

# Lecture 3: Instruction Set Architecture (ISA)

Hunjun Lee  
[hunjunlee@hanyang.ac.kr](mailto:hunjunlee@hanyang.ac.kr)

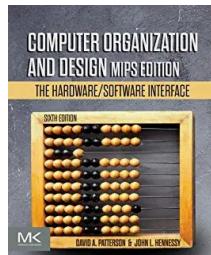
# About Textbook

## ◆ Textbook



5<sup>th</sup> edition  
(OK)

VS.



6<sup>th</sup> edition  
(recommended)

We use this ISA version!

P&H, 6e

// three versions: ARMS, RISC-V & **MIPS**

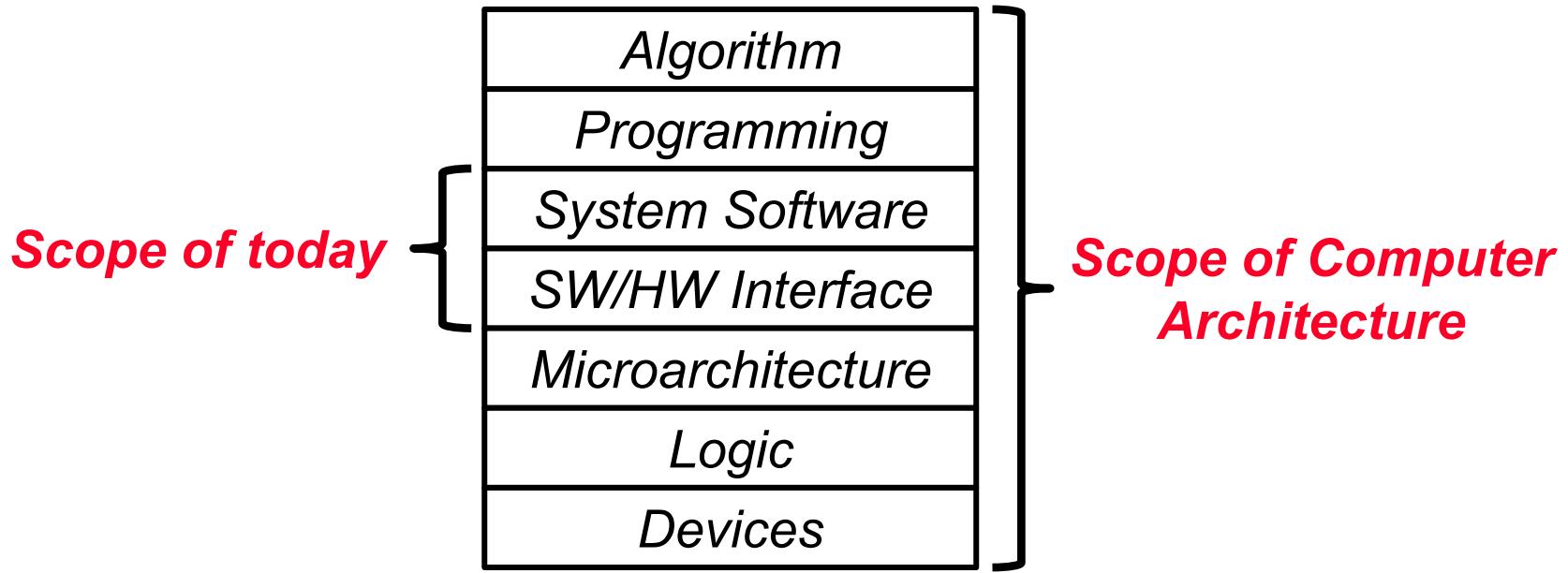
(computer organization & design by Patterson & Hennessy, 6<sup>th</sup> edition)

**MIPS**

**More importantly, “Lecture slides” are  
the most important readings in this class**

# Computing Problem Hierarchy

- ◆ Let's point out how your program run in computer!



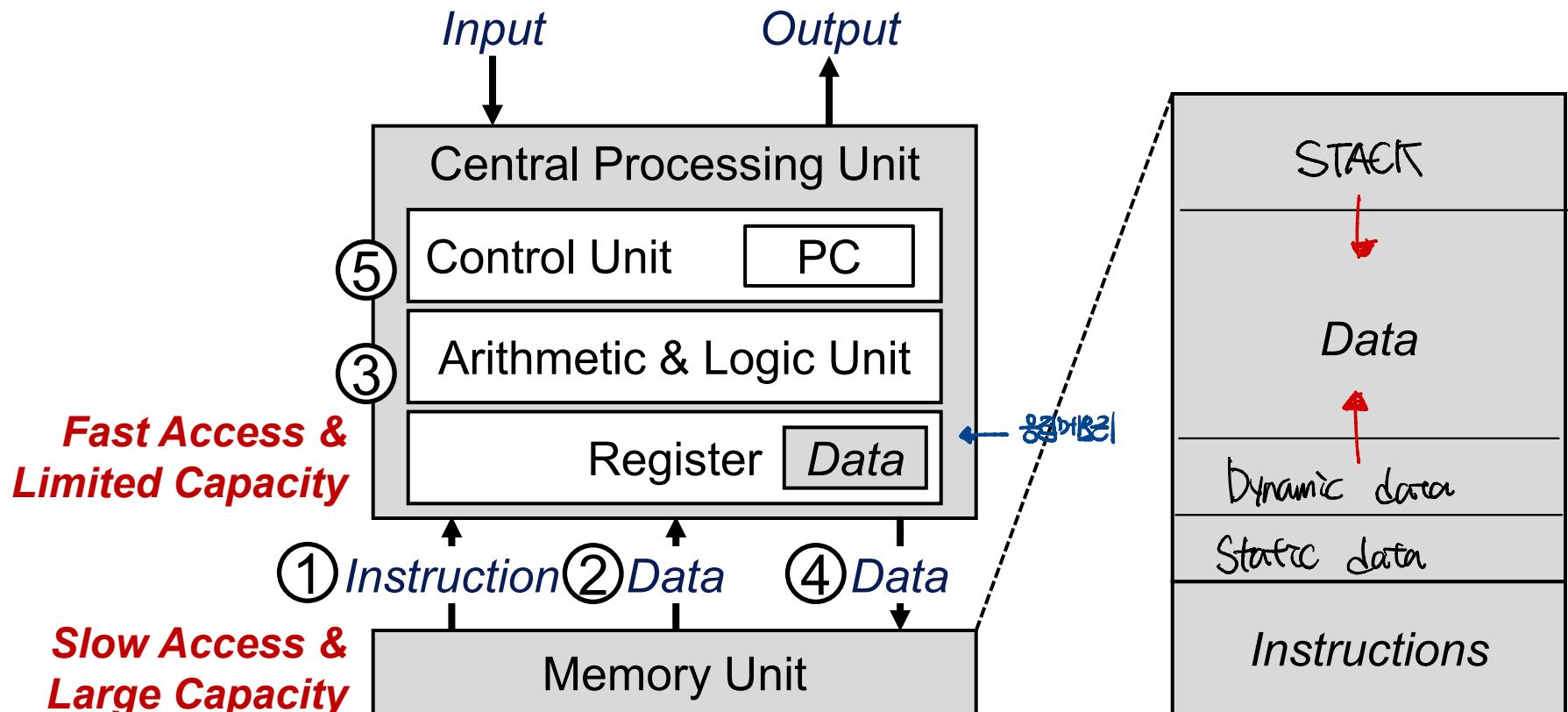
OK, very basics first.

“How does your program run on computer?”

# Basic components of a computer

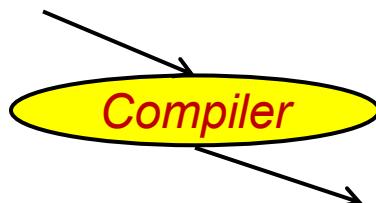
# Von Neumann Architecture

- ◆ Both instructions and data are stored in the memory
- ◆ Instructions dictate (1) which and how data are manipulated and (2) which instruction should be next



# How to load a program to your computer?

C program (\*.c)



Assembly language program (\*.s)

Text →



Object file (\*.o)

Binary →

P&H: Appendix B for further info

Bridges the gap

Library object (\*.o)



Executable binary file (\*.exe)

program-visible  
memory

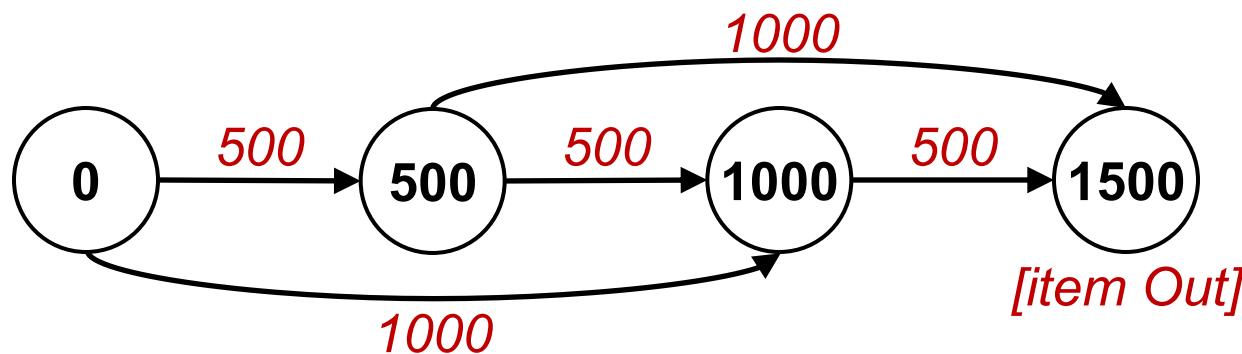


In MEMORY

# Recall: State machine

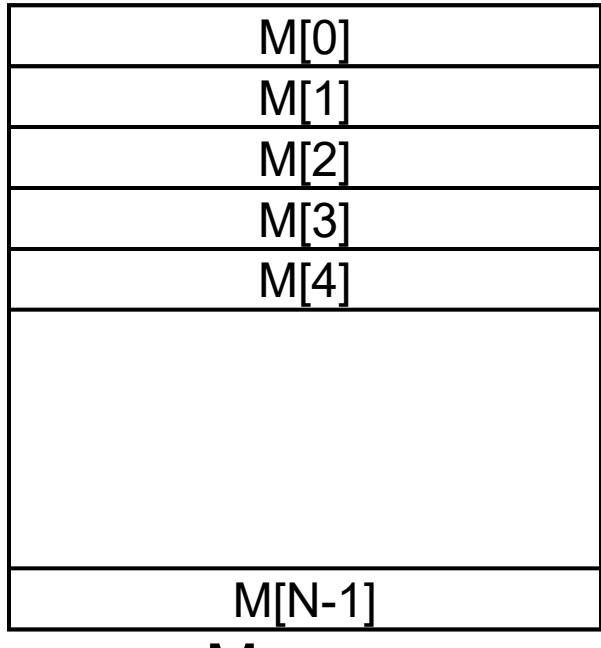
- ◆ Computers are essentially a complex state machine
- ◆ State machine (Ex. vending machine) 자판기
  - Condition: When a user inputs a total of 1500 ₩, the machine outputs a Coke
  - Input type: a user can input 500 ₩ coin and 1000 ₩ bill)
  - States: an amount of cash a user has inputted

↳ 누적 금액

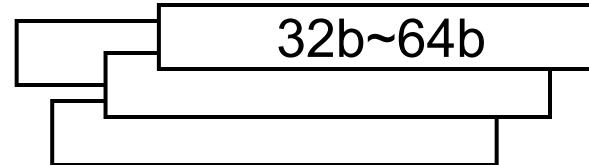


States in Computer  
: Memory, Register, PC +  
Register file

# Programmer visible state (a.k.a. architectural state)



Array of storage locations  
indexed by an address



## Registers

- Given special names in the ISA  
(as opposed to addresses)
- General vs. special purpose

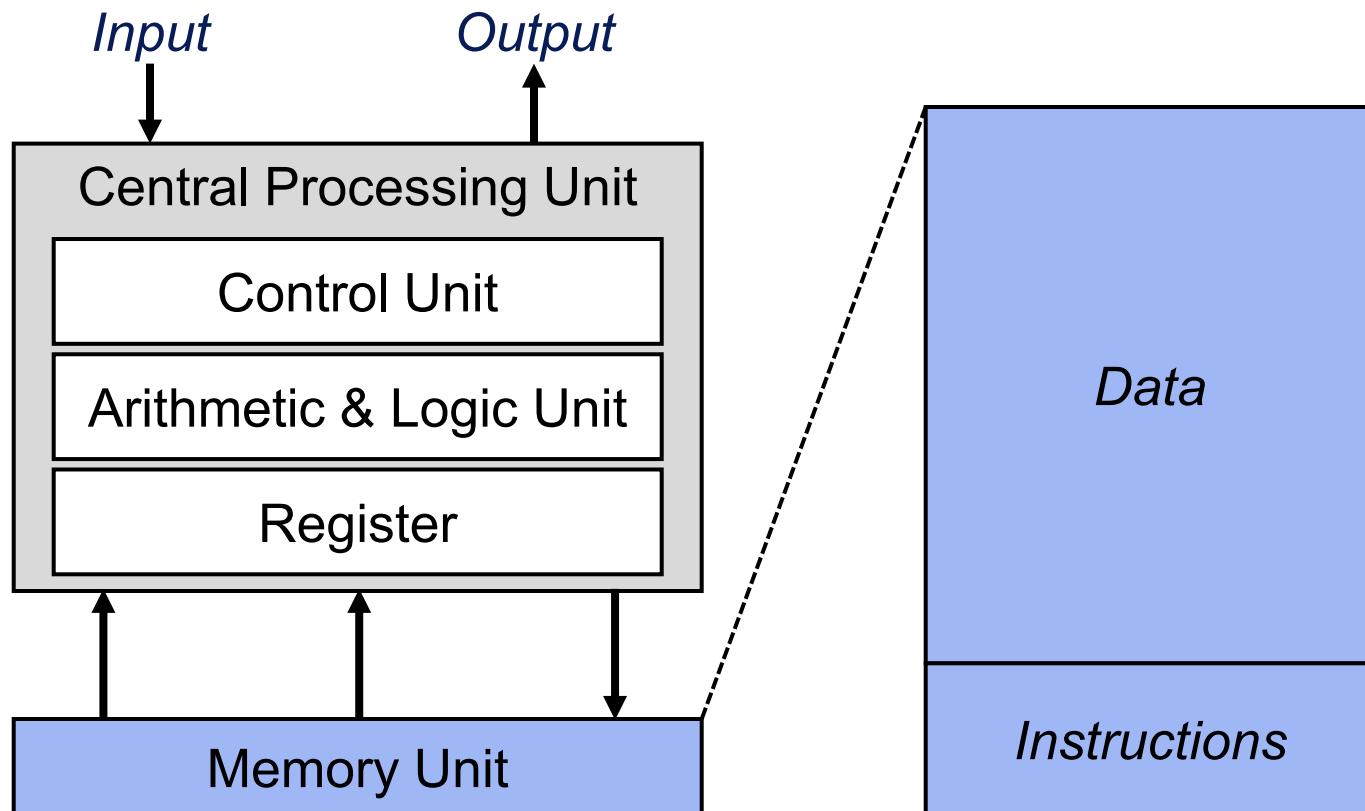
## Program Counter (32b~64b)

