



# Ch.3 Arithmetic for Computers

## ▼ 3.2 덧셈, 뺄셈

### ▼ Arithmetic for computers

- operations on integers → 사칙연산, overflow
- floating-point real numbers → 표현방법, 연산

#### • Number Representations

- 32-bit signed number(2의 보수)

0000 0000 0000 0000 0000 0000 0000 0000	$_{two} = 0_{ten}$	
0000 0000 0000 0000 0000 0000 0000 0001	$_{two} = +1_{ten}$	
...		
0111 1111 1111 1111 1111 1111 1111 1110	$_{two} = +2,147,483,646_{ten}$	
0111 1111 1111 1111 1111 1111 1111 1111	$_{two} = +2,147,483,647_{ten}$	<i>maxint</i>
1000 0000 0000 0000 0000 0000 0000 0000	$_{two} = -2,147,483,648_{ten}$	
1000 0000 0000 0000 0000 0000 0000 0001	$_{two} = -2,147,483,647_{ten}$	
...		
1111 1111 1111 1111 1111 1111 1111 1110	$_{two} = -2_{ten}$	
1111 1111 1111 1111 1111 1111 1111 1111	$_{two} = -1_{ten}$	<i>minint</i>

MSB (Most Significant Bit) is circled in red on the left. LSB (Least Significant Bit) is circled in red on the right.

- 2의 보수

① ■ 음수의 2의 보수 → 모든 bit complement 후 add 1

↳ 첫번째 1이 나온 이후 모두 바꾸기

② ■ n bit짜리 수를 n bit보다 더 많은 bit로 바꾸기 (부호확장)

- MIPS : 16 bit를 32 bit로 채움

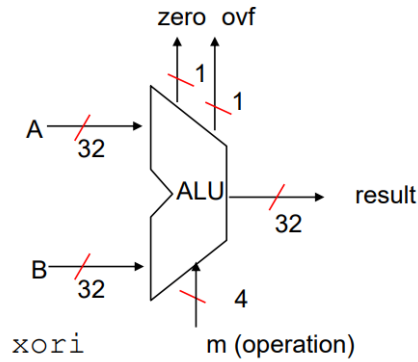
- 부호 비트를 copy해서 확장할 것 → sign extension

◦ 1bu (0으로 확장, unsigned)  $0100 \rightarrow 0000\ 0100$

◦ 1b (1로 확장 → signed)  $1001 \rightarrow 1111\ 1001$

- MIPS Arithmetic Logic Unit → ALU

- ISA 대부분의 연산을 지원



- arithmetic → sign extend.
  - add, addi, addiu, addu
  - sub, subu → overflow detection.
  - mult, multu, div, divu
  - sqrt → overflow(x)
- logic → zero extend.
  - and, andi, nor, or, ori, xor, xori
  - beq, bne, slt, slti, sltiu, sltu
- spcial handling
  - ▶ sign extend – addi, addiu, slti, sltiu
  - ▶ zero extend – andi, ori, xori
  - ▶ overflow detection – add, addi, sub

## ▼ Dealing with overflow

### ① Integer addition

ex) 7 + 6

7 : 0 0 0 1 1 1

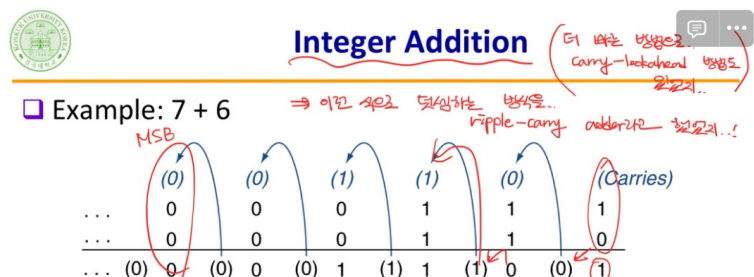
6 : 0 0 0 1 1 0

13 : 0 0 1 1 0 1

⇒ overflow

↳ overflow

- positive + negative → no overflow
- positive + positive = sign bit 0 | 1 → overflow



- negative + negative = sign bit이 0 → overflow

sign :  $0 + 0 = 1 \rightarrow \text{overflow}$   
 sign :  $1 + 1 = 0 \rightarrow \text{overflow}$

## Integer Subtraction

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

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

⇒ overflow

sign :  $0 - 1 = 1 \rightarrow \text{overflow}$

- overflow → 표현가능한 최솟값보다 작거나 최대값보다 크면.

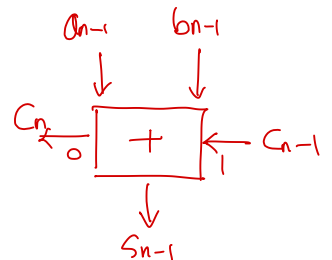
sign :  $1 - 0 = 0 \rightarrow \text{overflow}$

- positive - positive / negative - negative → no overflow
- positive - negative = sign bit이 1 → overflow
- negative - positive = sign bit이 0 → overflow

## Dealing with overflow ⇒ 계산결과가 임의성을 유지해야 함.

- 32 bit 내에서 표현할 수 없을 때 일어남
  - sign bit이 결과의 value bit 포함 and sign bit 불일치
- detecting overflow
  - 2의 보수 덧셈 → 두 수의 부호가 같은데 결과는 다를 때 발생

$$\begin{array}{r} a_{n-1} \ a_{n-2} \ \dots \ a_1 \ a_0 \\ + \ b_{n-1} \ b_{n-2} \ \dots \ b_1 \ b_0 \\ \hline = \ s_{n-1} \ s_{n-2} \ \dots \ s_1 \ s_0 \end{array}$$



이러한 식만으로는 overflow/underflow 알 수 없음.

$$V = a_{n-1} \times b_{n-1} \times s_{n-1} + a_{n-1} \times b_{n-1} \times s_{n-1}$$

$$V = c_n \otimes c_{n-1}$$

- carry in, carry out → MSB

→ carry in, carry out 서로 다를 때 overflow 발생.

