does
 - Addition, subtraction, multiplication, division, reciprocal, square-root
 - FP \leftrightarrow integer conversion
- Operations usually takes several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

Associativity of FP Addition

- Associativity: $a + (b + c) = (a + b) + c$
- Is FP addition associative?
 - No

FP Instructions in MIPS

- FP hardware is coprocessor 1
 - Adjunct processor that extends the ISA
- Separate FP registers
 - 32 single-precision: \$f0, \$f1, ... \$f31
 - Paired for double-precision: \$f0/\$f1, \$f2/\$f3, ...
 - Release 2 of MIPS ISA supports 32×64
- FP instructions operate only on FP registers
 - Programs generally don't do integer ops on FP data, or vice versa
 - More registers with minimal code-size impact
- FP load and store instructions
 - lwc1, ldc1, swc1, sdc1
 - e.g., ldc1 \$f8, 32(\$sp)



FP Instructions in MIPS

- Single-precision arithmetic
 - add.s, sub.s, mul.s, div.s
 - e.g., add.s \$f0, \$f1, \$f6
- Double-precision arithmetic
 - add.d, sub.d, mul.d, div.d
 - e.g., mul.d \$f4, \$f4, \$f6
- Single- and double-precision comparison
 - c.xx.s, c.xx.d (xx is eq, lt, le, ...)
 - Sets or clears FP condition-code bit
 - e.g. c.lt.s \$f3, \$f4
- Branch on FP condition code true or false
 - bc1t, bc1f
 - e.g., bc1t TargetLabel

Single and Double Precision operations

- A double precision register is really an even-odd pair of single precision registers, using the even register number as its name.

```
lwc1      $f4,c($sp)  # Load 32-bit F.P. number into F4  
lwc1      $f6,a($sp)  # Load 32-bit F.P. number into F6  
add.s    $f2,$f4,$f6  # F2 = F4 + F6 single precision  
swc1      $f2,b($sp)  # Store 32-bit F.P. number from F2
```

- What if add.d was used?

Floating-Point Instructions in MIPS

- Floating-point *addition*, *single* (add.s) and *addition, double* (add.d)
 - e.g., add.s \$f0, \$f4, \$f6 # $\$f2 = \$f4 + \$f6$
- Floating-point *subtraction*, *single* (sub.s) and *subtraction, double* (sub.d)
 - e.g., sub.d \$f2, \$f4, \$f6 # $\$f2 = \$f4 - \$f6$
- Floating-point *multiplication*, *single* (mul.s) and *multiplication, double* (mul.d)
 - e.g., mul.s \$f2, \$f4, \$f6 # $\$f2 = \$f4 \times \$f6$
- Floating-point *division*, *single* (div.s) and *division, double* (div.d)
 - e.g., div.d \$f2, \$f4, \$f6 # $\$f2 = \$f4 / \$f6$

Floating-Point Instructions in MIPS

- Floating-point *comparison*, *single* (c.x.s) and *comparison*, *double* (c.x.d),
where x may be *equal* (eq), *not equal* (neq), *less than* (lt), *less than or equal* (le), *greater than* (gt), or *greater than or equal* (ge)
 - e.g., c.lt.s \$f2, \$f4 # if (\$f2 < \$f4) cond = 1; else cond = 0
- Floating-point *branch*, *true* (bc1t) and *branch, false* (bc1f)
 - e.g., bc1t 25 # if (cond == 1) go to PC + 4 + 100

Floating-Point Instructions in MIPS

```
if (i==j)  
    f = g + h;  
else  
    f = g - h;
```

```

c.eq.s $f2, $f4
bc1f Else                                #go to Else if i ≠ j
add.s $f6,$f8,$f10                         #f = g + h
j Exit
Else: sub.s $f6,$f8,$f10                  #f = g - h
Exit:

```

MIPS FP Operands

MIPS floating-point operands

Name	Example	Comments
32 floating-point registers	\$f0, \$f1, \$f2, . . . , \$f31	MIPS floating-point registers are used in pairs for double precision numbers.
2^{30} memory words	Memory[0], Memory[4], . . . , Memory[4294967292]	Accessed only by data transfer instructions. MIPS uses byte addresses, so sequential word addresses differ by 4. Memory holds data structures, such as arrays, and spilled registers, such as those saved on procedure calls.

