

Computer Architecture

IN BA3 - Paolo IENNE

Notes by Ali EL AZDI

Introduction

This document is designed to offer a LaTeX-styled overview of the Computer Architecture course, emphasizing brevity and clarity. Should there be any inaccuracies or areas for improvement, please reach out at ali.elazdi@epfl.ch for corrections. For the latest version of the PDF, you can check the following link: <https://elazdi-al.github.io/comparch/index.html>. Feel free to send a pull request to propose any changes you think might be a useful addition to the course content or a modification.

<https://github.com/elazdi-al/comparch/blob/main/main.pdf>

Contents

Contents	3
1 Part I(a) - ISA Reminder, Assembly Language, Compiler - W 1.1	7
1.1 From High Level Languages to Assembly Language	7
1.1.1 High Level Languages	7
1.1.2 Assembly Language	7
1.2 Processors	8
1.3 Joint or Disjoint Program and Data Memories	9
1.4 The Encoding problem	10
1.4.1 The Stupid Solution	10
1.4.2 RISC-V Encoding (The Solution)	10
1.4.3 Automating this process	11
1.5 ISA (Instruction Set Architecture)	12
2 Part I(b) - ISA, Functions, and Stack - W 1.2	13
2.1 Arithmetic and Logic Instructions in RISCV	13
2.1.1 Constants must be encoded on 12 bits	13
2.1.2 Assembler Directives	13
2.1.3 The x0 Register	14
2.2 PseudoInstructions	14
2.2.1 Control flow instructions	15
2.2.2 If-Then-Else	15
2.2.3 Jumps and Branches	15
2.2.4 Comparaisons	16
2.2.5 Do-While	16
2.3 Functions	16
2.3.1 Jump to the Function/Retun control to the calling program	16
2.3.2 Jump Instructions	17
2.3.3 Register Conventions	18
2.3.4 Back to the good (not so good) approach	18
2.3.5 One simple solution (still not good)	18
2.3.6 Acquire storage resources the function needs (still not it)	19
2.3.7 The Stack	19
2.3.8 Spilling Registers to Memory	21
2.3.9 Register across functions	21
2.3.10 Preserving Registers	22
2.4 Passing Arguments in RISC-V	22
2.4.1 Option 1: Using Registers	22
2.4.2 Option 2: Using the Stack	23
2.4.3 The RISC-V Approach	23

2.5	Summary of RISC-V Register Conventions	23
3	Part I(c) - ISA Memory and Addressing Modes - W 2.1	24
3.1	Memory	24
3.1.1	Address and Data	24
3.2	Many Types of Memories	25
3.2.1	Functional Taxonomy of Memories	25
3.2.2	Taxonomy of Random Access Memories	25
3.2.3	Basic Structure	26
3.2.4	Write Operations	26
3.2.5	Read Operations	26
3.2.6	Practical SRAMs	26
3.2.7	DRAMs	27
3.2.8	Ideal Random Access Memory	27
3.2.9	Physical Organisation	27
3.2.10	Realistic ROM Array	28
3.2.11	Static Ram Typical Interface	28
3.3	Typical Asynchronous SRAM Read Cycle	28
3.4	Where is Memory in the Processor?	29
3.4.1	Arithmetic and Logic Instructions	29
3.5	More Addressing Modes? Not in RISC-V!	30
3.5.1	Word Adressed Memory	31
3.5.2	Loading Words (lw) and Instructions	31
3.5.3	Loading Bytes (lb)	31
3.5.4	A Few More Load/Store Instructions	31
3.5.5	Access as it is more suitable	32
3.5.6	Loading Bytes (lb)	33
4	Part I(d) - ISA Arrays and Data Structures - W 2.2	34
4.1	Arrays	34
4.1.1	Different Ways to Store Arrays	34
4.1.2	Adding Positive Elements	35
4.1.3	Pointer to Memory vs Index in Array	36
5	Part I(e) - ISA Arithmetic - W 3.1, 3.2	38
5.1	Notation	38
5.2	Numbers	38
5.2.1	Unsigned Integers	38
5.2.2	Signed Integers	39
5.2.3	Radix's Complement	39
5.2.4	Two's Complement Subtraction	40
5.2.5	Addition Is Unchanged from Unsigned	41
5.2.6	Sign Extension	41
5.2.7	Signed and Unsigned Instructions	41
5.3	Overflow	42
5.3.1	Overflow in 2's Complement	42
5.3.2	Overflow in Software	43
5.3.3	Detect Addition Overflow in Software	43
5.4	A Strange but Useful Property	43
5.4.1	Two's Complement Subtractor	44
5.4.2	Two's Complement Add/Subtract Unit	44