Memory address  
of the current instruction

Instructions (and programs) specify how to transform  
the values of programmer visible state

# Von Neumann Architecture

→ ଇନ୍ପୁଟ ଏରାଖିଲୁ କିମ୍ବା



# Runtime storage organization

- ◆ The memory stores the following:
  - (1) program (instructions) and (2) data required in program execution
- ◆ Memory contains bits
  - Logically grouped into bytes (8 bits) and words (e.g., 8, 16, 32 bits)
  - The word size determines the instruction width, register size, ...
- ◆ Address space: Total number of uniquely identifiable locations in memory (differs depending on the architecture)
  - In LC-3, the address space is  $2^{16}$  (16-bit addresses)
  - In MIPS, the address space is  $2^{32}$  (32-bit addresses)
  - In x86-64, the address space is (up to)  $2^{48}$  (48-bit addresses)

# Basics on how the data is stored

# Memory addressing

- ◆ Machines are either byte or word addressable
  - The addressing format determines how the data is aligned
  - A word addressable memory keeps the data word-aligned

*Adopted in memory*

0x00	'H'
0x01	'e'
0x02	'l'
0x03	'l'
0x04	'o'
0x05	'\0'

← 한 바이트

**Byte Addressing**

*Adopted in register file*

0x00	'H'	'e'	'l'	'l'
0x01	'o'	'\0'	PAD	PAD

한 워드 4바이트

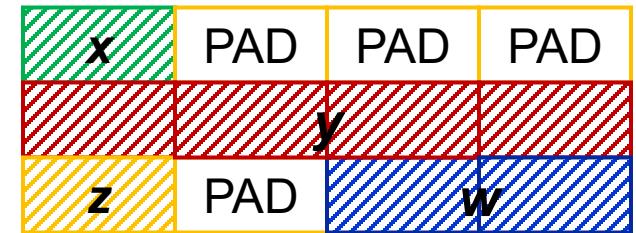
빈 공간 채움

**Word Addressing**

# Memory alignment

- ◆ An address of a variable is aligned based on the size of the variable
  - Char is byte aligned (any addr is fine)
  - Short is halfword aligned (LSB of byte addr must be 0)  
2의 부트 주소의 하드는 0이어야 한다.
  - Int is word aligned (2 LSBs of byte addr must be 0)  
4의 부트 주소의 하드는 00이어야 한다.

```
char x;           // size 1 byte
int y;            // size 4 byte
char z;           // size 1 byte
short w;          // size 2 byte
```



## Memory Alignment in C++?

CPU가 기본적으로 4-byte read하기 때문이다.

## 4 byte Chunk (word 단위)

- Memory bus는 4 byte 기준
  - Cache Line은 한 번에 처리 Byte(32 Byte ~ 64 Byte)를 일괄으로 처리하는 내부 버스 4 byte씩 처리한다.

문맥 Memory가 Align 되어 있지 않다면?

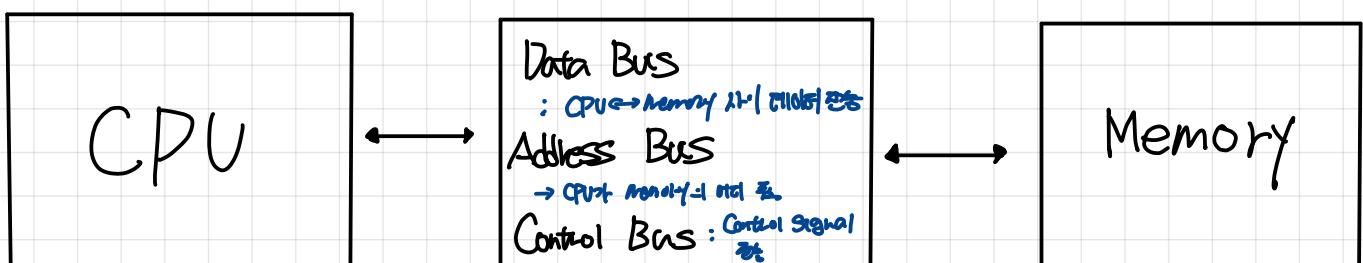
Bus와 Cache가 하나의 데バイ스로 통합되는 Bus, Cache integrated Bus, Cache interface는 전자기  
회로 가능성이 있어 성능이 저하된다.

\* 가장 좋은하게, MIPS architecture는 위 같은 상황을 용납하지 않는다.

→ 대나이션 반응시점

`LW` (`Load`) 도 32-bit를 가지으며, `Alignment` 되어있을 때 성능이 가장 좋다.

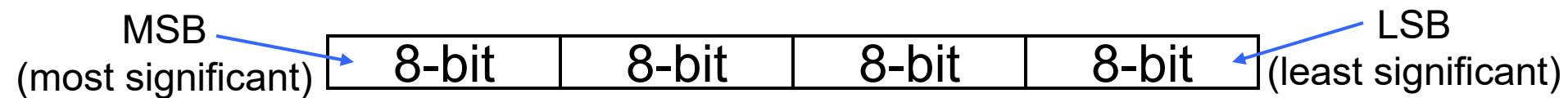
각각) MIPS는  $\text{LW}$ ,  $\text{Lh}$ ,  $\text{Lb}$  를 통해 주소해석 및 접근할 때  
 $2\text{ byte}$   $1\text{ byte}$



# How are the data stored?

## Big Endian vs. Little Endian

- ◆ 32-bit signed or unsigned integer comprises 4 bytes



- ◆ On a byte-addressable machine . . . . .
- When storing 0x12345678

### BigEndian

MSB	Addr	→	LSB
12	34	→	78
byte 7	byte 6	→	byte 4
byte 11	byte 10	→	byte 8
byte 15	byte 14	→	byte 12
byte 19	byte 18	→	byte 16

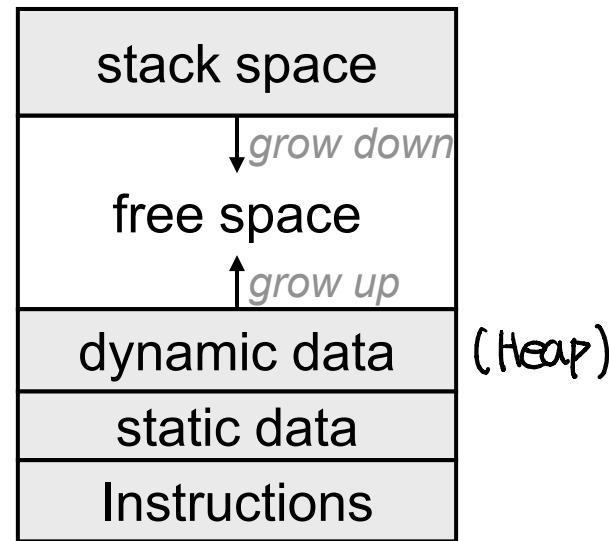
pointer points to the big end

### LittleEndian

LSB	Addr	→	MSB
78	56	→	12
byte 4	byte 5	→	byte 7
byte 8	byte 9	→	byte 11
byte 12	byte 13	→	byte 15
byte 16	byte 17	→	byte 19

pointer points to the little end

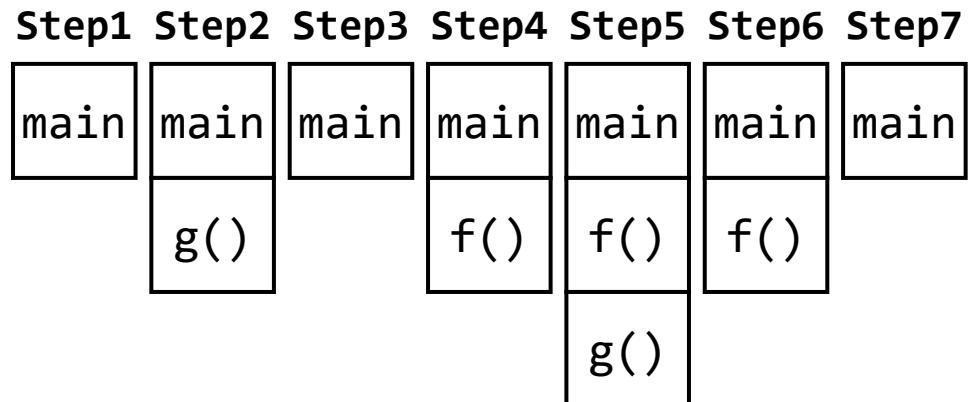
# There are four different parts in the storage!



# How to maintain data for the program? – using stack

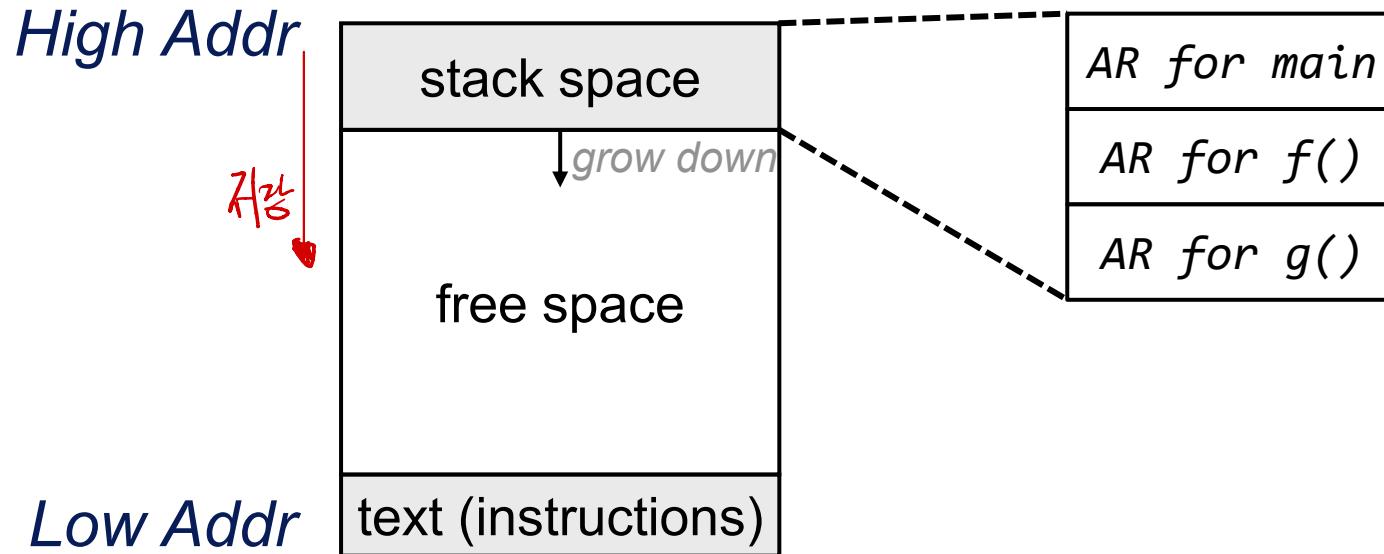
- ◆ We can utilize a “stack” to manage data for the active procedures (or functions) → 記憶域
- ◆ When P calls Q, then Q returns before P returns
  - Lifetimes of procedure activations (required data for the procedure) are properly nested
  - The activation depends on run-time behavior

```
int g() { return 1; }
int f() { return g(); }
void main() {
    g();
    f();
}
```