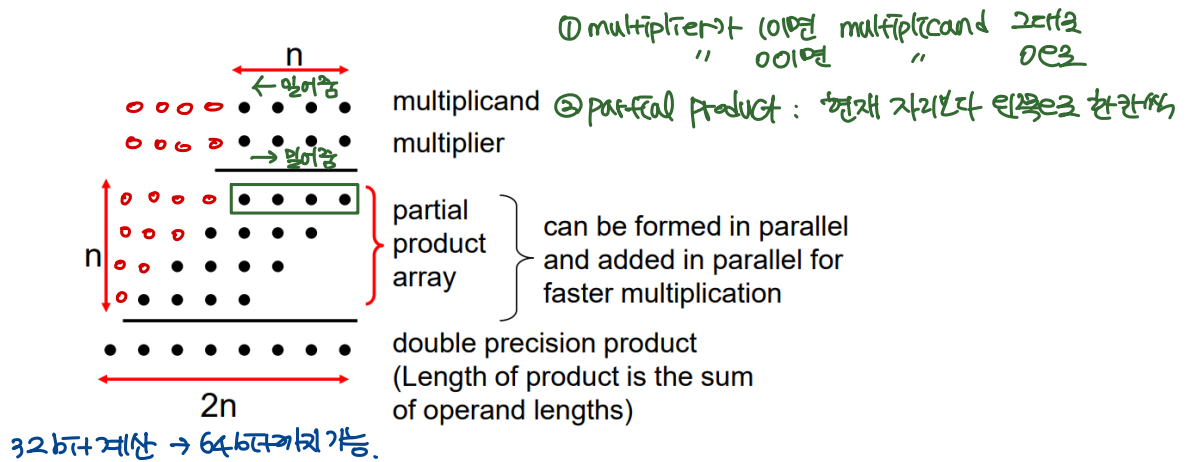
- C : ignore overflow

- Fortran : exception 발생

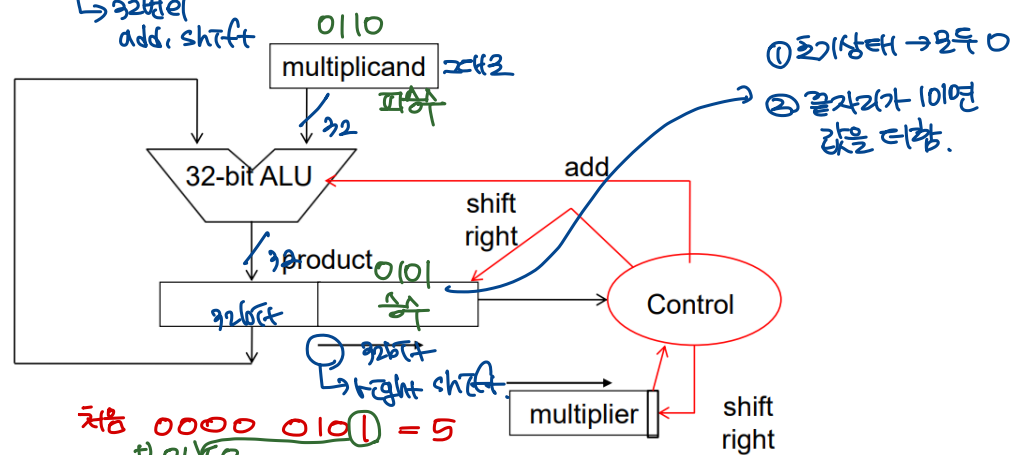
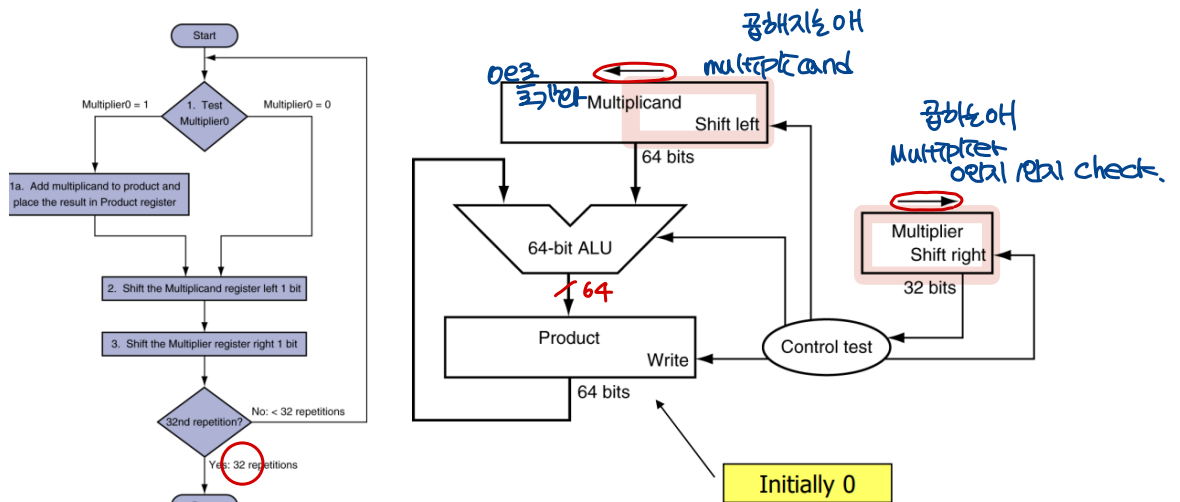
→ exception handler의 유무에 따라서 차가 있음.