MIPS FP Operations

MIPS floating-point assembly language

Category	Instruction	Example	Meaning	Comments
Arithmetic	FP add single	add.s \$f2,\$f4,\$f6	$\$f2 = \$f4 + \$f6$	FP add (single precision)
	FP subtract single	sub.s \$f2,\$f4,\$f6	$\$f2 = \$f4 - \$f6$	FP sub (single precision)
	FP multiply single	mul.s \$f2,\$f4,\$f6	$\$f2 = \$f4 \times \$f6$	FP multiply (single precision)
	FP divide single	div.s \$f2,\$f4,\$f6	$\$f2 = \$f4 / \$f6$	FP divide (single precision)
	FP add double	add.d \$f2,\$f4,\$f6	$\$f2 = \$f4 + \$f6$	FP add (double precision)
	FP subtract double	sub.d \$f2,\$f4,\$f6	$\$f2 = \$f4 - \$f6$	FP sub (double precision)
	FP multiply double	mul.d \$f2,\$f4,\$f6	$\$f2 = \$f4 \times \$f6$	FP multiply (double precision)
	FP divide double	div.d \$f2,\$f4,\$f6	$\$f2 = \$f4 / \$f6$	FP divide (double precision)
Data transfer	load word copr. 1	lwcl \$f1,100(\$s2)	$\$f1 = \text{Memory}[\$s2 + 100]$	32-bit data to FP register
	store word copr. 1	swcl \$f1,100(\$s2)	$\text{Memory}[\$s2 + 100] = \$f1$	32-bit data to memory
Conditional branch	branch on FP true	bcIt 25	if (cond == 1) go to PC + 4 + 100	PC-relative branch if FP cond.
	branch on FP false	bcIf 25	if (cond == 0) go to PC + 4 + 100	PC-relative branch if not cond.
	FP compare single (eq,ne,lt,le,gt,ge)	c.lt.s \$f2,\$f4	if ($\$f2 < \$f4$) cond = 1; else cond = 0	FP compare less than single precision
	FP compare double (eq,ne,lt,le,gt,ge)	c.lt.d \$f2,\$f4	if ($\$f2 < \$f4$) cond = 1; else cond = 0	FP compare less than double precision

MIPS FP Instruction Format

MIPS floating-point machine language

Name	Format	Example							Comments	
add.s	R	17	16	6	4	2	0		add.s	\$f2,\$f4,\$f6
sub.s	R	17	16	6	4	2	1		sub.s	\$f2,\$f4,\$f6
mul.s	R	17	16	6	4	2	2		mul.s	\$f2,\$f4,\$f6
div.s	R	17	16	6	4	2	3		div.s	\$f2,\$f4,\$f6
add.d	R	17	17	6	4	2	0		add.d	\$f2,\$f4,\$f6
sub.d	R	17	17	6	4	2	1		sub.d	\$f2,\$f4,\$f6
mul.d	R	17	17	6	4	2	2		mul.d	\$f2,\$f4,\$f6
div.d	R	17	17	6	4	2	3		div.d	\$f2,\$f4,\$f6
lwcl	I	49	20	2	100				lwcl	\$f2,100(\$s4)
swcl	I	57	20	2	100				swcl	\$f2,100(\$s4)
bclt	I	17	8	1	25				bclt	25
bclf	I	17	8	0	25				bclf	25
c.lt.s	R	17	16	4	2	0	60		c.lt.s	\$f2,\$f4
c.lt.d	R	17	17	4	2	0	60		c.lt.d	\$f2,\$f4
Field size		6 bits	5 bits	5 bits	5 bits	5 bits	6 bits	All MIPS instructions	32 bits	