# Stack management

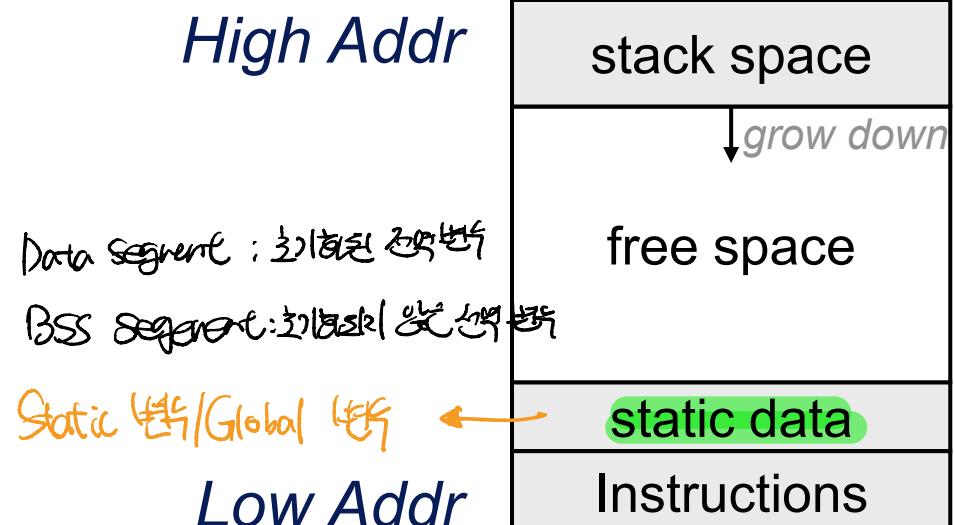
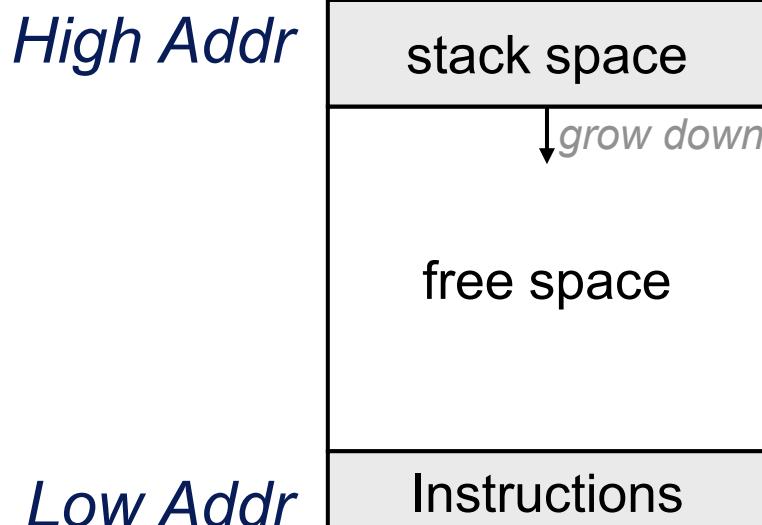
- ◆ Stack data is stored starting from the low address, which grows downwards
- ◆ The information to manage one function call is called activation record (AR) or frame



# Global variables

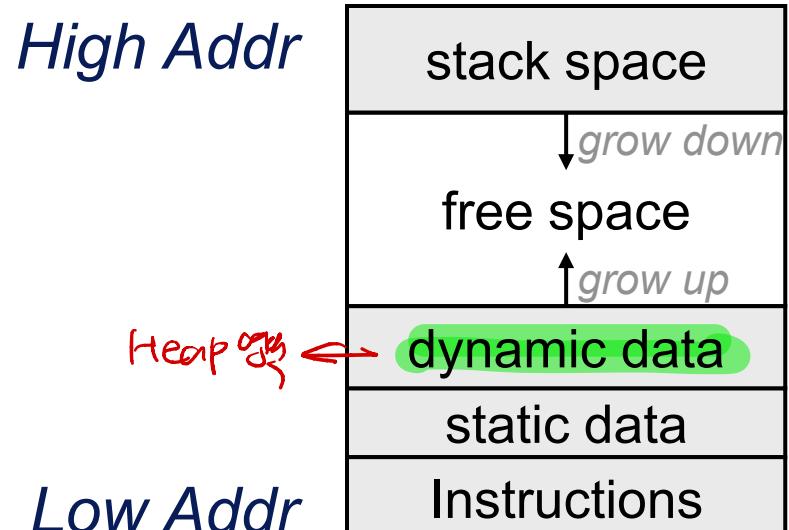
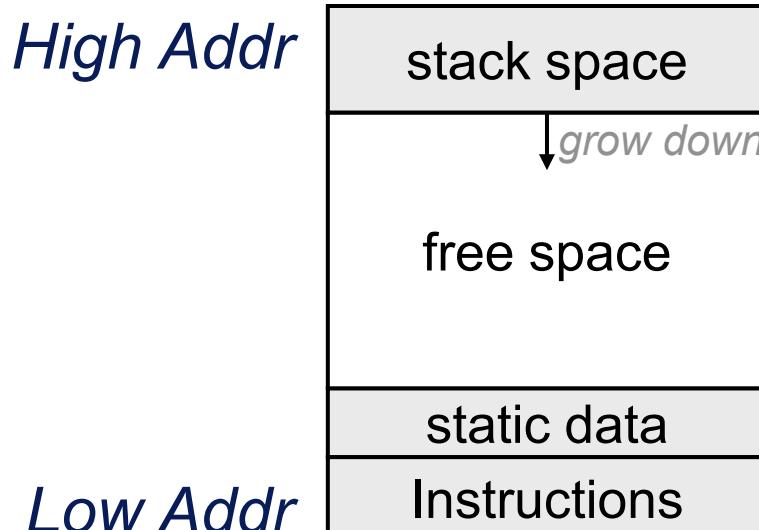
- ◆ All references to a global variable point to the same object
  - It would be impossible (or inefficient) to store a global activation in an activation record
- ◆ Global variables are assigned a fixed address once (statically allocated)

→ Data 영역 (Static data) on Computer memory

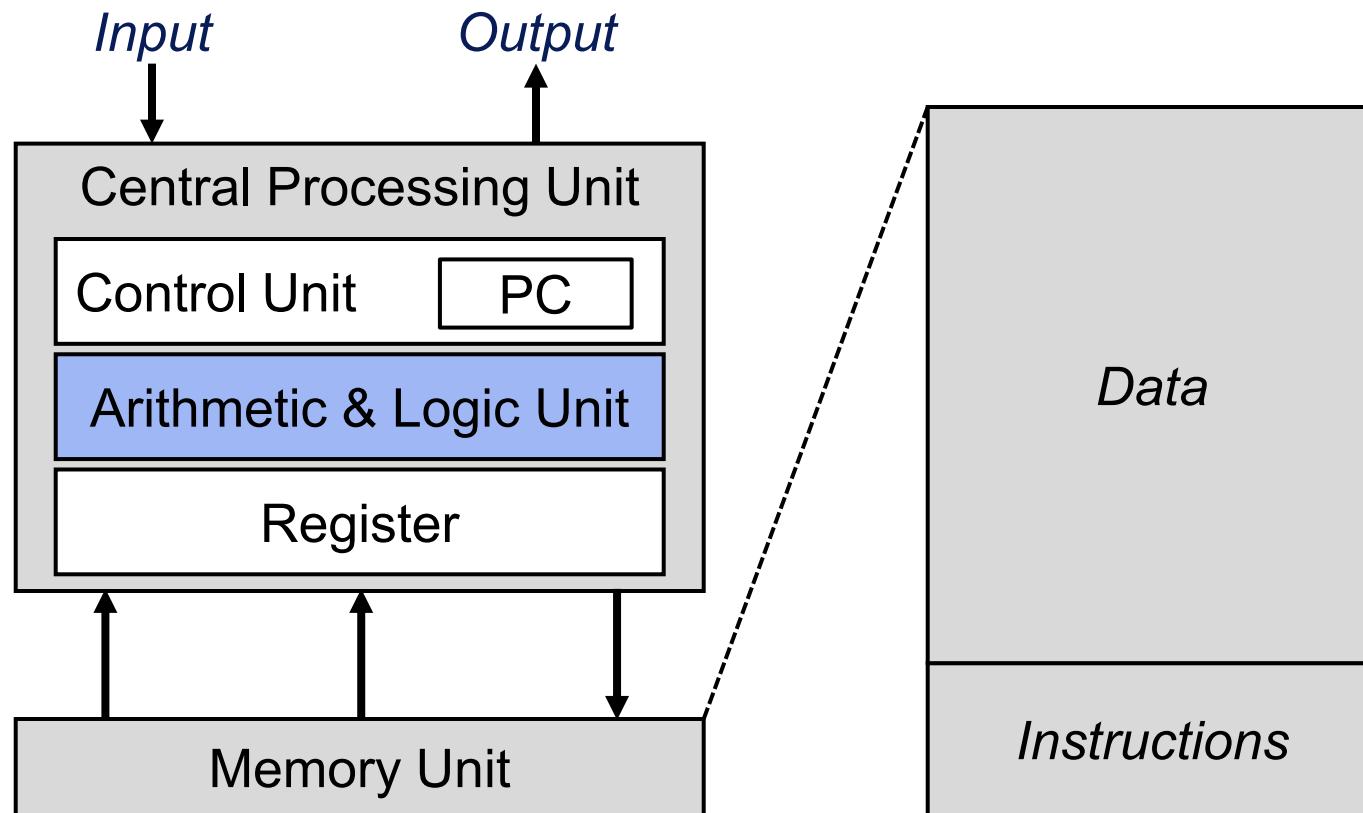


# Dynamic allocation

- ◆ The dynamically allocated value outlives the procedure that creates it (unless deleted beforehand)
- ◆ We rely on **heap** to store the dynamically allocated data



# Von Neumann Architecture



# Processing units

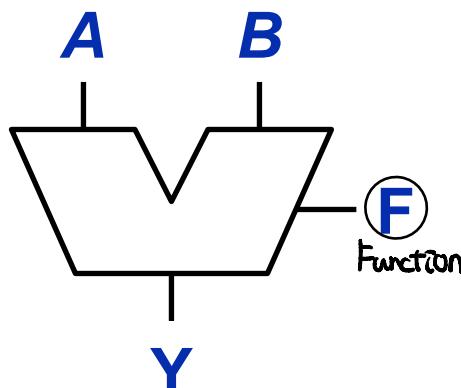
- ◆ Performs the actual computation(s)
- ◆ The processing unit can consist of many **functional units**
- ◆ We start with a simple Arithmetic and Logic Unit (ALU), which executes computation and logic operations
  - MIPS: add, sub, mult, and, nor, ...
- ◆ The ALU processes quantities that are referred to as **words**
  - Word length in MIPS is 32 bits (~~4 bytes~~)

ଅରିଥ୍ମେଟିକ୍ ଏଲ୍ୟୁନ୍ଟ୍

# ALU (Arithmetic logical unit)

→ 하위의 풀이된 학교에서 다양한 연습문제를 수강해 봄

- ◆ Combines a variety of arithmetic and logical operations into a single unit (that performs only one function at a time)
- ◆ Usually denoted with this symbol:



F[MSB to LSB 3-bit]      F[2]: MSB  
                                F[0]: LSB ) Endian 라는 상관없는

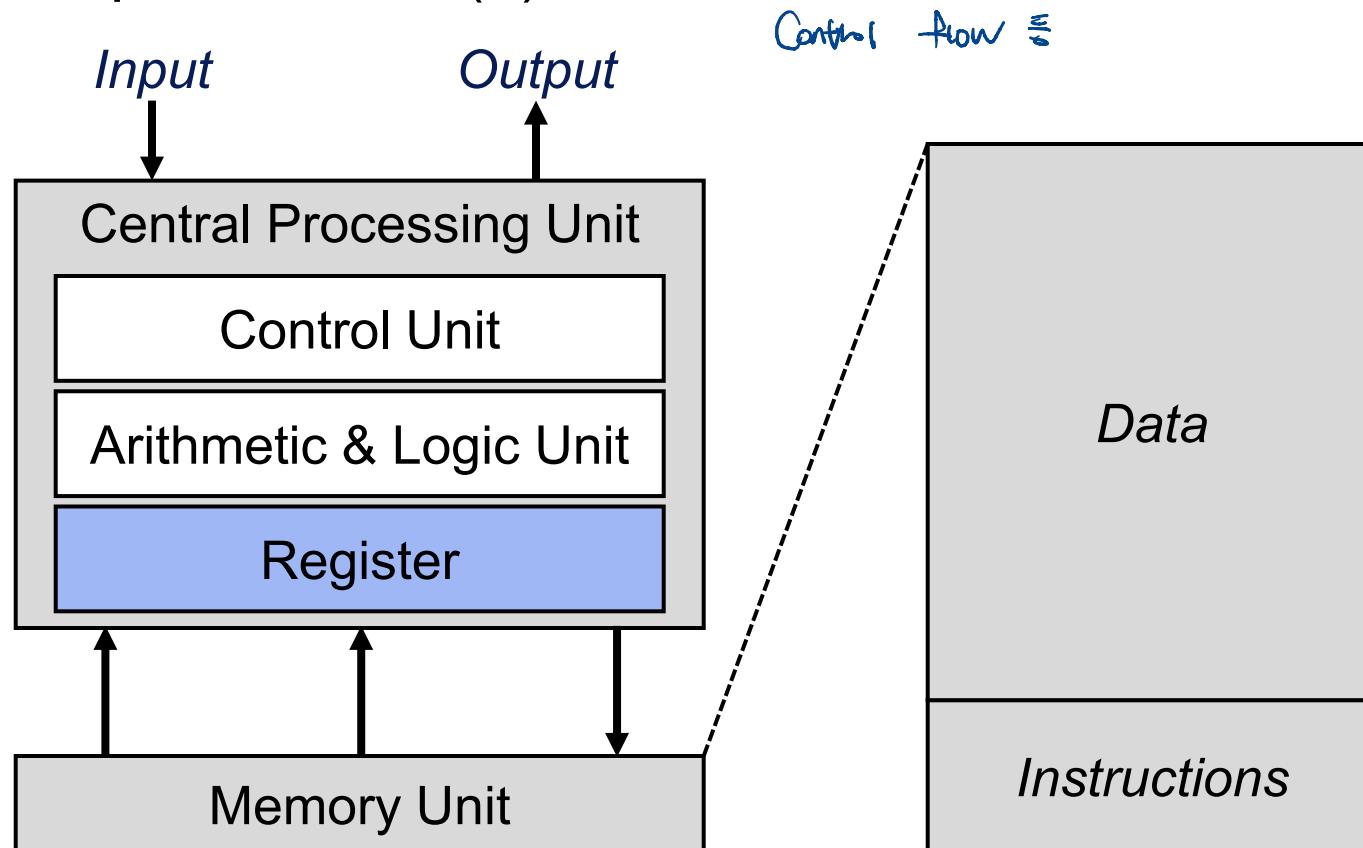
F[2:0]	Function
000	$Y = A \text{ and } B$
001	$Y = A \text{ or } B$
010	$A + B$
011	Not used <small>(다른 연산자를 위한  MSB-에서 LSB-로  예상되는 경우)</small>
100	$A - B$
101	$A * B$
110	$A / B$
111	SLT

One of the examples

만약 A가 1이라면 연산은 1, 아니면 0이다

# Von Neumann Architecture

- ◆ Both instructions and data are stored in the memory
- ◆ Each instruction dictates (1) which and how data are manipulated and (2) which instruction should be next



# Registers: fast storage

Synchronous write

→ 모든 write가 rising Clock 풀을 → System 전개를 예상하기 쉬워

Q) Reader writes? 동시에  
접촉되면??