## ▼ 3.3 곱셈

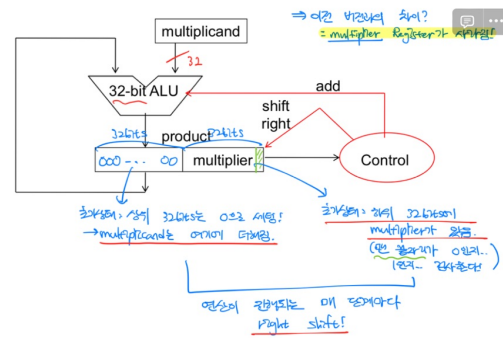
### ▼ Binary multiplication : right shift → add



## ▼ Multiplication hw



$$\begin{array}{r} \text{shift } 0000 \ 0101 = 5 \\ +) 0110 \\ \hline \text{add } 0110 \ 0101 \\ \hline \text{shift} \rightarrow 0011 \ 0010 \\ \hline \text{add } 0011 \ 0010 \\ \hline \text{shift} \rightarrow 0001 \ 1001 \\ +) 0110 \\ \hline \text{add } 0111 \ 1001 \\ \hline \text{shift} \rightarrow 0011 \ 1100 \\ \hline \text{add } 0011 \ 1100 \\ \hline \text{shift} \rightarrow 0001 \ 1110 \end{array}$$



## ▼ MIPS Multply instruction

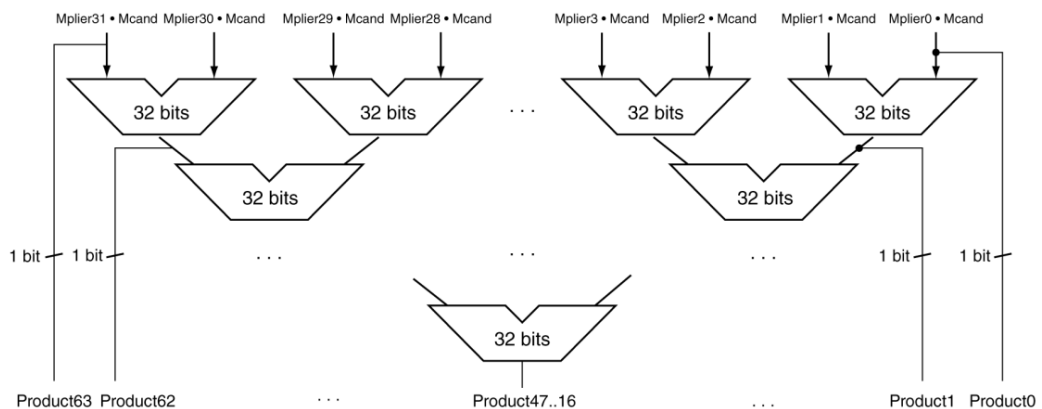
`mult $s0, $s1 # hi||lo = $s0 * $s1`  
 (or multu)

0	16	17	0	0	0x18
---	----	----	---	---	------

- 32-bit registers for product
  - **hi** : MSB → 32bits 상위 32비트 저장
  - **lo** : LSB → 32bits 하위 32비트 저장
- instructions : register에 product 옮기는 명령어
  - **mfhi rd** : move from lo
  - **mflo rd** : mover form hi

## ▼ Fast Multiplier ⇒ adder보다 빠름.

- 32번의 adder를 병렬구조로 시행 →  $\log 32 = 5$  정도로 걸림, 비용도 절약



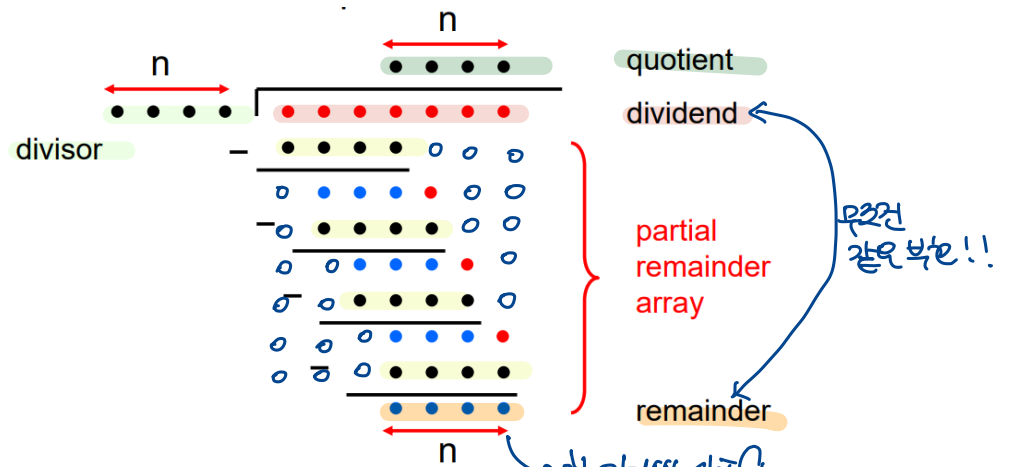
- 파이프라이닝 가능 → 다수의 곱셈도 병렬적으로 수행 가능 ⇒ 더 빨리 가능