FP Example: °F to °C

- C code:

```
float f2c (float fahr) {  
    return ((5.0/9.0)*(fahr - 32.0));  
}
```

- fahr in \$f12, result in \$f0, literals in global memory space

- Compiled MIPS code:

```
f2c:  lwc1  $f16, const5($gp)  
      lwc2  $f18, const9($gp)  
      div.s $f16, $f16, $f18  
      lwc1  $f18, const32($gp)  
      sub.s $f18, $f12, $f18  
      mul.s $f0,  $f16, $f18  
      jr    $ra
```

Accurate Arithmetic

- IEEE Std 754 specifies additional rounding control
 - Extra bits of precision (guard, round, sticky)
 - Choice of rounding modes
 - Allows programmer to fine-tune numerical behavior of a computation
- Not all FP units implement all options
 - Most programming languages and FP libraries just use defaults
- Trade-off between hardware complexity, performance, and market requirements

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
 - Signed and unsigned integers
 - Floating-point approximation to reals
- Bounded range and precision
 - Operations can overflow and underflow
- MIPS ISA
 - Core instructions: 54 most frequently used
 - 100% of SPECINT, 97% of SPECFP
 - Other instructions: less frequent

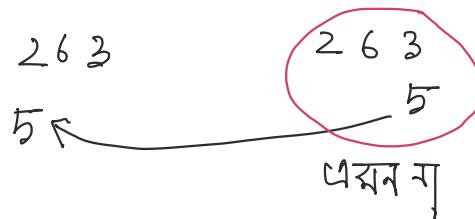
$$\begin{array}{r}
 & 5 | 263 | 52 \\
 & \underline{25} \\
 & \underline{\underline{13}} \\
 & \underline{\underline{10}} \\
 & \underline{3} \rightarrow \text{remainder}
 \end{array}$$