- ◆ Memory is large but slow while registers are fast and small

- ◆ Processing unit utilizes registers during operation

- Ensure fast access to values to be processed in the ALU (stores temporary values)  
ALU 예상, 필요한 load 시 순간적으로 register에 저장
- Combinational read & synchronous write → can execute read + ALU  
+ write in a single cycle

Clock Signal 예상 일정

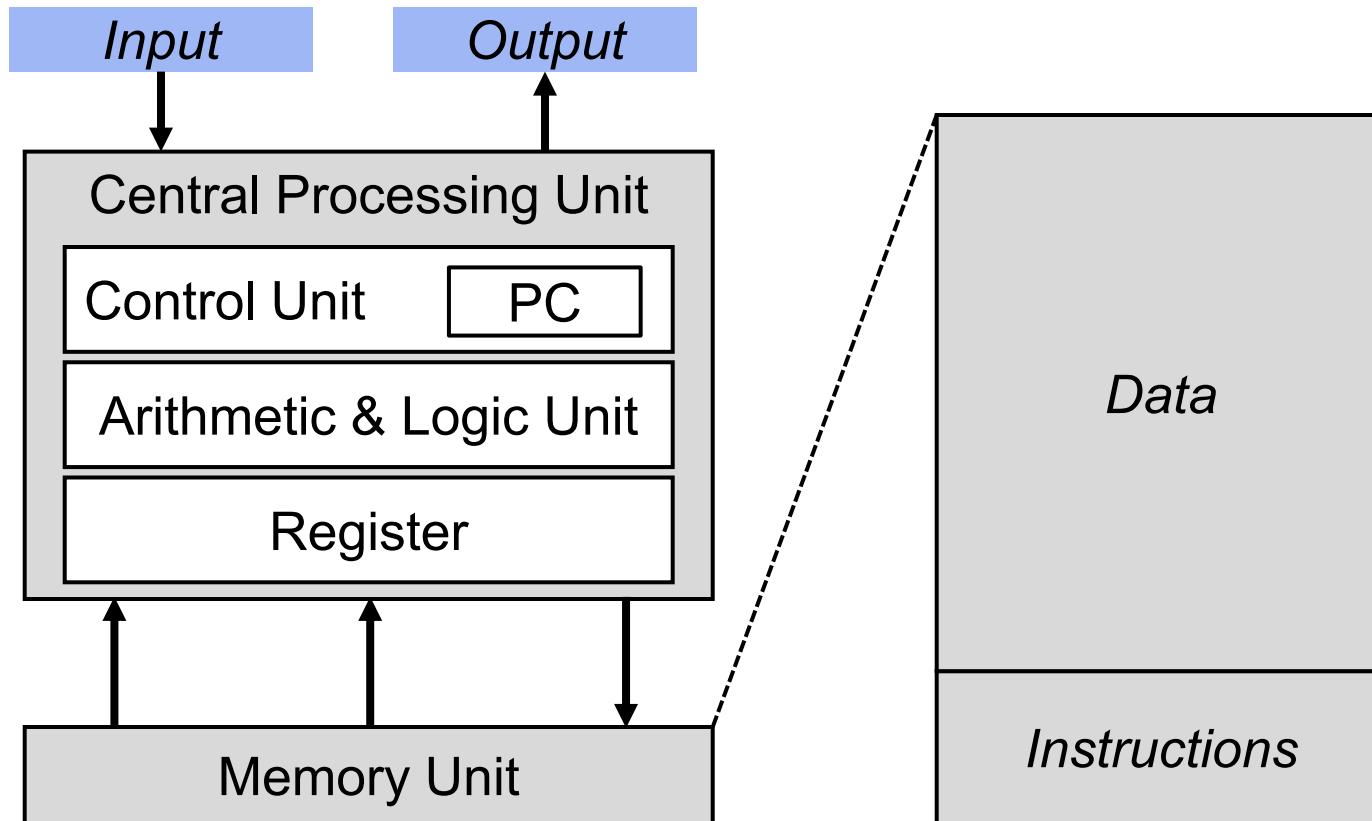
read+ALU: Clock 예상 일정  
write: Clock 예상 일정  
→ 즉 Cycle 예상 일정

ALU + Register 사용 (시스템마다 다르겠음)  
Clock 속도 edge 때는 write 가능

- ◆ Register Set or Register File (Single register file or 32H registers)

- Set of registers that can be manipulated by instructions
- MIPS has 32 registers  
 $\rightarrow 31 = 11111$  (5-bit)
  - R0 to R31: 5-bit register number (or Register ID)
  - Register size = Word length = 32 bits
- There are some special-purpose registers (e.g., \$fp, \$sp, ...)

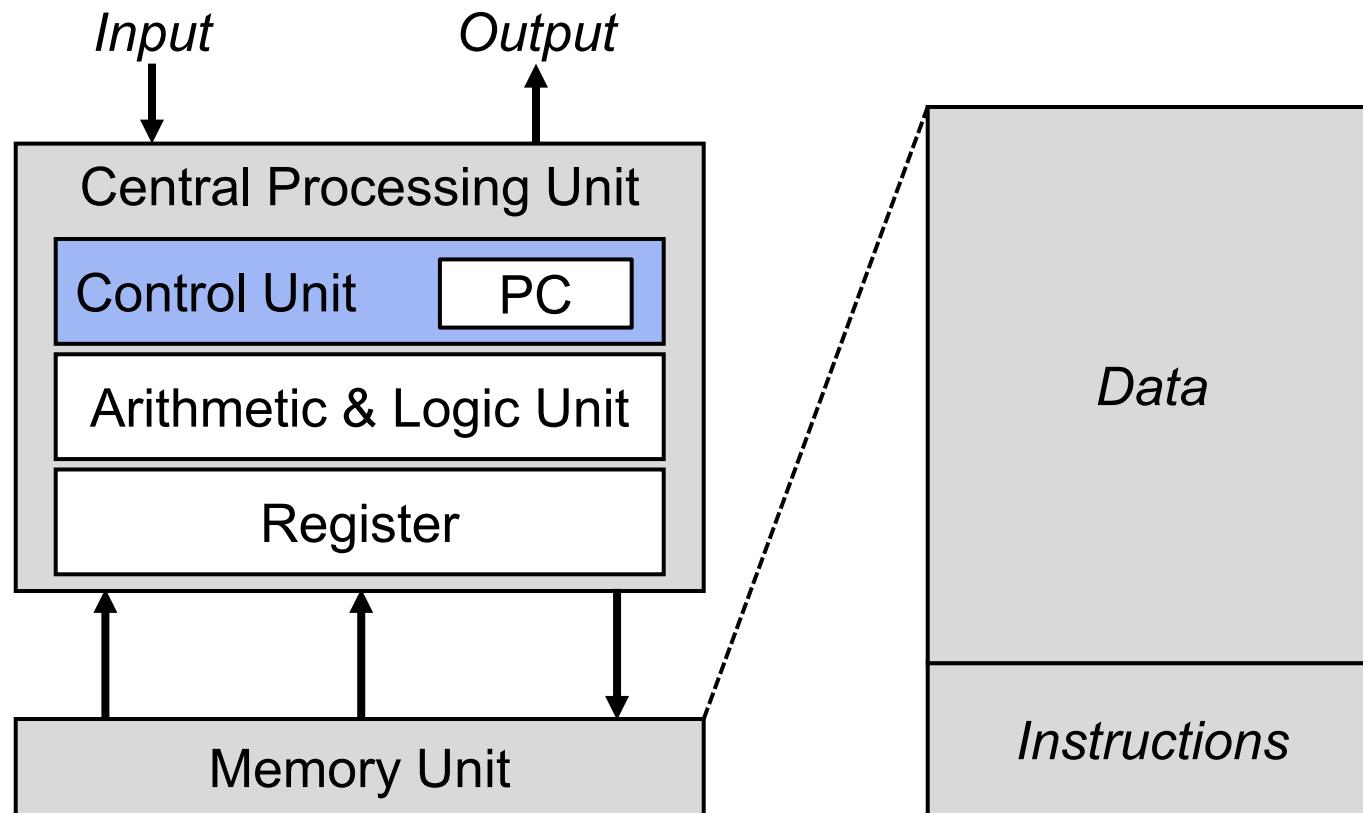
# Von Neumann Architecture



# Input and output

- ◆ Enable **information** to get into and out of a computer
- ◆ There are many input and output devices
- ◆ Input examples
  - keyboard, mouse, scanner, disks, ...
- ◆ Output examples
  - monitor, printer, disks, ...

# Von Neumann Architecture



# Control unit

- ◆ Enables a **step-by-step execution** of a program
  - Proceeds through each instruction in a program in sequence  
*순차로 보자!*
- ◆ Keeps track of which instruction is being processed, via an instruction register  
*Lo ~ L31 을 통하여 어떤 CPU 내부의 전용 메모리  
주소는 Instruction Register 라고  
설정된 주소에 해당하는 주소를 주소부에 넣어주면*  
*Fetch → Decode → Execute*  
*Load Load 명령어와 같은 명령어를 주소부에 넣어주면  
실행부에서 실행된다.*
- ◆ Keeps track of which instruction to process next, using a program counter (PC) and determines which instruction to process next  
*Instruction's address를 저장*  
*→ PC + 4 Control flow instruction의 다음은  
PC + 4*

# Wrap up

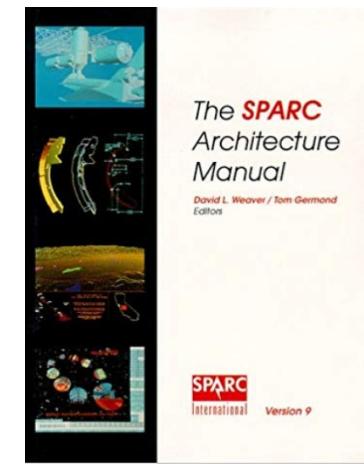
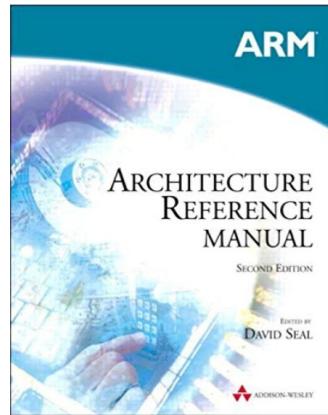
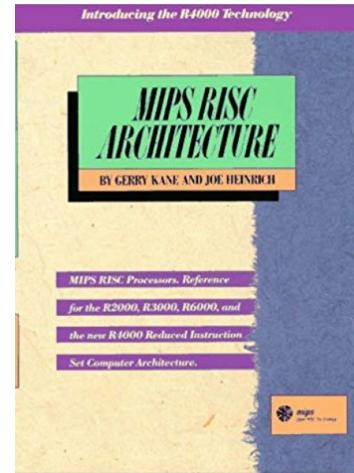
- ◆ **Memory Unit:** stores (1) program (instructions) and (2) data required in program execution
- ◆ **Arithmetic & Logical Unit:** performs the actual computation(s)
- ◆ **Input & Output:** enables information to get in and out of the computer
- ◆ **Control Unit:** enables a step-by-step execution of a program

# Instruction Set Architecture (ISA)

“User’s manual for the computer”

— 사용자의 매뉴얼

# Architecture Manuals



*Each vendor specifies the instruction sets in a different way*

# Architecture\*

- ◆ “The term *architecture* is used here to describe the attributes of a system as seen by the programmer, i.e., the **conceptual structure** and **functional behavior**, as distinct from the organization of the data flow and controls, the logical design, and the physical implementation.”

--- footnote on page 1, Architecture of the IBM System/360,  
Amdahl, Blaauw and Brooks, 1964.

# How to specify what a computer does?

- ◆ **Architecture Level** : User-level에 가까움
  - Car: driving manual & operation manual  
*// you don't have to be a car mechanic to drive a car.*
  - Computer : program manual  
*// you don't have to be a circuit designer to program a computer.*
  - 헌법학연 풍경지침
- ◆ **Microarchitecture (implementation) Level** : Backstage 돌아가는 권리  
Architecture의 구현
  - A particular car design has a certain configuration of (설계자)  
electrical/mechanical components (e.g., v8 engine vs. v4 engine)
  - A particular computer design has a certain configuration of datapath and control logic units (e.g., adder type, cache, ...)
  - 구현자 - 성능측정에서는 Architecture보다 중요하다.