## ▼ 3.4 나눗셈

### ▼ Binary Division :

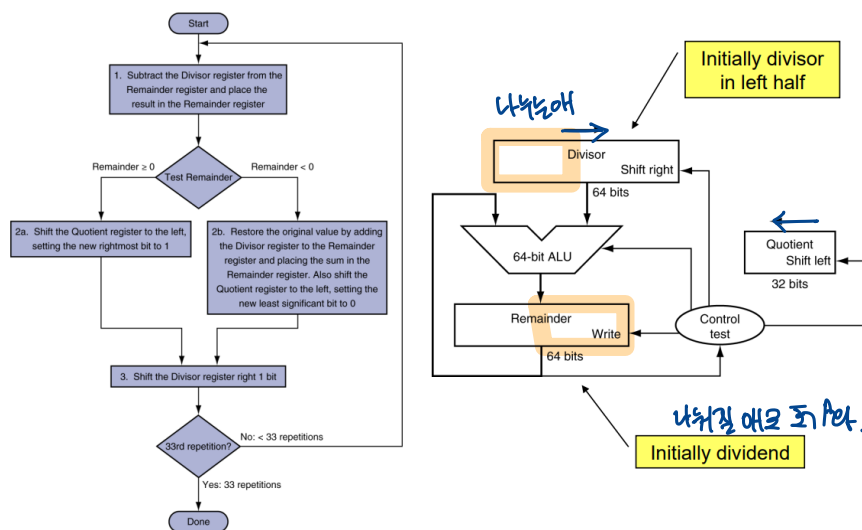
- dividend(나뉘지는 수) = quotient(몫) X divisor(나눌 수) + remainder

연산 반복 ⇒ 몫이되면 복귀해줌.