5.5	Bounds Check Optimization	45
5.6	Floating Point Representation	45
5.6.1	Sign-and-Magnitude Addition	47
6	Part II(a) - I/O - Exceptions Multicycle Processor W - 3.2, 4.1	49
6.1	Processor	49
6.1.1	Unified Memory	49
6.1.2	Single-Cycle Processor	50
6.2	Propagation Time	50
6.2.1	Increasing the Frequency	51
6.2.2	Two-Cycle Processor	51
6.2.3	Not All Paths Are Born Equal	51
6.2.4	Asynchronous/Synchronous Memories	52
6.3	Multicycle Processor	52
6.4	Mealy or Moore?	53
6.5	Processor - Building the Circuit	53
6.5.1	Adding the Instruction Register	54
6.5.2	Adding functionality	55
6.5.3	I-Type Instructions Need RF and ALU	55
6.5.4	R-Type Instructions and Second Operand Selection	56
6.5.5	And More, and More...	57
6.5.6	Guidelines for Writing Verilog	57
6.5.7	Detailing Complex Combinational Modules (ALU)	58
6.5.8	Verilog - Sticking to Basic Patterns	58
7	Part II(b) - Processor, I/Os, and Exceptions W - 4.1 - 4.2	59
7.1	The CPU	59
7.2	Physical Memory Map	60
7.2.1	Connecting CPU and Memory	60
7.3	Input/Output (I/O) Devices	61
7.3.1	Accessing I/Os: Port-Mapped I/O (PMIO)	61
7.3.2	Memory Mapped I/O (MMIO)	62
7.4	Example - A/D Converter	63
7.4.1	Bus Interface	63
7.4.2	Memory Mapping	63
7.4.3	Assembling everything	64
7.5	What do these tri-state buffers do?	65
7.5.1	A Classic UART	66
8	Part II(c) - Interrupts	67
8.1	I/O Polling	67
8.2	I/O Interrupts	67
8.2.1	The Basic Concept of I/O Interrupts	68
8.2.2	Interrupt Cycle Description	69
8.2.3	I/O Interrupt Priorities: Daisy Chain Arbitration	70
8.3	Direct Memory Access (DMA)	70
8.3.1	Timer and Interrupt Mechanism	72

9 Part II(d) - Processor, I/Os, and Exceptions	73
9.1 Exceptions, Interrupts, Faults, Traps, and Checks	73
9.1.1 Undefined Instruction	73
9.1.2 Optional <code>fadd.s</code> Instruction	74
9.1.3 Outline of an Undefined Instruction Handler	74
9.2 Exceptions and Interrupts	75
9.2.1 A Possible Classification of Exceptions	75

Chapter 1

Part I(a) - ISA Reminder, Assembly Language, Compiler - W 1.1

hum...welcome back

In the first part of the course, professor introduced (for motivational purposes) how computer architecture, specifically processors, have become essential to our lives, and how the field is growing exponentially. (didn't think it was essential to mention here...)

1.1 From High Level Languages to Assembly Language

1.1.1 High Level Languages

When talking about programming we usually think of programs that look like this...

```
1 int data = 0x00123456;
2 int result = 0;
3 int mask = 1;
4 int count = 0;
5 int temp = 0;
6 int limit = 32;
7 do {
8     temp = data & mask;
9     result = result + temp;
10    data = data >> 1;
11    count = count + 1;
12 } while (count != limit);
```

name	value
data	0x00123456
result	0
mask	1
count	...
temp	
limit	
...	
my_float	3.141529
a_string	Hello world!