# What are specified/decided in an ISA?

- ◆ Data format and size
  - character, binary, decimal, floating point, ...
- ◆ “Programmer Visible State” (a.k.a. architectural state)
  - memory, registers, program counter (PC), etc.
- ◆ Instructions: how to “change” the programmer visible state?
  - What to perform and what to perform next
  - Where the operands are
- ◆ How to interface with the outside world? ISA와 외부 세계의 인터페이스
- ◆ Protection and privileged operations
- ◆ Software conventions

***Often, you compromise performance for future scalability and compatibility***

# General Instruction Classes

## ◆ Arithmetic and logical operations (e.g., add, sub, and, or)

- 1) Load operands from specified locations  
→ **Registere에 지정된 값을 불러온다.**
- 2) Compute a result as a function of the operands
- 3) Store result to a specified location
- 4) Update PC to the next sequential instruction

## ◆ Data movement operations (e.g., load, store)

- 1) Fetch operands from specified locations  
→ **메모리에서 지정된 위치에서 값을 가져온다.**
- 2) Store operand values to specified locations  
→ **값을 지정된 위치에 저장한다.**
- 3) Update PC to the next sequential instruction

## ◆ Control flow operations (e.g., branch, jump)

- 1) Fetch operands from specified locations
- 2) Compute a **branch condition** and a **target address**
- 3) If “**branch condition is true**”      then  $PC \leftarrow \text{target address}$   
    else  $PC \leftarrow \text{next seq. instruction}$

**Generally defined to be atomic**

# Instructions for operations

- ◆ Different ISAs utilize different instruction for operations
- ◆ Number of Operands

- **Monadic**

OP in2

- **Binatic**

OP inout, in2

- Save memory (smaller instruction size)

- **Triadic**

OP out, in1, in2 ← MIPS32 ISA

→ 이는 ALU operand 2개를 동시에 메모리에서 불러와서Register에 동시에 저장하는 것이다.

- ◆ Can ALU operands be in memory? MIPS Register file에서 주소는 2~15까지!

- **No!**

ADD r1 r2 r3 (r1~r3 are registers)

- You should load memory to the register before ALU operations

- **Yes!**

ADD r2, r1, [2000]

→ MIPS는址 x, 다른 메모리에는 가능할 수도

**Different methods depending on the ISA**

# Instructions for memory addressing

[W \$ft, offset \$hs: ft off [rst offset] z13]

- ◆ Absolute  $lw\ rt, \text{10000}$   
-  $lw$ : load word use immediate value as address
  - ◆ Register Indirect  $lw\ rt, (r_{base})$   
- use Register[ $r_{base}$ ] as address
  - ◆ Displaced or based  $lw\ rt, \text{offset}(r_{base})$   
- use offset+ Register[ $r_{base}$ ] as address
  - ◆ Indexed  $lw\ rt, (r_{base}, r_{index})$   
- use Register[ $r_{base}$ ]+ Register[ $r_{index}$ ] as address
  - ◆ Memory Indirect  $lw\ rt ((r_{base}))$   
- use value at  $M[\text{Register}[r_{base}]]$  as address

MIPS 2진 X → Register Indirect 2진 A로 1진을  
가져온다.

# *Complicated memory addressing modes simplify program*

# MIPS RISC

- ◆ Simple operations
  - 2-input, 1-output arithmetic and logical operations
  - Only few alternatives exist to do the same thing
- ◆ Simple data movements
  - ALU ops are register-to-register (**need a large register file**)
  - Memory can be accessed by only load and store instructions  
→ “**Load-store architecture**”
- ◆ Simple branches
  - Limited varieties of branch conditions and targets
- ◆ Simple instruction encoding
  - All instructions encoded in the same number of bits  
*4 byte*
  - Only a few formats



***Such ISA intended for compiler advances  
rather than assembly programmers***

# Evolution of ISA

- ◆ Why were the earlier ISAs so simple?
  - Technology limitation
  - Inexperience, lack of precedence
- ◆ Why did it get so complicated later?
  - Complex instruction set architecture (CISC)
  - Ease of assembly programming 
  - Lack of memory size and performance (in 1970s~80s)
  - Micro-programmed implementation
- ◆ Why did it become simple again?
  - Reduced instruction set architecture (RISC)
  - Memory size and speed (cache!)
  - **Compilers**

***What about x86 (Intel and AMD)?***

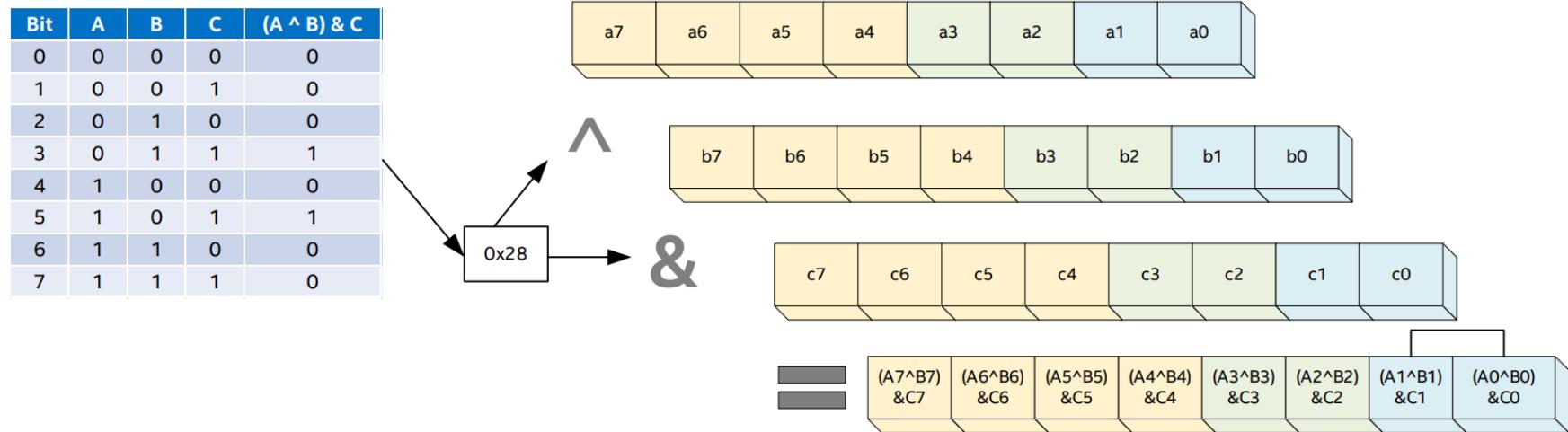
# RISC vs. CISC

- ◆ RISC (Reduced instruction set architecture)
  - The hardware exposes only basic operations as an ISA
  - Simpler instruction (+ simple decoding)
  - Fixed-width instruction (one word)
- ◆ CISC (Complex instruction set architecture)
  - Exposes more complex operations (combination of multiple RISC-style operations)
  - E.g., A single instruction may ... load the data from both memory and register, perform addition, and write back the result
    - Remember ADD r2, r1, [2000]

*Modern CPUs are RISC, but x86 (Intel + AMD) are CISC  
Why??*

# ISA extension in modern processors - 1

- ◆ Intel AVX enables **vector operations** by defining an Intel AVX-512 ISA



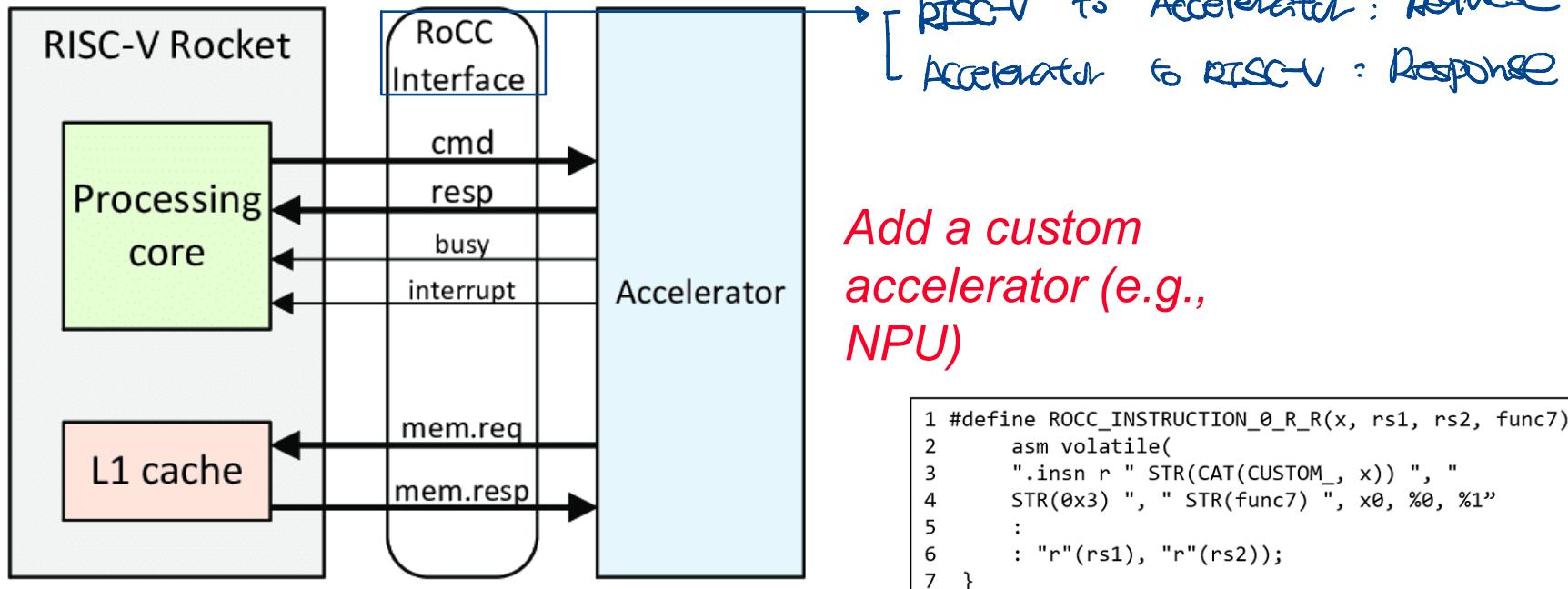
## INSTRUCTION SET      C INTRINSIC FORM OF INSTRUCTION

AVX-512	<code>_m512i_mm512_ternarylogic_epi64 (_m512i a, _m512i b, _m512i c, int imm8)</code>
AVX-512	<code>_m512i_mm512_mask_ternarylogic_epi64 (_m512i src, __mmask8 k, _m512i a, _m512i b, int imm8)</code>
AVX-512	<code>_m512i_mm512_maskz_ternarylogic_epi64 (_mmask8 k, _m512i a, _m512i b, _m512i c, int imm8)</code>

# ISA extension in modern processors - 2

- ◆ Rocket custom coprocessor (RoCC) interface in RISC-V
  - The computer architects can add a custom coprocessor
  - The RoCC interface controls the coprocessor

*Controls the accelerator using RoCC*



# Wrap-up: Terminologies

- ◆ Instruction Set Architecture (ISA)
  - The machine behavior as observable and controllable by the programmer
- ◆ Instruction Set
  - The set of commands understood by the computer
- ◆ Assembly Code
  - A collection of instructions expressed in “textual” format
    - e.g. Add r1, r2, r3
  - Converted to machine code by an assembler
  - One-to-one correspondence with machine code
- ◆ Machine Code
  - A collection of instructions encoded in binary format  
0101000...
  - Directly consumable by the hardware

# Let's dive into a concrete example

## MIPS ISA (32bit)

16216 → 224 457t  
→ ½ 321at

三

## ବ୍ୟକ୍ତିଗତ : ଜୀବନ ରୀତି

0x00000000	M[0] (8b)	= 1 byte	Word 1
0x00000001	M[1] (8b)		0x00000001 저장
0x00000002	M[2] (8b)	→ 8 비트와 같은 여기 메모리에서 저장	
	M[3] (8b)		여기 메모리에서 저장
	M[4] (8b)	→ Big Endian, Small Endian이 다른 결과.	
32:	16진수 8 bit → $2^{32}$		
		→ 32-bit address space	
	M[N-1] (8b)		
	Memory	→ Address 32	

# address space

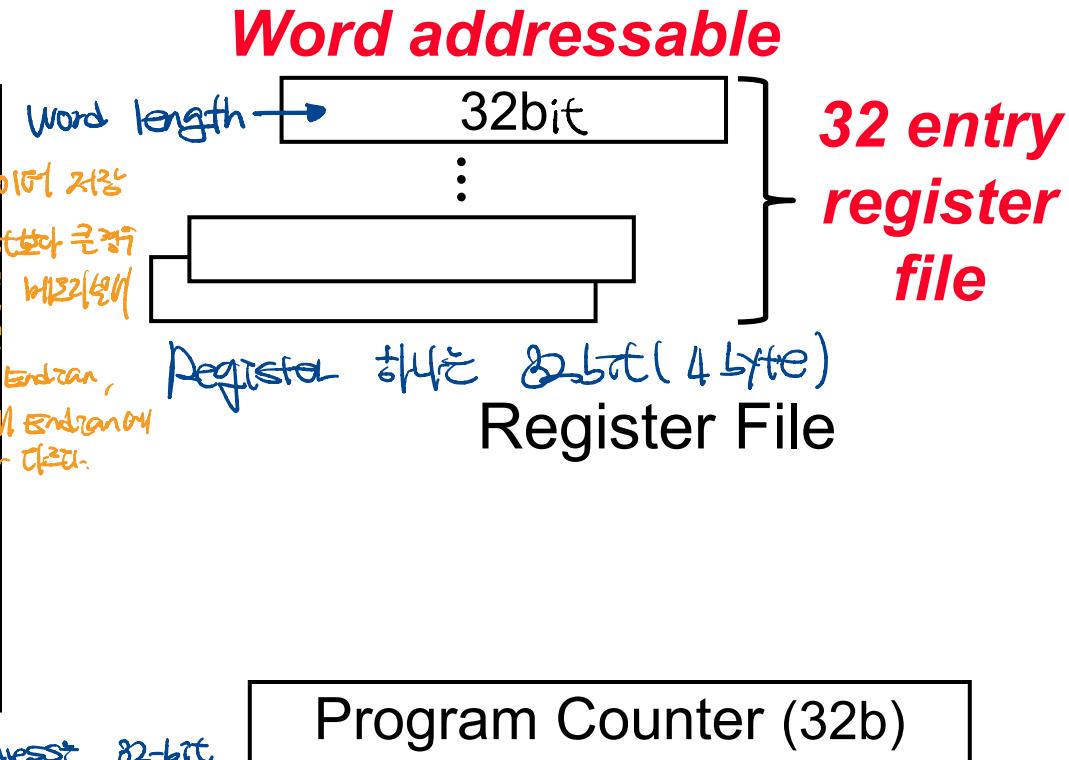
**+ Byte addressable**

→ 82 bit architecture

## Program Counter (32b)

# **32bit instruction**

**MIPS architecture has (1) 32-bit word size and  
(2) a 32-entry register file**



# MIPS Data format

- ◆ Most of the data are in 32 bits (= word)
  - Instruction and data addresses
  - Signed and unsigned integers
- ◆ Unsigned data format (all values are positive)

$$(b_{31}b_{30}\dots b_0)_2 = 2^{31} \cdot b_{31} + 2^{30} \cdot b_{30} + \dots + 2^0 \cdot b_0$$