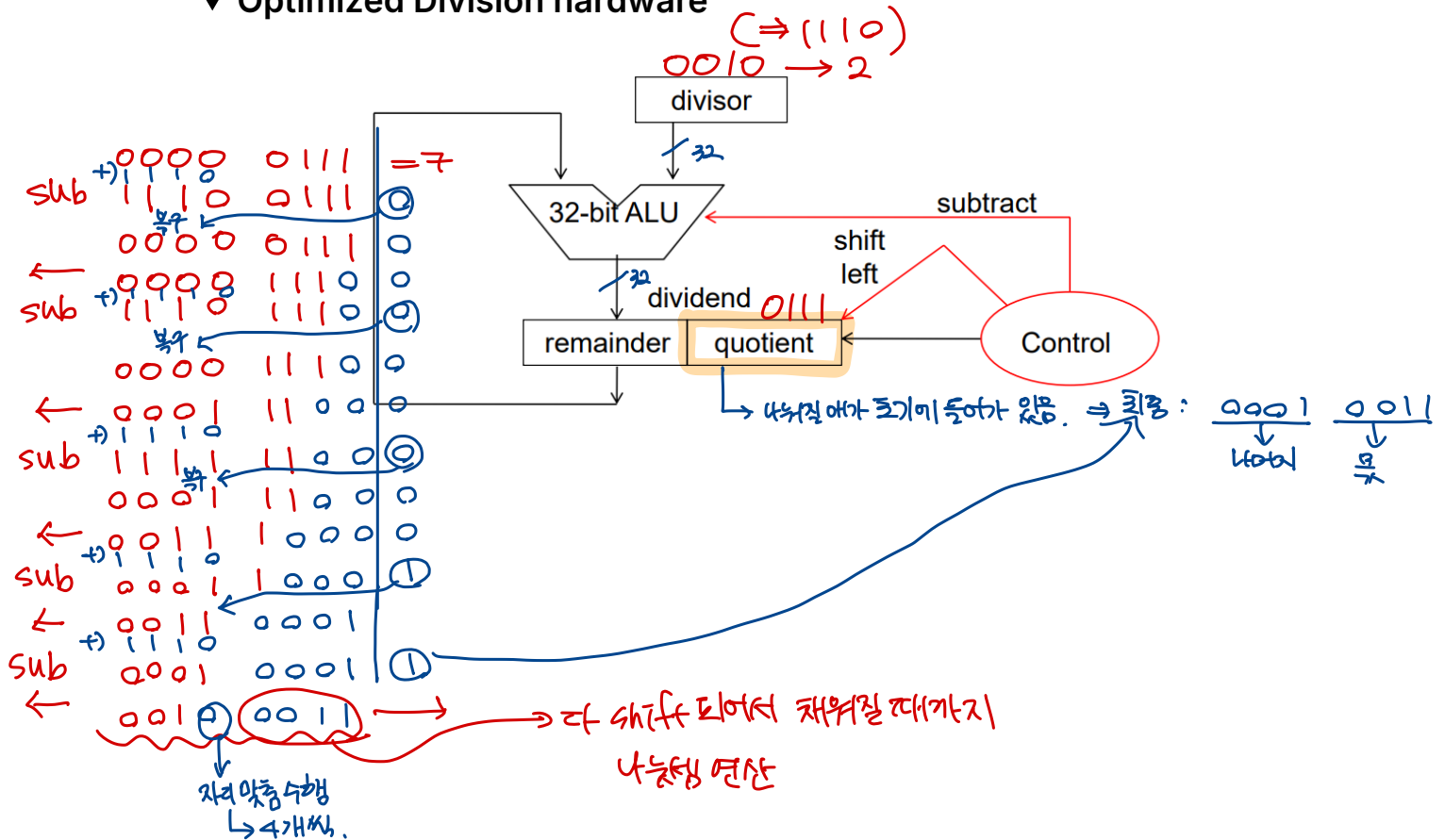


- division approach → 몫이 모두 사라진걸까? 캐워질 때까지 연산!
  - 나누는 애가 0인가?
  - long division
    - 나누는 애 bit ≤ 나눠지는 애의 bit
      - 몫 = 1, 나눠지는 애에서 빼줌
    - otherwise
      - 몫 = 0, 다음 나눠지는 애에서 가져옴
  - restoring division
    - 뺄셈 후 남은 놈이 0보다 작다 → 다시 나누는 애 더해줌  
\* restore.
  - signed division
    - 절대값으로 나눠야 함
    - MIPS : 나눠지는 애랑 나눠지고 있는 남은 애들은 다 같은 부호여야 함

## ▼ Division hardware



## ▼ Optimized Division hardware



## ▼ MIPS Divide Instruction

```
div $s0, $s1 # lo = $s0 / $s1
               # hi = $s0 mod $s1
```

0	16	17	0	0	0x1A
---	----	----	---	---	------

### • mfhi rd, mflo rd

- 결과에 접근 가능 → register 안에 있는 몫, remainder 옮길 수 있음
- 몫이 너무 크면 divide는 overflow 무시함.
- **(SW)**: 0으로 나뉘지는 것을 방지하기 위해 반드시 check함.

## ▼ 3.5 부동 소수점

### ▼ Floating point

- 너무 큰 수, 너무 작은 수 32 bit로 표현 → 부동 소수점 사용
- notation

◦ **normalized** : 소수점 앞에 한 자리만 + 0이 아닌 숫자로 표현

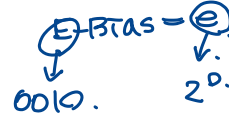
▪ ex)  $(-2)^{34} \times 10^{56}$

▪ 유효 숫자 몇 개까지 있는지 표현 ★

◦ **not normalized** :  $\pm (1.xxxxxx_2 \times 2^{yyyy})$   
frac exp

▪ ex)  $0.23 \times 10^9, 234.3 \times 10^7$

• type : float, double in C



## ▼ Floating point representation

↳ 2진수 표현법

(F) → 유효 숫자

s	E (exponent)	F (fraction)
---	--------------	--------------

⇒ 어떻게 표현하는지에 따라서 표현의 범위가 달라짐.

• fraction(F)와 exponent(E) 크기 사이에서 타협점 찾아야 함

◦ 정밀도/표현 범위 타협

• floating point standard → IEEE, single/double

◦ IEEE 754 FP Standard

	s	E (exponent)	F (fraction)
• single	1 bit	8 bits	23 bits
• double	1 bit	11 bits	52 bits

$$(-1)^{\text{sign}} \times (1 + F) \times 2^{(e - \text{E-bias})}$$

▪ **(F)** normalized format ⇒ hidden bit를 더해줘야 함

$$1.xxxx_2 \times 2^{yyyy}$$

▪ **(E)** excess (biased) notation 표현 ⇒ FP 구분 단순화

• single : Bias = 127 ⇒ -126 ~ 127 사용 ⇒ Smallest Value : 1

• double : Bias = 1023 ⇒ -1022 ~ 1023 사용

• floating point example

$$-0.75 = \frac{-3}{2^2} = -11 \times 2^{-2} = -1.1 \times 2^{-1} \quad (2)$$

$S = 1$  (hidden bit)  
 $F = 1.1000 \dots 0$   
 $\text{Exponent} = e + \text{Bias}$   
 single :  $-1 + 127 = 126$   
 double :  $-1 + 1024 = 1023$

⇒ single	1 bit	8 bits	23 bits	32 bits
Double	1 bit	11 bits	52 bits	64 bits



single

◦ 11000100101000 ... 00  
 $\rightarrow 129$

$S = 1$

$e = E - Bias = 2$

$F = .0100 \dots 0$

$\Rightarrow (-1)^1 \times 1.01 \times 2^2 = (-1)^1 \times 101_{(2)} = -5$

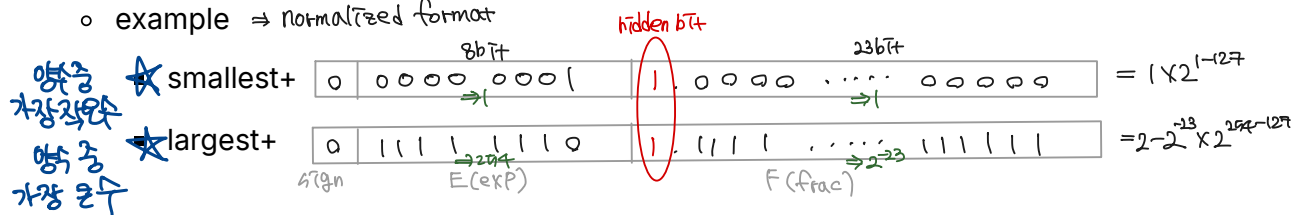
# IEEE 754 FP Normalized Form

◦ Normalized form (single) with a hidden bit

▪ E : 0000 0001 ~ 1111 1110  
 $\begin{matrix} 126 & 127 \end{matrix}$

▪ F : Any

◦ example  $\Rightarrow$  normalized format



ex)  $1.0_2 \times 2^{-1} = 0.1_2 = 0.01111101000000000000000$

$0.75_{10} \times 2^4 = \frac{3}{4} \times 2^4 = 3 \times 2^2 = 11 \times 2^2 = 1.1 \times 2^3 = 0.1000001011000000$