1.1.2 Assembly Language

We use this code because it enables us to build a *Finite State Machine*, which isn't feasible with C code. This language provides a more rigid format with a sequence of numbered instructions, an *opcode*, predefined variable names, and the ability to **jump between lines**.

```

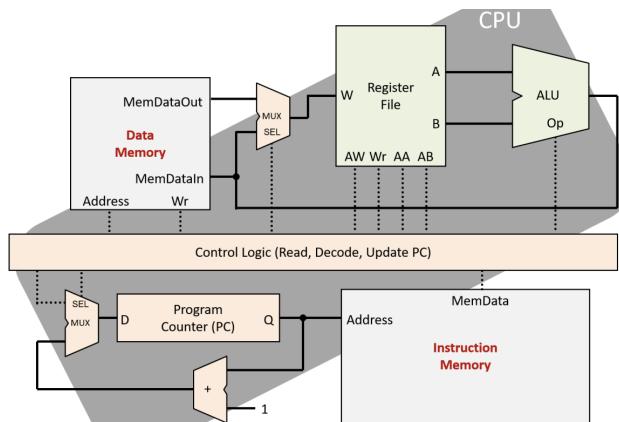
1 li x1, 0x000123456
2 li x2, 0
3 li x3, 1
4 li x4, 0
5 li x5, 0
6 li x6, 32
7 loop: and x5, x1, x3
8     add x2, x2, x5
9     srl x1, x1, 1
10    addi x4, x4, 1
11    bne x4, x6, loop

```

1.2 Processors

Remember, a processor can be decomposed into five components:

- **ALU (Arithmetic and Logic Unit)**: Performs arithmetic and logical operations.
- **Register File**: Stores data temporarily for quick access during processing.
- **Memory**: Holds data and instructions needed by the processor.
- **Control Logic**: Directs the operation of the processor by coordinating the other components.
- **PC (Program Counter)**: Keeps track of the address of the next instruction to be executed.
- **Instruction Memory**: Stores the program instructions that the processor will execute.



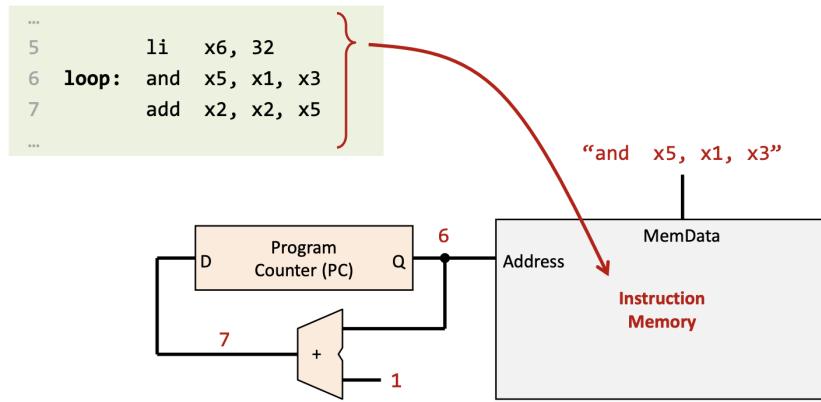
We may distinguish three types of general operations made by the processor:

Encoding

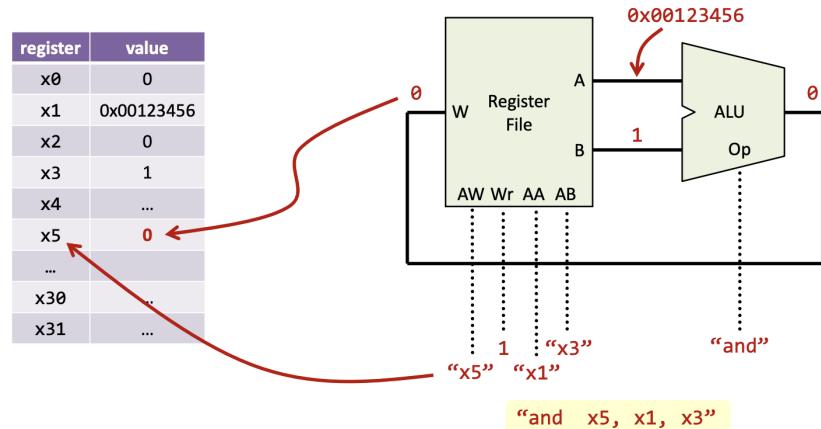
add x1, x1, x1	0 = 0000 0000 0000 0000 0000 0000 0000 0000
add x1, x1, x2	1 = 0000 0000 0000 0000 0000 0000 0000 0001
add x1, x1, x3	2 = 0000 0000 0000 0000 0000 0000 0000 0010
add x1, x1, x4	3 = 0000 0000 0000 0000 0000 0000 0000 0011
add x1, x1, x5	4 = 0000 0000 0000 0000 0000 0000 0000 0100
...	...
and x1, x1, x1	32768 = 0000 0000 0000 0000 1000 0000 0000 0000
and x1, x1, x2	32769 = 0000 0000 0000 0000 1000 0000 0000 0001
and x1, x1, x3	32770 = 0000 0000 0000 0000 1000 0000 0000 0010
and x1, x1, x4	32771 = 0000 0000 0000 0000 1000 0000 0000 0011
and x1, x1, x5	32772 = 0000 0000 0000 0000 1000 0000 0000 0100
...	...

of opcodes x # destinations x # source 1 x # source 1 ≤ 2³² combinations

Fetching



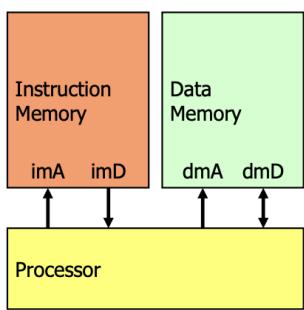
Executing



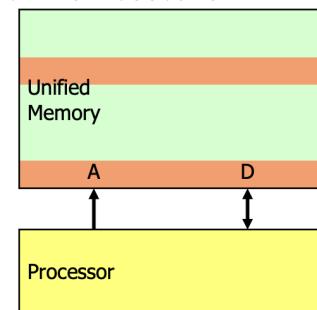
1.3 Joint or Disjoint Program and Data Memories

There are two main types of architectures one called the Harvard Architecture (Where the data and the memory are separate) and one called Unified Architecture (where data is shared with the program memory)

Harvard Architecture



Unified Architecture

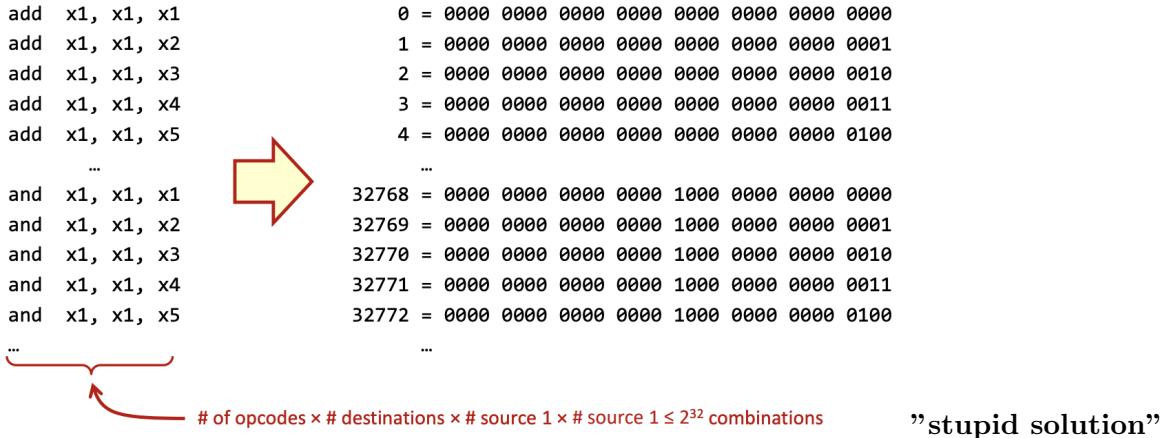


1.4 The Encoding problem

We may ask ourselves how we encode assembly written instructions into actual 0s and 1s.