- ◆ Signed data format (There exists negative values)
  - 1's complement: change sign by flipping bits 전체 부호를 뒤집
    - Ex) neg(0111) → 1000, neg(1100) → 0011, ...
    - Complicated zero value (i.e., 1111 and 0000 are all 0)
  - 2's complement: change sign by inverting bits and add 1  
(Adopted in MIPS)
    - Ex) neg(0111) -> 1000 + 1 → 1001, neg(1100) → 0011 + 1 → 0100
    - Sign extension

# Data format - 2

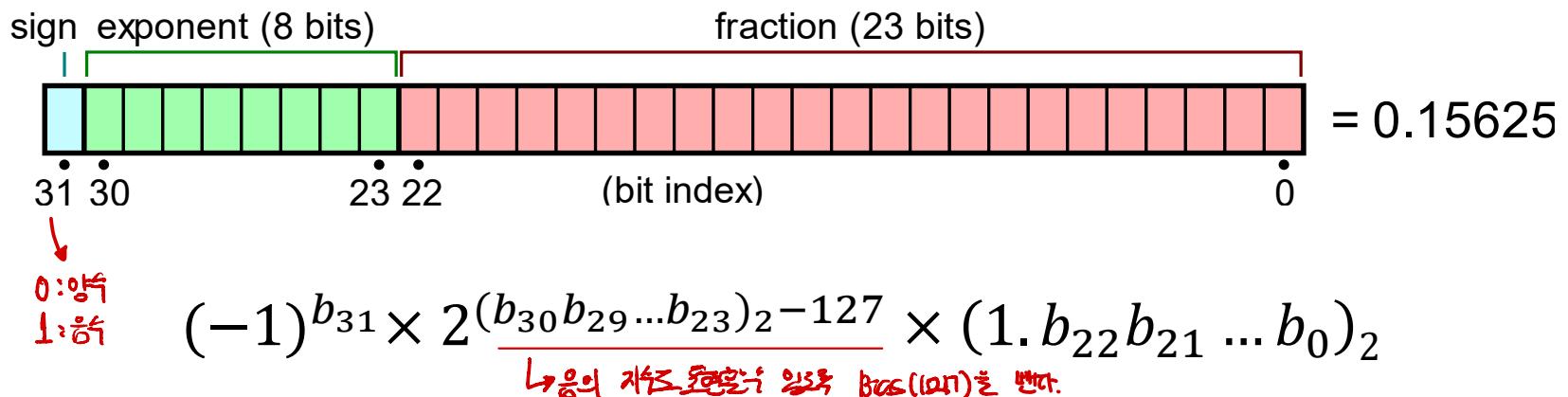
## ◆ Floating-point numbers

- Single precision:

1-bit sign, 8-bit exponent, 23-bit fraction (Adopted in MIPS)

- Double precision:

1-bit sign, 11-bit exponent, 52-bit fraction



*There can be precision errors (A floating point numbers may not represent exact value!)*

# Data format (Sidenote)

- ◆ We can further reduce the precision in some **fault tolerant applications** (e.g., neural networks)

- IEEE 754 half-precision

1-bit sign, 5-bit exponent, 10-bit fraction

→ Suffers from high accuracy drop due to small range

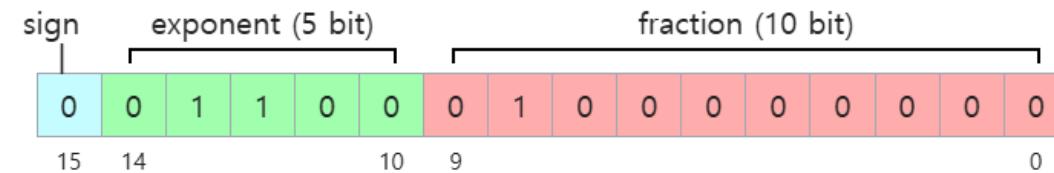
- bfloat (brain floating point) ← 인공지능(딥러닝)에서  
장르별 월드컵

1-bit sign, 8-bit exponent, 7-bit fraction

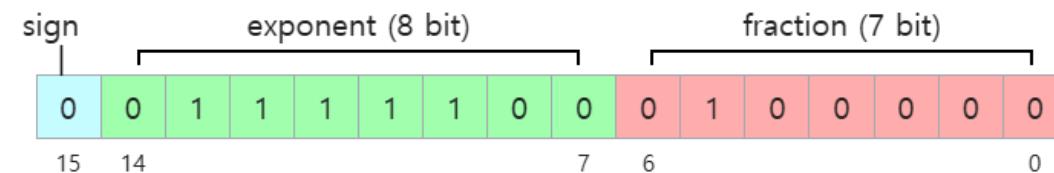
7 계쓰 | 정호드는 드래프트에 드는다

→ Improves accuracy by maintaining the range 

## IEEE half-precision 16-bit float



## **bfloat16**



# MIPS instruction formats

## ◆ Three simple formats

- R-type, 3 register operands

0	rs	rt	rd	shamt	funct	R-type
6-bit	5-bit	5-bit	5-bit	5-bit	6-bit	

- I-type, 2 register operands and 16-bit immediate operand

opcode	rs	rt	immediate	I-type
6-bit	5-bit	5-bit	16-bit	

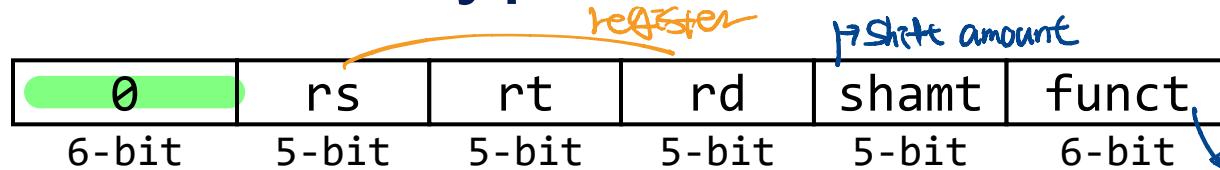
- J-type, 26-bit immediate operand

opcode	immediate	J-type
6-bit	26-bit	

## ◆ Simple Decoding

- 4 bytes per instruction, regardless of format (fixed size)
- Must be 4-byte aligned (2 LSB of PC must be 2b'00)
- Format and fields readily extractable

# R-type Instructions



R-type

→ R-type의 핵심 특징

어떤 operation  
제공되는 것인가

- ◆ funct: type of ALU operation

- Arithmetic: {signed, unsigned} x {ADD, SUB, MULT, DIV, ...}
- Logical: {AND, OR, XOR, NOR, ...}
- Shift: {Left, Right-Logical, Right-Arithmetic}

↳ 0으로 채우기

↳ Sign-bit3 채우기

- ◆ shamt: shift amount (only used for shift operation)

- ◆ Assembly and semantics

- Corresponding assembly: opcode rd rs rt
  - $GPR[rd] = GPR[rs] \text{ op } GPR[rt]$
  - PC = PC + 4 (use the next instruction)

rs와 rt를 연산하고 rd에 저장

No Branch or Jump

Memory, ALU, Control

# I-type Instructions – ALU ver.

opcode	rs	rt	immediate	I-type
6-bit	5-bit	5-bit	16-bit	

- ◆ opcode: there are immediate ALU instructions
  - addi, addiu, andi, ori, xori, slti, sltiu, lui, ...
    - ↑ immediate
    - ↓ unused
- ◆ Assembly and semantics
  - Corresponding assembly: **opcode rt rs immediate**
    - $\text{GPR}[rt] = \text{GPR}[rs] \text{ op sign-extend(immediate)}$  → sign-extended 32-bit
      - sign-extend(1000...000) = concat(11...1, 1000...000)
    - $\text{PC} = \text{PC} + 4$  (use the next instruction)
      - ↳ 32-bit PC
      - ↳ immediate → 16-bit
      - bit 7 is always 1 for 32-bit

## ★ What if you need to perform 32-bit immediate?

- lui at 0xABCD // store the upper 16 bits
- ori at 0x1234 // store the lower 16 bits

# Assembly Programming - ALU

- ◆ Break down high-level program constructs into a sequence of elemental operations
- ◆ E.g. High-level Code

```
f = ( g + h ) - ( i + 100 )
```

- ◆ Assembly Code
  - suppose  $f$ ,  $g$ ,  $h$ ,  $i$  are in  $r_f$ ,  $r_g$ ,  $r_h$ ,  $r_i$
  - suppose  $r_{temp}$  is a free (available) register

```
add rtemp rg rh          # temp = g+h
addi rf ri 100        # f = i+100
sub rf rtemp rf      # f = temp-f
```

# I-type Instructions – Memory ver.

opcode	rs	rt	offset	I-type
6-bit	5-bit	5-bit	16-bit	

- ◆ opcode: there are immediate memory instructions

- **lw**, **lh**, **lhu**, **lb**, **lbu**, **sw**, **sh**, **sb**, ...  
    half(16-bit)
- // lw indicates load, sw indicates store

- ◆ Assembly and semantics

- Corresponding assembler: **load/store rt offset(rs)**
  - **GPR[rt] = MEM[GPR[rs] + sign-extend(offset)] // load**
  - **MEM[GPR[rs] + sign-extend(offset)] = GPR[rt] // store**
  - **PC = PC + 4 (use the next instruction)** sign-extend from  
16-bit to 32-bit

# Assembly Programming - Memory

- ◆ E.g. High-level Code

$$A[8] = h + A[0]$$

where **A** is an array of integers (4-byte each)

- ◆ Assembly Code

- suppose **&A**, **h** are in **r<sub>A</sub>**, **r<sub>h</sub>**
- suppose **r<sub>temp</sub>** is a free (available) register

```
lw rtemp 0(rA)      # temp = A[0]
add rtemp rh rtemp  # temp = h + A[0]
sw rtemp 32(rA)    # A[8] = temp
                           # note A[8] is 32 bytes
                           #           from A[0]
```