# IEEE 754 FP Standard Encoding

◦ special encodings : unusual event 표현

- 0 : 모든 bit가 0
- 무한대 : 어떤 수를 0으로 나눔
- NAN : 0/0과 같은 invalid operations의 결과

Single Precision		Double Precision		Object Represented
E (8)	F (23)	E (11)	F (52)	
0000 0000	0	0000 ... 0000	0	true zero (0) $\Rightarrow$ 모든 것이 0임.
0000 0000	nonzero	0000 ... 0000	nonzero	$\pm$ denormalized number $(0.x \times x)$ $\rightarrow$ 정수 범위를 초과한 수 (hidden 사용)
0000 0001 to 1111 1110	anything	0000 ... 0001 to 1111 ... 1110	anything	$\pm$ floating point number $\Rightarrow$ normalized
1111 1111	+ 0	1111 ... 1111	- 0	$\pm$ infinity
1111 1111	nonzero	1111 ... 1111	nonzero	not a number (NaN)

## ▼ Denormal Numbers

- Exponent = 000 ... 0  $\rightarrow$  hidden bit is 0  
 $\uparrow$  0인 것임

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

single → -126  
double → -1022

- smaller than normal numbers ⇒ normal 한 숫자보다 더 작은 숫자로 표현
  - precision 점점 작아져서 0으로 다가감 → 0으로 부호 표현 가능
- Fraction = 000... 0

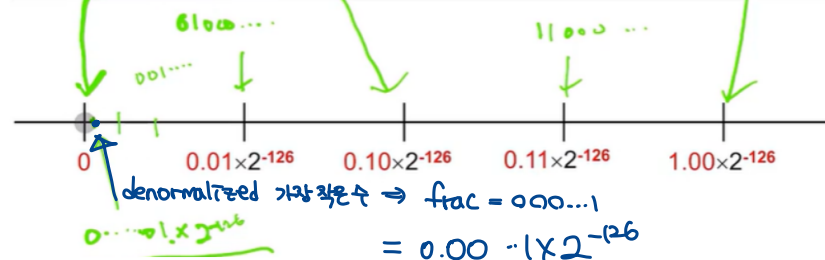
$$x = (-1)^S \times (0 + 0) \times 2^{1-\text{Bias}} = \pm 0.0$$

Two representations of 0.0!

- Denormalization

▶ True zero is the bit string all zero. ★

Single Precision		Double Precision		Object Represented
E (8)	F (23)	E (11)	F (52)	
0000 0001	0000...0000	000000000001	0000...0000	Smallest Normalized No.
0000 0000	Nonzero	000000000000	Nonzero	±Denormalized No.
0000 0000	0000...0000	000000000000	0000...0000	true zero



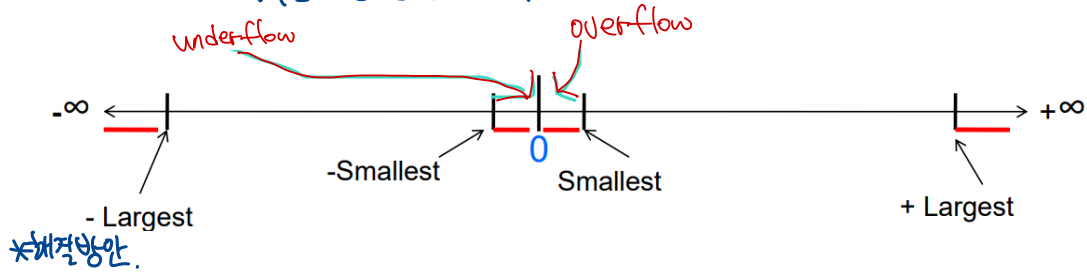
## ▼ Infinities and NaNs

- Exp = 111 ... 1, Frac = 000...0
  - +/- infinity
  - subsequent 계산에 사용 가능 → overflow check 필요 x
- Exp = 111 .... 1, Frac ≠ 0000..00
  - Nan
  - undefiend result
  - subsequent 계산에 사용 가능

## ▼ Exception Events in FP

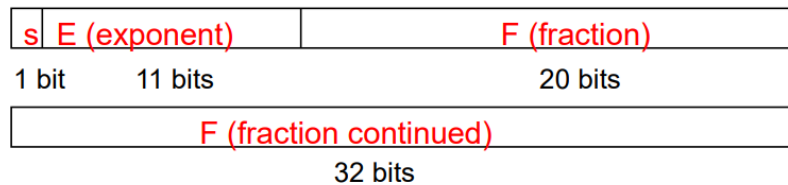
- overflow : 절댓값이 너무 클 때 발생

- underflow : 음수를 표현하는데 너무 작은 수를 표현할 때  $\Rightarrow$  exp의 절대값 너무 커짐!