1.4.1 The Stupid Solution

Now, the professor throws out the "stupid idea" (his words) of just counting all possible instructions, assigning a number to each one, and writing the numbers in binary. The problem with such a method is that the number of instructions could grow exponentially, requiring an unmanageable number of bits to represent each one, leading to inefficiency.



1.4.2 RISC-V Encoding (The Solution)

Instead, the chosen solution is to use an instruction set encoding where instructions are grouped into classes, each with a fixed format optimizing both memory usage and processing speed by limiting the number of bits required to represent instructions.

Instruction	Pseudocode	Type	funct7	funct3	opcode
Shift					
sll rd,rs1,rs2	rd ← rs1 ≪ rs2	R	0x00	0x1	0x33
slli rd,rs1,imm	rd ← rs1 ≪ imm	I	0x00	0x1	0x13
srl rd,rs1,rs2	rd ← rs1 ≫ _u rs2	R	0x00	0x5	0x33
srli rd,rs1,imm	rd ← rs1 ≫ _u imm	I	0x00	0x5	0x13
sra rd,rs1,rs2	rd ← rs1 ≫ _s rs2	R	0x20	0x5	0x33
srail rd,rs1,imm	rd ← rs1 ≫ _s imm	I	0x20	0x5	0x13
Arithmetic					
add rd,rs1,rs2	rd ← rs1 + rs2	R	0x00	0x0	0x33
addi rd,rs1,imm	rd ← rs1 + sext(imm)	I		0x0	0x13
sub rd,rs1,rs2	rd ← rs1 - rs2	R	0x20	0x0	0x33
lui rd,imm	rd ← imm				
auipc rd,imm	rd ← pc				
Logical					
xor rd,rs1,rs2	rd ← rs1	R	funct7	rs2	rs1 funct3 rd opcode
xori rd,rs1,imm	rd ← rs1	I		imm[11:0]	rs1 funct3 rd opcode
or rd,rs1,rs2	rd ← rs1	I	funct7	imm[4:0]	rs1 funct3 rd opcode
ori rd,rs1,imm	rd ← rs1	S	imm[11:5]	rs2	rs1 funct3 imm[4:0] opcode
and rd,rs1,rs2	rd ← rs1	B	imm[12-10:5]	rs2	rs1 funct3 imm[4:1-11] opcode
andi rd,rs1,imm	rd ← rs1	U			rd opcode
		J	imm[20-10:1-11-19:12]		rd opcode

Instruction types

31	25	24	20	19	15	14	12	11	7	6	0
R	funct7		rs2	rs1	funct3		rd		opcode		
I		imm[11:0]		rs1	funct3		rd		opcode		
I	funct7		imm[4:0]	rs1	funct3		rd		opcode		
S	imm[11:5]		rs2	rs1	funct3	imm[4:0]		opcode			
B	imm[12-10:5]		rs2	rs1	funct3	imm[4:1-11]		opcode			
U		imm[31:12]					rd		opcode		
J	imm[20-10:1-11-19:12]						rd		opcode		

Register-Register
Register-Immediate
Register-Immediate Shift
Store
Branch
Upper Immediate
Jump

RISC-V encoding

1.4.3 Automating this process

Now to automate the processes of decoding assembler code into machine code we use an **Assembler**, and to automate the process of decoding a higher level language to assembler we use a **Compiler**.

Assembler

The program that does this is called an assembler. It takes the assembly code and converts it into machine code.

```

0      li    x1, 0x00123456
1      li    x2, 0
2      li    x3, 1
3      li    x4, 0
4      li    x5, 0
5      li    x6, 32
6  loop: and  x5, x1, x3
7      add  x2, x2, x5
8      srl  x1, x1, 1
9      addi x4, x4, 1
10     bne  x4, x6, loop

```



```

0101 0101 0101 0000 0100 0111 1010 1110
0001 0100 1001 1101 0011 0000 1100 1001
1101 1100 1101 0110 0000 1101 0001 0111
0010 0011 1101 0110 0010 0000 0001 1001
1100 1010 1011 1010 0111 0100 0000 0110
1111 0010 1001 0011 1001 1110 1001 1101
0011 0000 0010 0111 1111 0000 0100 0011
0111 1001 0101 1101 1000 1000 0111 1011
1100 1010 1011 0000 0100 0100 0110 0101
0111 1001 0010 0110 0000 0011 0001 0010
0101 1100 1000 0101 0000