29.10.24

divisor → 5
64bit
dividend → 52
quotient → 13
remainder → 3

subtraction করে left ।

most of the time divisor ছোট হবে।

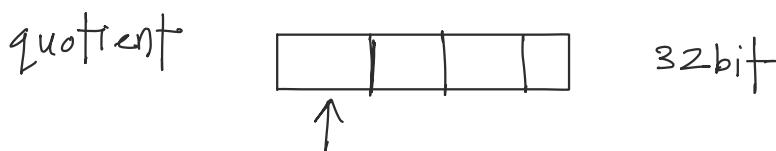


5 কে full shift করা লাগবে।

আমলে full shift নাকরে 1 bit করা shift লাগে।
but এটি আমরা করবে না।

divisor 64 bit

dividend 64 bit। নাহলে subtract করা যাবে না।



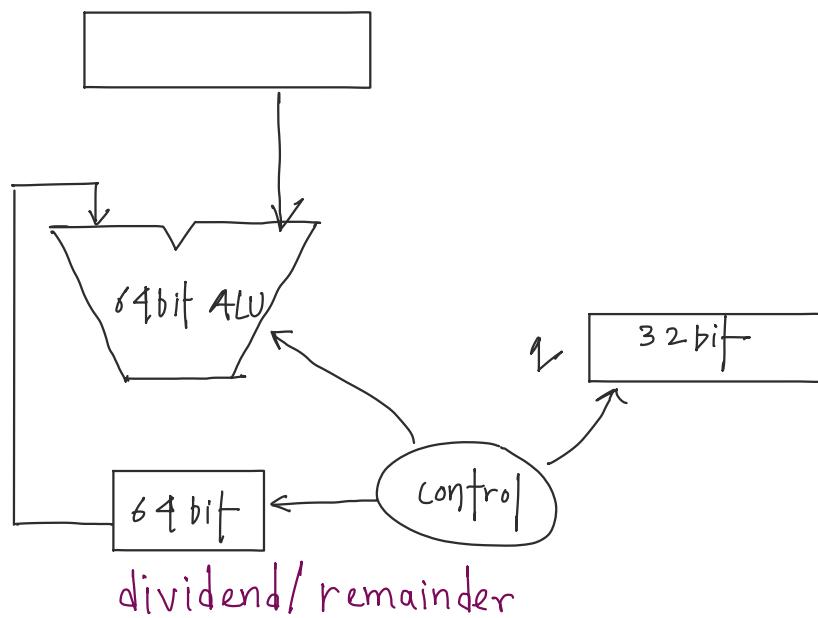
প্রথমটি এখানে আমবে।
এরপর right।

option 1: এখন কোথায় বনবে মেই count বাথার জন্য আলাদা register

option 2: last এ insert then shift.

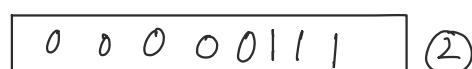
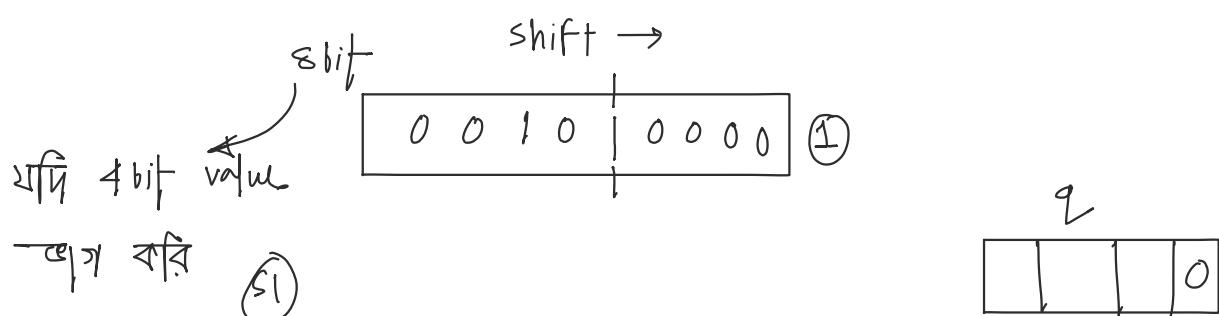
dividend বনে কয়ে reminder হবে। তাই reminder এর জন্য আলাদা register লাগে না।

Divisor 64 bit



way 1 → এটি extra register ব্যবহার করে subtract
এর result থাববে। sign দেওয়ে ...

way 2 → একটি extra operation লাগবে। কোটি এটি follow.



② - ① তাহা

প্রথমে করা যাবে না

সুo quotient 0 0 1

1 ঘর shift

(S2) $\begin{array}{r} 0001 \quad 0000 \\ \rightarrow \\ 0000 \quad 0111 \end{array}$

	1	0	0
--	---	---	---

(S3) $\begin{array}{r} 00001 \quad 000 \\ 00000 \quad 111 \end{array}$

$\begin{array}{r} 00000111 \\ 00000100 \\ \hline 00000011 \end{array}$

(S4) $\begin{array}{r} 000000111 \\ \rightarrow \\ 0000000111 \end{array}$

0	0	1	0	1
---	---	---	---	---

(S5) $\begin{array}{r} 0000000001 \\ 0000000001 \end{array}$

0	1	0	1	1
---	---	---	---	---

right ১ আয়ত

then left shift.

$n+1$ step লাগলো \rightarrow full shift করায় (32 shift লাগে)

n step এ ইটা \rightarrow 1 bit কর শift করা হবে

0011

$n-1$ shift দিয়ে এস্বা

$$\begin{array}{c|c} 0001 & 1000 \\ \hline & \end{array}$$

0011 নিয়ে

shift করা লাগবে 31 shift

algorithm এ লাগে 32 shift.