$\Rightarrow$  다른 format이 더 큰 exp field 가지고 있으면 방지 가능

Double precision – takes two MIPS words



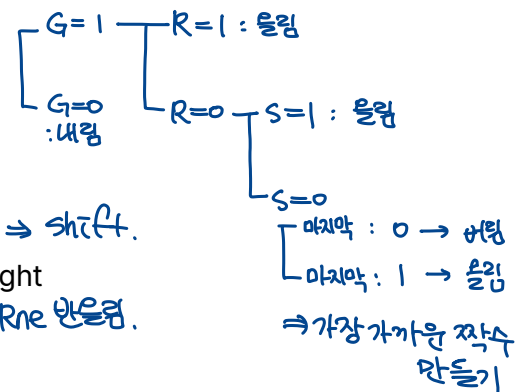
### ▼ Support for accurate arithmetic

- IEEE 754 rounding modes(자리맞춤)
  - 항상 올림    항상 내림    버림    반올림
  - round up, round down, truncate, round to nearest even
- rounding : extra F bit 필요함
  - Guard bit : result normalize  $\rightarrow$  shift left 1bit
  - Round bit : rounding accuracy 향상
  - Sticky bit : round to nearest even(가장 가까운 짝수로의 자리맞춤)
    - 1 bit shift 할 때마다 1로 setting

**F = 1 . XXXXXXXXXXXXXXXXXXXXXXXXXX G R S**

### ▼ Floating point addition

- addition(subtraction)
  - F1, F2  $\rightarrow$  hidden bit restore
  - 둘 중 작은 수의 **E**를 큰 것과 일치할 때까지 shift right  
 $\rightarrow$  큰 것에 맞춰줌!!  $\Rightarrow$  shift.  
 $\rightarrow$  shift  $\rightarrow$  bit 수 늘어나면 Rne 만들림.
  - F1 + F2 = F3
  - F3 normalize  $\Rightarrow$  1, xxx
    - F1, F2 같은 부호 +  $0 < F3 < 4$  : 1 bit right shift F3 + increment E3  
 (overflow) ✱
    - F1, F2 다른 부호 : many left shifts each time decrementing E3  
 (underflow) ✱



→ 표현할 수 없으면 0을 넣어서 반올림

5. Round F3, normalize F3

6. rehide the MSB → 표현 변경

◦ example

$$0.7 = (1.0000 \times 2^{-1}) + (-0.4375 = -1.100 \times 2^{-2})$$

o) add the hidden bit

1)  $-1.100 \rightarrow$  shift right  $\Rightarrow 2^{-1}$ 로 맞추기  
 $\Rightarrow -0.111 \times 2^{-1}$

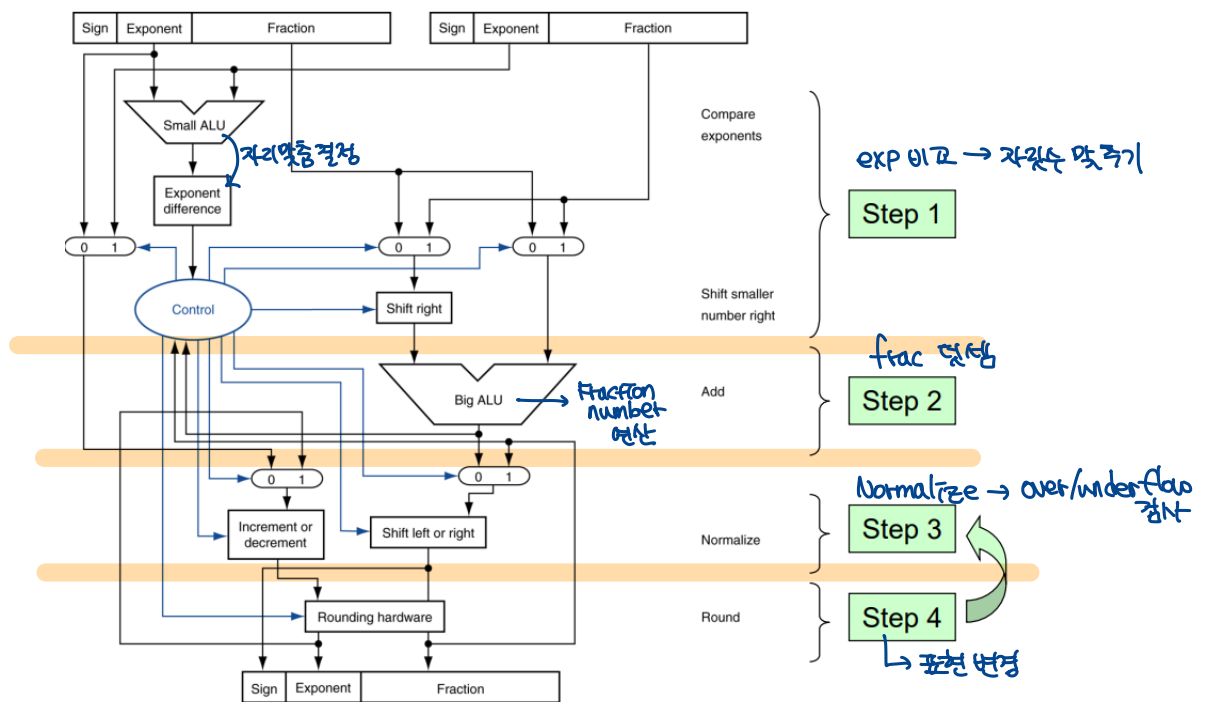
2) add significands :  $1.0000 + (-0.111) = 0.001$

3) normalize the sum, checking exp over/underflow  
 $0.001 \times 2^{-1} = 0.0010 \times 2^{-2} = \dots = 1.000 \times 2^{-4}$

4)  $1.000 \times 2^{-4} \Rightarrow$  shift left 4칸  $\Rightarrow 0.0001000 \Rightarrow 0.0001000$