```

A fairly trivial job

Assembly

Compiler

A compiler is a program that translates high-level source code written in languages like C or Java into machine code or an intermediate representation.

```

int data  = 0x00123456;
int result = 0;
int mask   = 1;
int count  = 0;
int temp   = 0;
int limit  = 32;
do {
    temp  = data & mask;
    result = result + temp;
    data   = data >> 1;
    count  = count + 1;
} while (count != lim

```

A pretty hard job!...

```

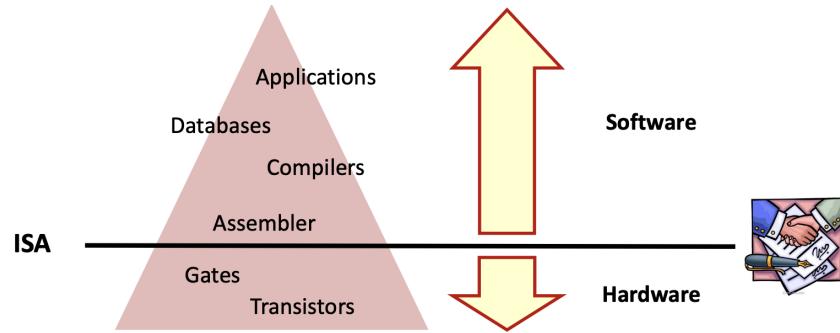
0      li    x1, 0x00123456
1      li    x2, 0
2      li    x3, 1
3      li    x4, 0
4      li    x5, 0
5      li    x6, 32
6  loop: and  x5, x1, x3
7      add  x2, x2, x5
8      srl  x1, x1, 1
9      addi x4, x4, 1
bne  x4, x6, loop

```

Compilation

1.5 ISA (Instruction Set Architecture)

The ISA is the interface between the hardware and the software. It defines the instructions that a processor can execute, as well as the format of those instructions.



Chapter 2

Part I(b) - ISA, Functions, and Stack - W 1.2

2.1 Arithmetic and Logic Instructions in RISCV

Below some examples of RISCV instructions:

Two Operands Instructions

```
1 sll  x5, x5, x9
2 add  x6, x5, x7
3 xor  x6, x6, x8
4 slt  x8, x6, x7
```

Shift $x5$ left by $x9$ positions $\rightarrow x5$
Add $x5$ and $x7 \rightarrow x6$
Logic XOR bitwise $x6$ and $x8 \rightarrow x6$
Set $x8$ to 1 if $x6$ is lower than $x7$, otherwise to 0

Arithmetic Instructions

```
1 slli x5, x5, 3
2 addi x6, x5, 72
3 xori x6, x6, -1
4 slti x8, x6, 321
```

Shift $x5$ left of 3 positions $\rightarrow x5$
Add 72 to $x5 \rightarrow x6$
Logic XOR bitwise $x6$ and 0xFFFFFFFF $\rightarrow x6$
Set $x8$ to 1 if $x6$ is lower than 321, to 0 otherwise

Here, you may ask yourself, why are all immediates (constants) written on a maximum of 12bits?

2.1.1 Constants must be encoded on 12 bits

As you may see here, all instructions encode immediates on 12 bits.

	31	25	24	20	19	15	14	12	11	7	6	0	
R	funct7		rs2		rs1		funct3		rd		opcode		Register-Register
I		imm[11:0]			rs1		funct3		rd		opcode		Register-Immediate
I	funct7		imm[4:0]		rs1		funct3		rd		opcode		Register-Immediate Shift
S	imm[11:5]		rs2		rs1		funct3		imm[4:0]		opcode		Store
B	imm[12–10:5]		rs2		rs1		funct3		imm[4:1–11]		opcode		Branch
U		imm[31:12]							rd		opcode		Upper Immediate
J		imm[20–10:1–11–19:12]							rd		opcode		Jump

2.1.2 Assembler Directives

Assembler directives help write cleaner and more readable code. The code snippets on the left and right below are equivalent.