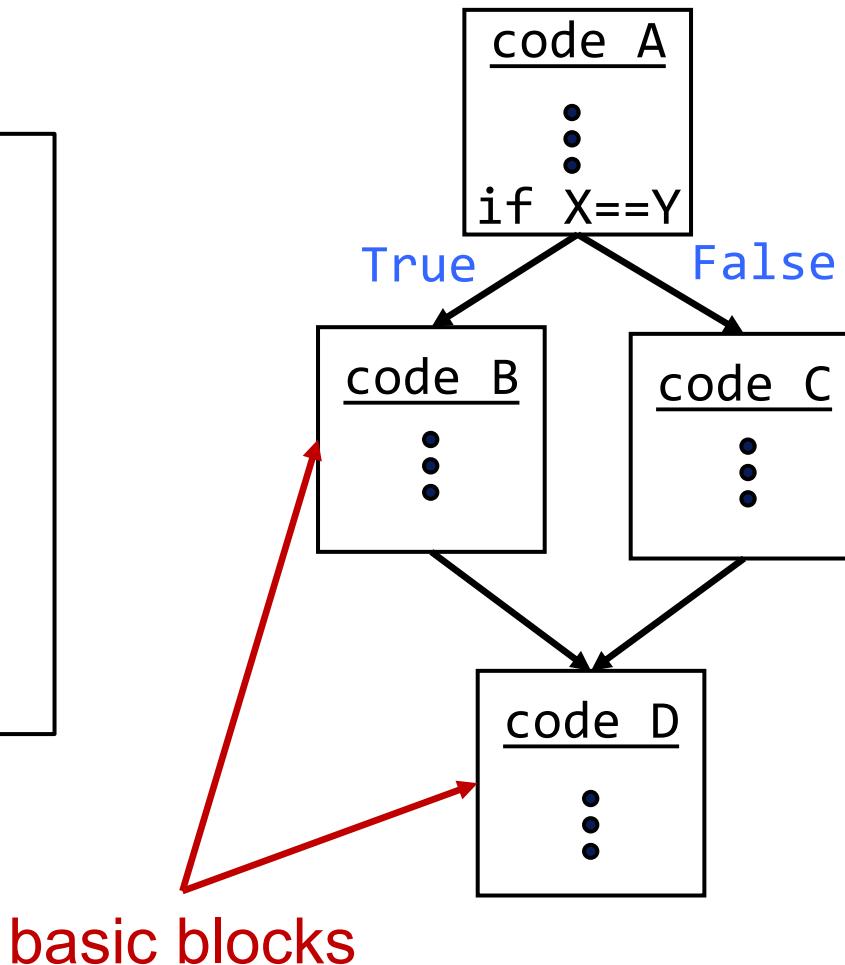
~~byte~~

4 byte x8

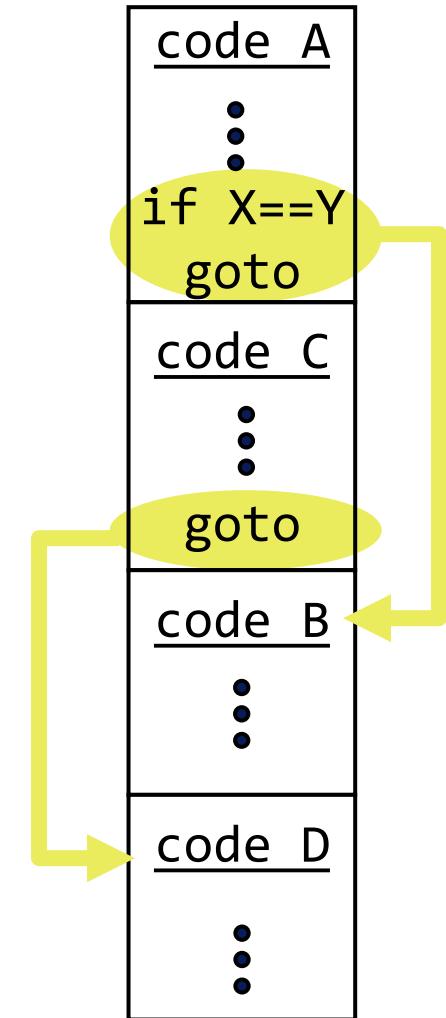
# Control Flow Instructions

## Control Flow Graph

```
Example code:  
code A  
if (X == Y){  
    code B  
}  
else {  
    code C  
}  
code D
```



## Assembly Code (linearized)



I-type 향수 PC에서 상대주소 이동 → 두 가지 경우 16-bit Sign-Extend

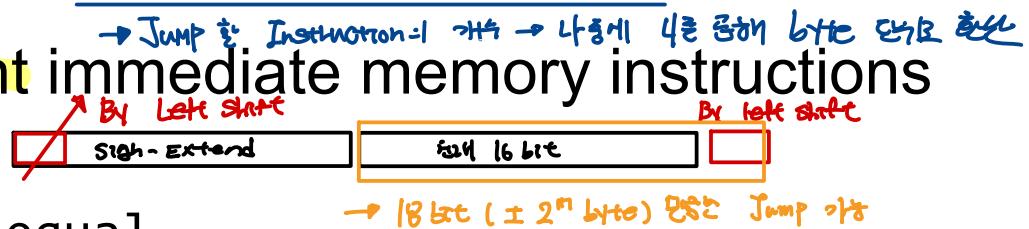
# I-type Instructions – Control ver.

opcode	rs	rt	label	I-type
6-bit	5-bit	5-bit	16-bit	

- ◆ opcode: there are **eight immediate memory instructions**

- bne, beq, ...
- branch equal or not equal ...

↳ rs는 rt와 비교

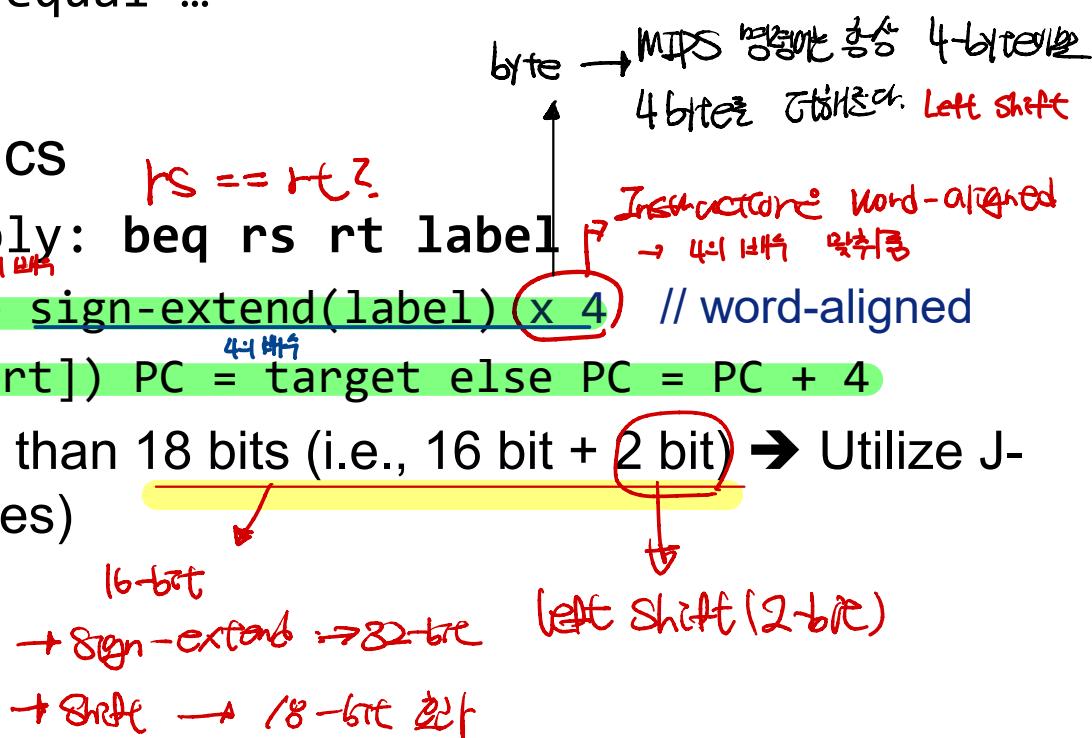


- ◆ Assembler and semantics

- Corresponding assembly:  $\text{beq } rs \text{ } rt \text{ } label$

• target =  $(PC + 4) + \text{sign-extend}(label) \times 4$  // word-aligned  
• if ( $GPR[rs] == GPR[rt]$ )  $PC = \text{target}$  else  $PC = PC + 4$

If you want to jump more than 18 bits (i.e., 16 bit + 2 bit) → Utilize J-type (in the following slides)



# J-type Instructions

opcode	label
6-bit	26-bit

J-type

- ◆ opcode: there are eight immediate memory instructions

- j, jal
 

MIPS는 J-type은 32-bit 전체 주소 표현하지 못한다.  
 $\rightarrow (PC+4) \div 4$  상위 4bit로 사용한다.

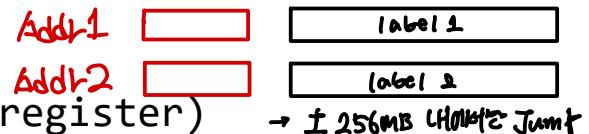


- ◆ Assembler and semantics

- target =  $(PC+4)[31:28] \times 2^{28}$  | bitwise-or zero-extend(label)  $\times 4$
- // use the first four bits of PC and append label  $\times 4$
- Corresponding assembly: j label
  - PC = target
- Corresponding assembly: jal label
  - GPR[ra] = PC + 4
    - save the next PC to ra (a dedicated register)
  - PC = target

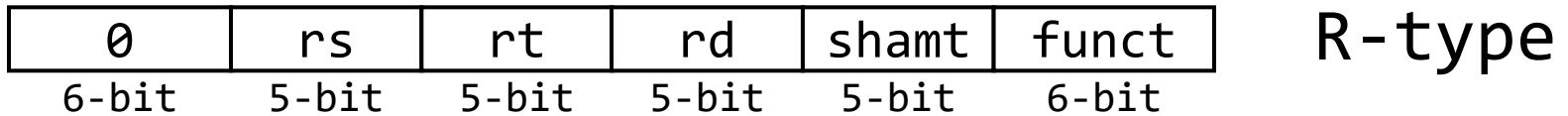
Jump PC에서 상위 4bit로  
 $\rightarrow$  현재 PC에서 256MB 범위로 Jump

Jump



zero extend을 사용하는 이유  
 $\rightarrow$  절대 주소 사용, 초기 설정이 어렵다.

# R-type Instructions – Control ver.



- ◆ There exists an R-type control instruction
  - jr
- ◆ Generally used along with jal label (upon a function call)
- ◆ Assembler and semantics
  - Corresponding assembly: **jr rs**
    - PC = GPR[rs]

# Assembly Programming - Control

- ◆ E.g. High-level Code

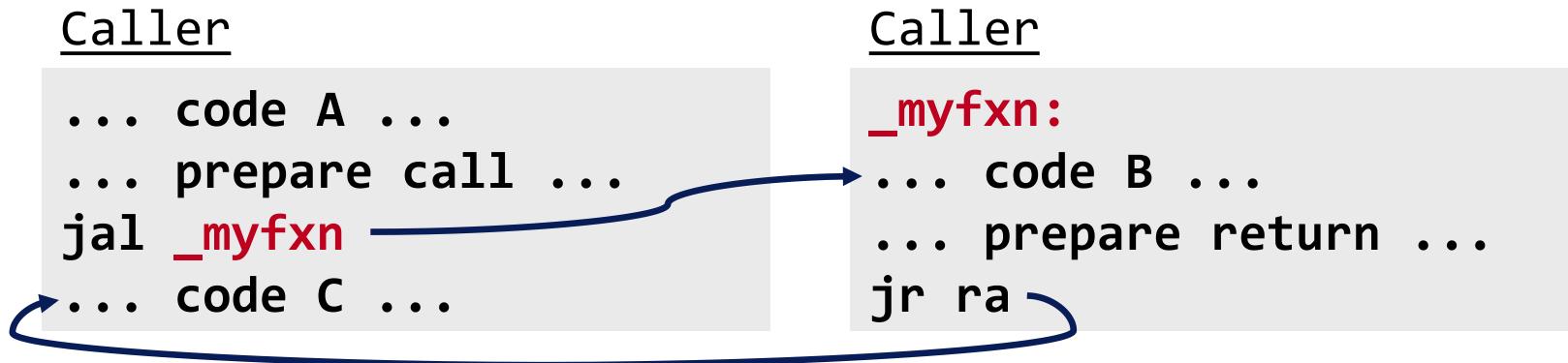
```
if (i == j) then
    e = g;
else
    e = h;
f = e;
```

- ◆ Assembly Code

- suppose `e`, `f`, `g`, `h`, `i`, `j` are in `re`, `rf`, `rg`, `rh`, `ri`, `rj`

```
bne ri rj L1          # L1 and L2 are addr labels
                           # assembler computes offset
add re rg r0 # e = g
j L2
L1: add re rh r0 # e = h
L2: add rf re r0 # f = e
```

# Assembly Programming – Function Call



- ◆ sequence: **A** → `call` **B** → `return` **C** ...
- ◆ There is a prologue and epilogue before and after a function call
  - How to pass an argument to the caller?
  - If **A** set `r10` to 1, what is the value of `r10` when **B** returns to **C**?
  - What registers can **B** use?
  - What happens to `ra` if **B** calls another function

*Let's dive into the register file usages in MIPS*

# R2000 Register Usage Convention

Register	Functionality	Alias
r0	Always 0	\$zero
r1	Reserved for the assembler	\$at
r2, r3	Function return values	\$v0, \$v1
r4 ~ r7	Function call arguments	\$a0 ~ \$a3
r8 ~ r15	“Caller-saved” temporaries	\$t0 ~ \$t7
r16 ~ r23	“Callee-saved” temporaries	\$s0 ~ \$s7
r24 ~ r25	“Caller-saved” temporaries	\$t8 ~ \$t9
r26, r27	Reserved for the OS (e.g., interrupts)	\$k0 ~ \$k1
r28	Global pointer	\$gp
r29	Stack pointer	\$sp
r30	Frame pointer	\$fp
r31	Return address	\$ra

# R2000 Register Usage Convention

Register	Functionality	Alias
r0	Always 0 → 많이 사용하지 않아 <del>전체</del> 전체적으로 사용하지 않음	\$zero
r1	Reserved for the assembler	\$at
r2, r3	<b>Stores 32-bit zero value</b>	
r4 ~ r7	Function call arguments	\$a0 ~ \$a3
r8 ~ r15	“Caller-saved” temporaries	\$t0 ~ \$t7
r16 ~ r23	“Callee-saved” temporaries	\$s0 ~ \$s7
r24 ~ r25	“Caller-saved” temporaries	\$t8 ~ \$t9
r26, r27	Reserved for the OS (e.g., interrupts)	\$k0 ~ \$k1
r28	Global pointer	\$gp
r29	Stack pointer	\$sp
r30	Frame pointer	\$fp
r31	Return address	\$ra

# R2000 Register Usage Convention

Register	Functionality	Alias
r0	Always 0	\$zero
r1	Reserved for the assembler	\$at
r2, r3	Function return values	\$v0, \$v1
r4 ~ r7	<b>Temporary register to enable 32-bit immediate operation</b>	
r8 ~ r15	Caller-saved temporaries	\$t0 ~ \$t7
r16 ~ r23	“Callee-saved” temporaries	\$s0 ~ \$s7
r24 ~ r25	“Caller-saved” temporaries	\$t8 ~ \$t9
r26, r27	Reserved for the OS (e.g., interrupts)	\$k0 ~ \$k1
r28	Global pointer	\$gp
r29	Stack pointer	\$sp
r30	Frame pointer	\$fp
r31	Return address	\$ra

# R2000 Register Usage Convention

Register	Functionality	Alias
r0	Always 0	\$zero
r1	Reserved for the assembler	\$at
r2, r3	Function return values	\$v0, \$v1
r4 ~ r7	Function call arguments	\$a0 ~ \$a3
r8 ~ r15	"Caller-saved" temporaries	
r16 ~ r23		\$t0 ~ \$t7
r24 ~ r25		
r26, r27		
r28		
r29	Stack pointer	\$sp
r30	Frame pointer	\$fp
r31	Return address	\$ra

**Used to pass values between function calls!**

- Q1. We need to pass more than 4 values** *→ for argument*
- Q2. The function returns non-primitive data (e.g., struct)** *→ 메모리 주소를 다른 사용하지  
않음*

# R2000 Register Usage Convention

Register	Functionality	Alias
r0	Always 0	\$zero
r1	Reserved for the assembler	\$at
r2, r3	Function return values	\$v0, \$v1
r4 ~ r7	Function call arguments	\$a0 ~ \$a3
r8 ~ r15	“Caller-saved” temporaries	\$t0 ~ \$t7
r16 ~ r23	“Callee-saved” temporaries	\$s0 ~ \$s7
r24 ~ r25	“Caller-saved” temporaries	\$t8 ~ \$t9
r26, r27		
r28		
r29		
r30		
r31		