∴ total 63 shift
 32? 33?

divisor 0 0 0 0 1 |
 dividend 1 0 1 1 0 0 0 0 0 0 0 0 | 0 0 0 0 0 0
 0 0 0 0 0 0 | 1 0 1 1 0 0

(S1)

0 0 0 0 0 1 | 1 0 0 0 0 0



2

1 1 1 0

0 0 0 0 0 0 | 1 0 1 1 0 0

(S2)

0 0 0 0 0 0 | 1 1 0 0 0 0



1 1 0 0

0 0 0 0 0 0 | 1 0 1 1 0 0

(S3)

0 0 0 0 0 0 | 0 1 1 0 0 0

1 0 0 0

0 0 0 0 0 0 | 1 0 1 1 0 0

(S4)

0 0 0 0 0 0 | 0 0 1 1 0 0

0 0 0 1

0 0 0 0 0 0 | 0 1 0 1 0 0 (sub)

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

$$\begin{array}{r}
 010100 \\
 001100 \\
 \hline
 001000
 \end{array}$$

(55)

$$\begin{array}{r}
 000000 | 000110
 \end{array}$$

$$\begin{array}{r}
 00011
 \end{array}$$

$$\begin{array}{r}
 000000 | 001000
 \end{array}$$

(56)

$$\begin{array}{r}
 000000 | 000011
 \end{array}$$

$$\begin{array}{r}
 000111
 \end{array}$$

$$\begin{array}{r}
 000000 | 000010
 \end{array}$$

(57)

$$\begin{array}{r}
 000000 | 000011
 \end{array}$$

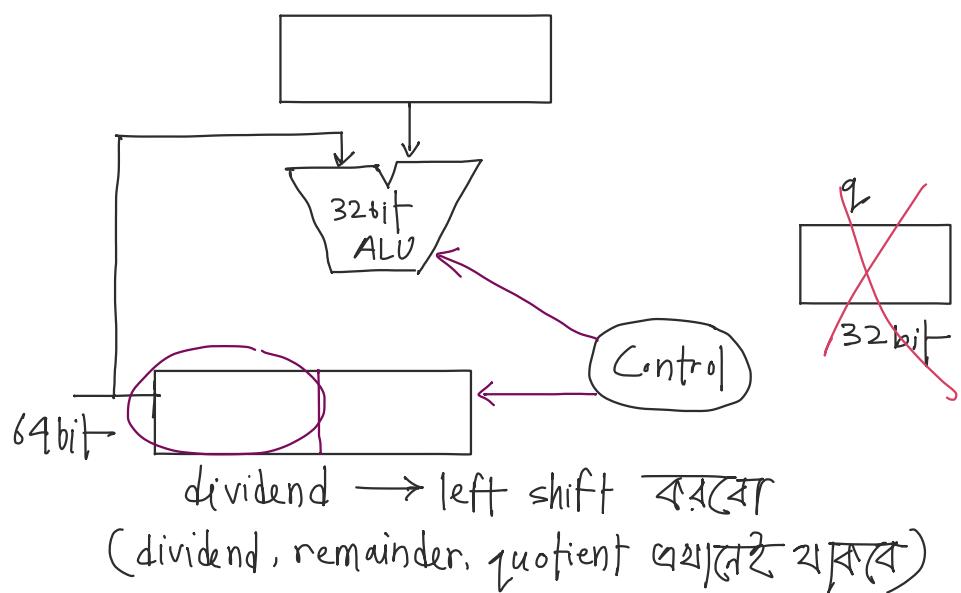
$$\begin{array}{r}
 000000 | \underline{000010}
 \end{array}$$

$$\begin{array}{r}
 00110 \\
 \hline
 \text{Quotient}
 \end{array}$$

remainder

এই hardware real life আ feasible না কারণ 64 bit ALU.