5) rehide the hidden bit before storing =  $\frac{1}{16}$

• FP adder hw



◦ integer adder보다 훨씬 복잡

◦ 한 clock cycle에서 수행 → integer operations보다 오래 걸림

▪ slower clock : 모든 instructions 병렬적으로 수행

◦ FP adder : 여러 cycle에 걸쳐 수행 → 병렬적으로 처리 가능

→ 할 수 있으면 Integer도

## ▼ Floating point Multiplication

• multiplication

1. F1, F2 → hidden bit restore
2. biased 수를 더한 뒤 bias 빼기 ⇒ **자승제산, 복원결정 (sign bit 같은면 양수)**  

$$127 + e_1 \quad (27 + e_2)$$

$$a. (E1) + E2 - 127 = E3 = e_1 + e_2 + 127$$

3. F1 \* F2 = F3(double precision)

4. F3 normalize ⇒ 1, xxx

a. F1, F2 normalized +  $0 < F3 < 4$  : 1 bit right shift F3 + increment E3 (overflow)

b. overflow, underflow check

5. Round F3, normalize F3 (**반올림**)

6. rehide the MSB

◦ example  $0.5 = (1.0000 \times 2^{-1}) \times (-0.4375 = -1.1100 \times 2^{-2})$

o) add the hidden bit

1) add the two (biased) exponents and subtract the bias from the sum,

$$E_1 = -1, E_2 = -2$$

$$E_3 (\text{biased}) = -1 - 2 + 127 = 124$$

2) multiplication  $F_1 \times F_2 \rightarrow 1.110000 \times 2^{-3}$

3) normalization → multiplication 결과

$$\text{over, underflow} \rightarrow -126 \leq -3 \leq 127$$

4)  $1.110 \times 2^{-3} \Rightarrow$  **자라맞춤**.

5) **복원가 다음** ⇒ **결과가 음수** →  $-1.110 \times 2^{-3}$

$$\Rightarrow \text{10진수로 변환} \quad 1.110 \times 2^{-3} = 0.001110 = \frac{1}{8} + \frac{1}{16} + \frac{1}{32} = -\frac{7}{32}$$

$$\begin{array}{r} 1.000 \\ 1.110 \\ \hline 1.0000 \\ 1.0000 \\ \hline 1.0000 \\ 1.1100 \\ \hline 1.1100 \end{array}$$

GRS  
↓  
내림

#### • FP Arithmetic hw

- FP multiplier : FP adder와 복잡도가 유사
- FP arithmetic hw → 사칙연산, reciprocal(역수), square-root ↔ integer conversion
- operation → 여러 사이클에 걸쳐 병렬적으로 수행됨

### ▼ FP instructions in MIPS

- FP h/w : coprocessor 1 (**ISA 확장**)
- MIPS : 분리된 FP register file 가지고 있음
  - 32 single-precision : \$f0, \$f1, \$f2, ... \$f31
  - paired double-precision : \$f0/\$f1, \$f2/\$f3, ... → **register 두개씩 사용**.
- FP instructions : FP register에서만 operate

- MIPS FP instructions

- load and store

```
lwcl    $f1, 54($s2)    # $f1 = Memory[$s2+54]
swcl    $f1, 58($s4)    # Memory[$s4+58] = $f1
```

- supports IEEE 754

- single

```
add.s   $f2, $f4, $f6    # $f2 = $f4 + $f6
```

- double

```
add.d   $f2, $f4, $f6    # $f2 || $f3 =
                        $f4 || $f5 + $f6 || $f7
```

- similarly for

```
similarly for sub.s, sub.d, mul.s, mul.d, div.s,
div.d
```

- floating point instructions

```
c.x.s   $f2, $f4          # if ($f2 < $f4) cond=1;
                        else cond=0
```

- double precision comparison

```
c.x.d   $f2, $f4          # $f2 || $f3 < $f4 || $f5
                        cond=1; else cond=0
```

- floating point branch

```
bclt    25                # if (cond==1)
                        go to PC+4+25

bclf    25                # if (cond==0)
                        go to PC+4+25
```

## ▼ 3.5 Subword Parallellism

### ▼ subword parallellism

- 그래픽, 오디오 어플리케이션 : 데이터의 벡터에 같은 연산을 반복 수행
  - example : 128-bit adder
    - 8비트 연산자 16개 or 16비트 연산자 8개 or 32bit 연산자 4개 동시 연산 가능
- 데이터 수준 병렬성, 벡터 병렬성, SIMD라고도 함

## ▼ 3.9 오류 및 함정

### ▼ Fallacy : 1bit shift left = 2를 곱해준 것 $\Rightarrow$ 1 bit shift right = 2로 나눈 것

$\Rightarrow$  부호 없는 정수에서만 가능

$\Rightarrow$  부호 있으면 오류 발생

e.g.,  $-5 / 4$

•  $11111011_2 \gg 2 = 11111110_2 = -2$

• Rounds toward  $-\infty$

↓ logical shift 2sb.

c.f.  $11111011_2 \gg 2 = 00111110_2 = +62$

### ▼ Pitfall : floating-point addition $\rightarrow$ 결합법칙 성립 x $(c + (a+b)) = (c+a)+b \Rightarrow$ FP에서는 성립 x

- 병렬프로그램  $\rightarrow$  예상되지 않은 순서로 interleaved operations

결합법칙 성립 x.

### ▼ Fallacy : integer type에서 사용되는 병렬 수행 방식은 FP type에도 적용

- 결합 법칙이 성립하지 않으므로 옳지 않은 가정

### ▼ Pitfalls : MIPS addiu $\rightarrow$ 16비트의 immediate field를 부호확장하여 사용

x sign extends ☆.