**Temporary values are stored in the register file, instead of using memory**

→ Register 파일에 대체로 메모리에 저장하기

**Q. The # of temporaries exceeds 18**

**Q. What are caller-saved and callee-saved?**

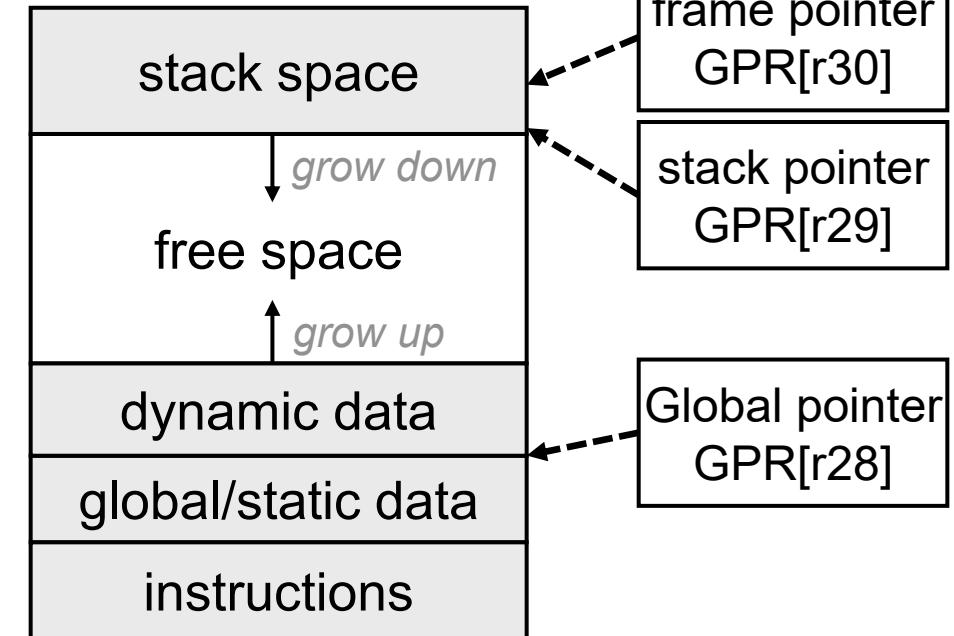
# R2000 Register Usage Convention

Register	Functionality	Alias
r0	Always 0	\$zero
r1	Reserved for the assembler	\$at
r2, r3	Function return values	\$v0, \$v1
r4 ~ r7	Function call arguments	\$a0 ~ \$a3
r8 ~ r15	“Caller-saved” temporaries	\$t0 ~ \$t7
r16 ~ r23		
r24 ~ r25	“Caller-saved” temporaries	\$t8 ~ \$t9
r26, r27	Reserved for the OS (e.g., interrupts)	\$k0 ~ \$k1
r28	Global pointer	\$gp ↳ <i>내부</i> <i>spcall</i>
r29	Stack pointer	\$sp ↳ <i>스택</i>
r30	Frame pointer	\$fp → <i>\$k0</i> & <i>\$k1</i> <i>지정</i>
r31	Return address	\$ra ↳ <i>리턴</i>

# R2000 Registers

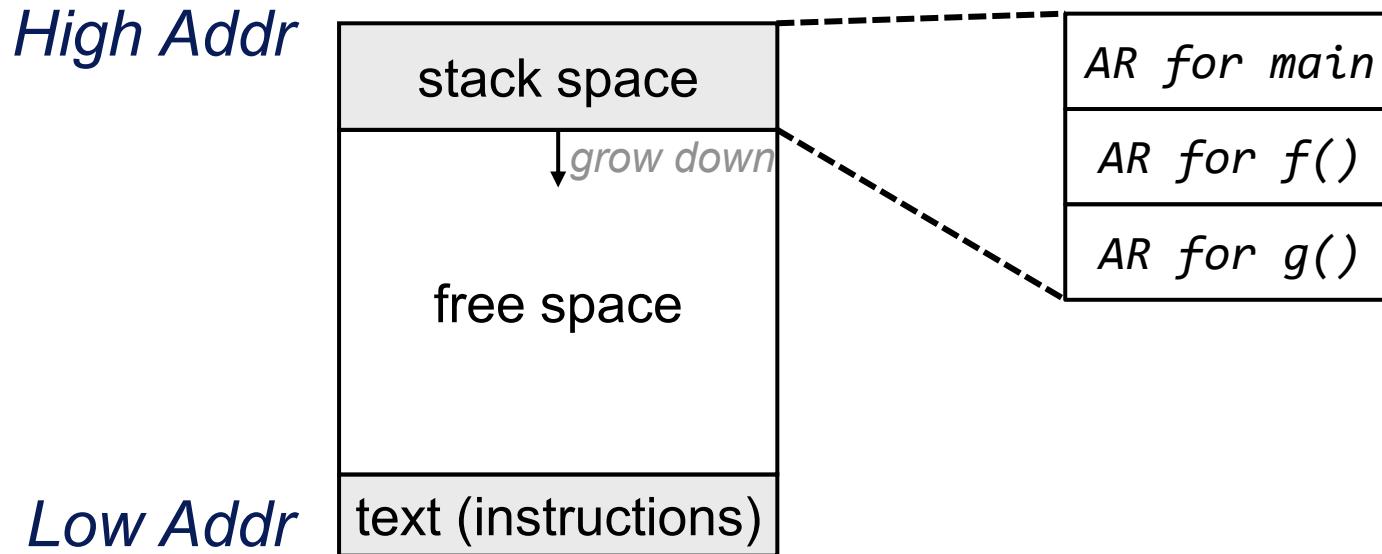
Register	Description	Format
r0		
r1	Reserved	
r2, r3	Function	
r4 ~ r7	Function	
r8 ~ r15	“Caller”	
r16 ~ r23	“Callee”	
r24 ~ r25	“Caller”	
r26, r27	Reserved for the OS (e.g., interrupts)	\$R0 ~ \$R1
r28	Global pointer	\$gp
r29	Stack pointer → stack ↗	\$sp
r30	Frame pointer	\$fp
r31	Return address	\$ra

**Pointers to address main memory**



# Recall: Stack management

- ◆ Stack data is stored starting from the low address, which grows downwards
- ◆ The information to manage one function call is called activation record (AR) or frame



# Calling a function - 1

```
int f(int x) {  
    reg1 = 11;  
    reg2 = 16;  
    val = reg1+reg2;  
    ...  
    return 20;  
}  
  
void main() {  
    reg1 = 10;  
    reg2 = 15;  
    val = f(5);  
    ...  
}
```

$\$fp \rightarrow 0x0$	AR for main
$\$sp \rightarrow 0x100$	
①	Argument x (i.e., 5)
②	Caller-Saved Regs (reg1 = 10) → Main()에서의 초기 설정값을 넘기기 위해 Caller가 저장하는 registros
③	Return addr (i.e., 0x20)

\* 컴파일러가 정한거 봄  
→ 풀려면 알면

# Calling a function - 2

```
0x00:  
0x04:  
0x08:  
0x0c:  
0x10:  
  
0x14:  
0x18:  
0x1c:  
0x20:
```

```
int f(int x) {  
    reg1 = 11;  
    reg2 = 16;  
    val = reg1+reg2;  
    ...  
    return 20;  
}  
  
void main() {  
    reg1 = 10;  
    reg2 = 15;  
    val = f(5);  
    ...  
}
```

reg1 = Caller-Saved register  
reg2 = Callee-Saved register

⑤ \$fp → frame[5]

④

⑥

⑦

⑧

\$sp → 0x11c registers → val

AR for main

Argument x (i.e., 5)

Caller-Saved Regs  
(reg1 = 10)

Return addr (i.e., 0x20)

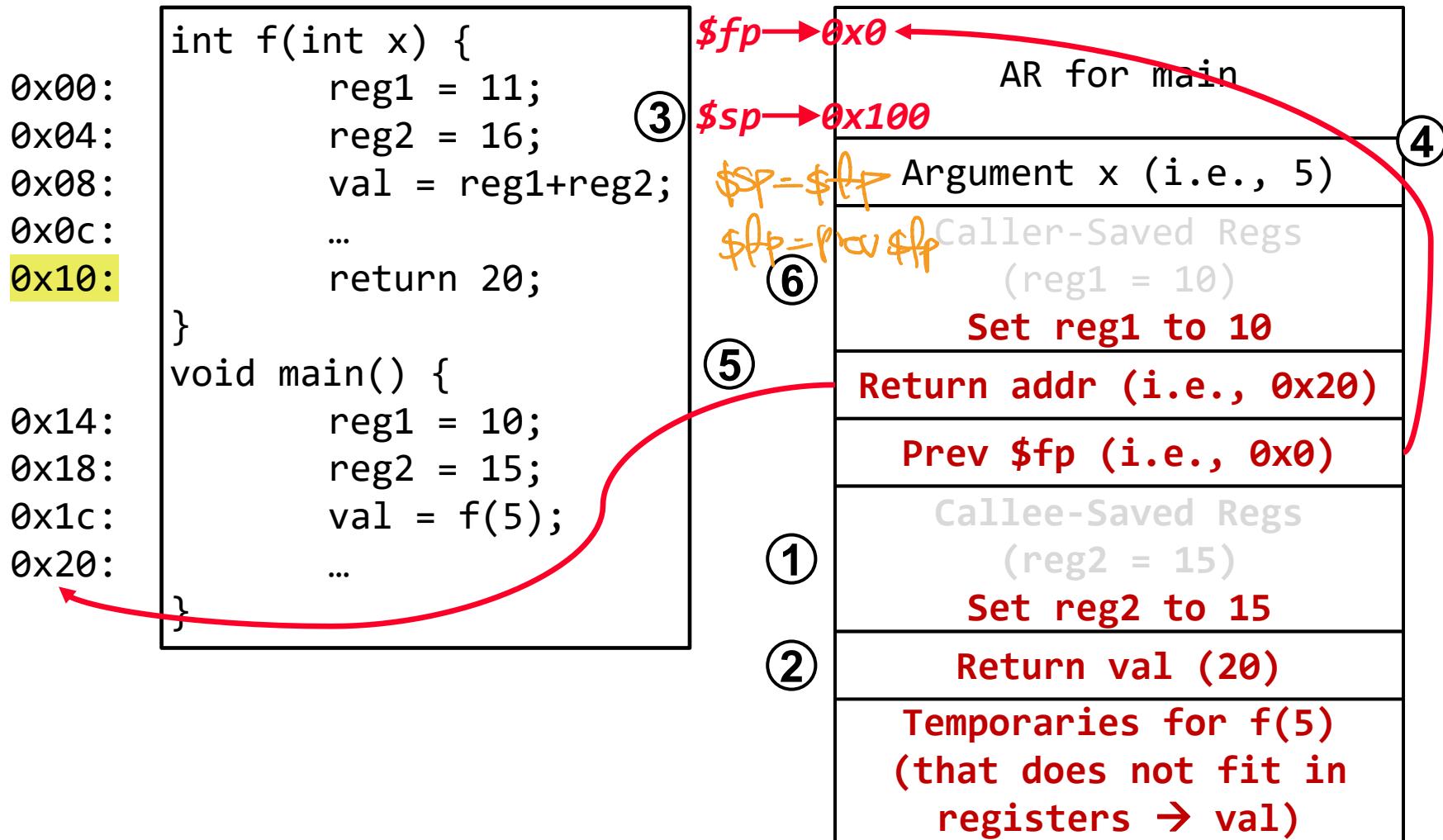
Prev \$fp (i.e., 0x0)

Callee-Saved Regs  
(reg2 = 15)

Return val (20)

Temporaries for f(5)  
(that does not fit in

# Returning a function



# Wrap up - 1

- ◆ Code before call instruction

- Push each actual parameter
- Push caller-saved registers
- Push return address (current PC) and jump to callee code

- ◆ Prologue = code at function entry

- Push dynamic link (i.e., FP)
- Old stack pointer becomes new frame pointer
- Push callee-saved registers
- Push local variables

Arguments
Caller-saved regs
Return Addr
Prev FP
Callee-saved regs
Local Variables

# Wrap up - 2

- ◆ Epilogue = code at return instruction
  - Pop (restore) callee-saved registers
  - Store return value at appropriate place
  - Restore old stack pointer (pop callee frame)
  - Pop old frame pointer
  - Pop return address and jump to that address
- ◆ Code after call
  - Pop (restore) caller-saved registers
  - Use return value

Arguments
Caller-saved regs
Return Addr
Prev FP
Callee-saved regs
Local Variables

***You will learn further details in the compiler class  
(next semester)***

# Something to think about!

- ◆ Why do we need to store return value to the stack?  
Why not use registers (e.g., r2, r3)?
  - A. The program may return complex values  
*h2, h3 는 복잡한 값 return 할 때 뭇 뒤에 써야 한다.*
- ◆ Why do we store arguments to the stack? Why not use registers (e.g., r4~r7)?
  - A. The function may need multiple arguments
  - A. The function may call another function (which will overwrite r4~r7)
- ◆ Why are some registers caller-saved while others are callee-saved?

# Caller and Callee Saved Registers

→ function call 와서 Register 를 2.1 번까지 쓸때는 2번으로

- ◆ **Caller-saved registers (AKA volatile registers, or call-clobbered):** used to hold temporary quantities that need not be preserved across calls
  - It's the caller's responsibility to save / restore the registers (if the caller wants)
  
- ◆ **Callee-saved registers (AKA non-volatile registers, or call-preserved):** used to hold long-lived values that should be preserved across calls
  - It's the callee's responsibility to save / restore the registers (necessary)

# Question?



*Announcements:* Will skip chapter 3 (arithmetic)  
and jump to Chapter 4 (edition 6) in the next lecture

*Reading:* **P&H Ch 2**

*Handouts:* None