Divisor 32



quotient এ যে value টা ধারণ করা থি shift এ যেই value
দিয়ে fill করবে, তাহলে quotient এর register লাগছে না,

0000 11

divisor 000011

dividend 101100

(51)

000000	101100
--------	--------

After sub

000000	101100
--------	--------

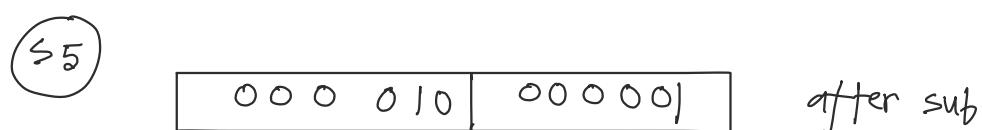
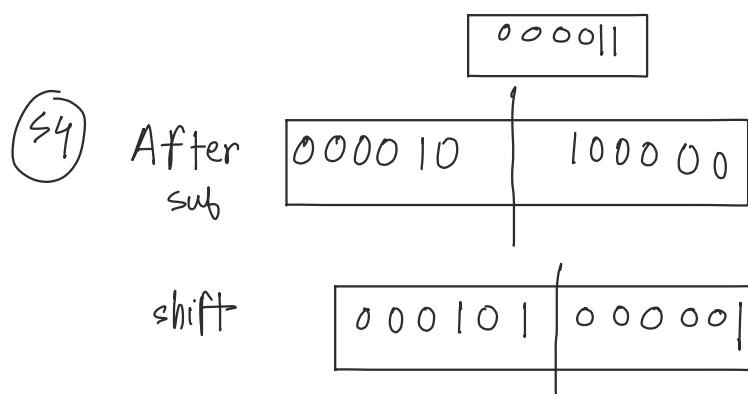
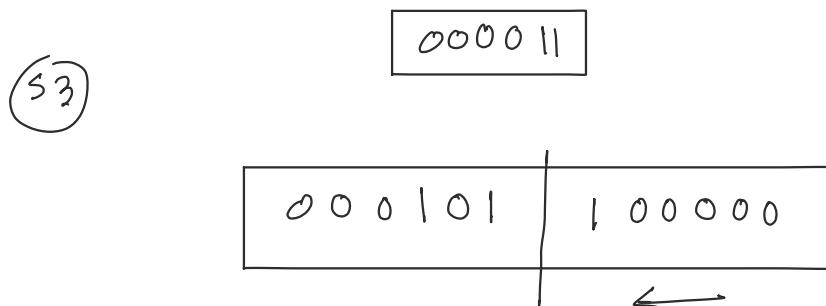
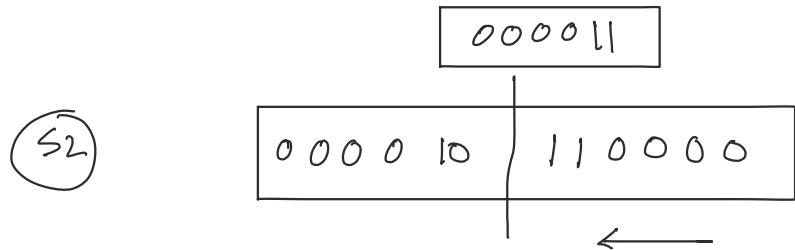
After shift

000001	011000
--------	--------



ultimately

000001	011000
--------	--------



000100		000011
--------	--	--------

sub : 000001 000011

(S6)

000010 000111

after shift

(S7)

000100	001110
--------	--------

remain der portion (extra right shift लगाए तो)
∴ remain der 000010

∴ last step-এ একটি extra কাল বয় লাগবে, right shift
শুধু ধৰ্মী register-এ shift কৰছি High, low থাকবে (register-ও)
সো এটি বয়া possible.

{ xm এ 8 bit এর পূজ্য থাকবে।
full steps- দেখানো লাগবো (after sub, after shift*)
show the calculation কোনো থাকলেও ক্ষব লেখা লাগবে,

multiplication এর addition-আলাদা করে দেখানো লাগবেনা,
divisor এ subtraction.