<pre> lui x5, 0x12345 addiu x5, x5, 0x678 xor x6, x6, x5 </pre>		<pre> .equ something, 0x12345678 lui x5, %hi(something) addiu x5, x5, %lo(something) xor x6, x6, x5 </pre>
--	--	--

The left-hand side code snippet shows an assembly sequence where a 32-bit constant value (0x12345678) is loaded into a register (x5). Since immediate values are 16-bit limited, this requires splitting the 32-bit value into two instructions:

- The first instruction, `lui`, loads the upper 20 bits (0x12345) into the register `x5`.
- The second instruction, `addiu`, adds the lower 12 bits (0x678) to `x5`, completing the full 32-bit value in the register.

This approach, while functional, can become cumbersome when dealing with multiple constants, making the code less readable and harder to maintain.

The right-hand side shows the same functionality but makes use of assembler directives, specifically the `.equ` directive to define a label (`something`) for the constant 0x12345678. Using the `%hi()` and `%lo()` pseudo-instructions, the assembler automatically splits the constant into its upper and lower parts:

- The `%hi(something)` loads the upper 20 bits into `x5`.
- The `%lo(something)` adds the lower 12 bits to `x5`.

This method enhances code clarity and maintainability, especially when working with multiple constants, by using human-readable labels instead of raw numeric values. The assembler handles the details of splitting the 32-bit constant into its upper and lower parts.

Directive	Effect
<code>.text</code>	Store subsequent instructions at next available address in <i>text</i> segment
<code>.data</code>	Store subsequent items at next available address in <i>data</i> segment
<code>.asciiz</code>	Store string followed by null-terminator in <code>.data</code> segment
<code>.byte</code>	Store listed values as 8-bit bytes
<code>.word</code>	Store listed values as 32-bit words
<code>.equ</code>	Define constants

2.1.3 The x0 Register

The `x0` register is hardwired to 0 and cannot be changed. Any attempt to write into `x0` will have no effect.

Why is this useful?

One common application is in introducing wait delays during program execution. By leveraging the fixed nature of `x0`, it simplifies certain instructions that require an immediate zero value.

2.2 PseudoInstructions

PseudoInstructions simplify commands involving the `x0` register by creating easier-to-use alternatives.

Pseudoinstruction	Base Instruction(s)	Meaning
nop	addi x0, x0, 0	No operation
li rd, immediate	Myriad sequences	Load immediate
mv rd, rs	Myriad sequences	Copy register
not rd, rs	xori rd, rs, -1	One's complement
neg rd, rs	sub rd, x0, rs	Two's complement
seqz rd, rs	sltiu rd, rs, 1	Set if = zero
snez rd, rs	sltu rd, x0, rs	Set if \neq zero
sltz rd, rs	slt rd, rs, x0	Set if \downarrow zero
sgtz rd, rs	slt rd, x0, rs	Set if \downarrow zero

The term *myriad sequences* refers to a series of instructions that together achieve the functionality of a single pseudoinstruction, such as using lui and addi to implement li rd, immediate. According to the professor li should be called mvi (as move immediate).

2.2.1 Control flow instructions

Control flow instructions are used to change the order of execution of instructions are a kind of pseudo-instructions.

```

1  li x1, 0x00123456
2  li x2, 0
3  li x3, 1
4  li x4, 0
5  li x5, 0
6  li x6, 32
7  loop: and x5, x1, x3
8    add x2, x2, x5
9    srli x1, x1, 1
10   addi x4, x4, 1
11   bne x4, x6, loop

```

2.2.2 If-Then-Else

```

1  if (x5 == 72) {
2    x6 = x6 + 1;
3  } else {
4    x6 = x6 - 1;
5 }

```

```

1  .text
2    li x7, 72
3    beq x5, x7, then_clause
4    else_clause:
5      addi x6, x6, -1
6      j end_if
7    then_clause:
8      addi x6, x6, 1
9    end_if:

```

As seen here, beqi does not exist in RISCV, instead we use beq and li to achieve the same result.

2.2.3 Jumps and Branches

A common but not universal distinction exists between *jumps* and *branches*. In RISC-V (inherited from MIPS and used by SPARC, Alpha, etc.), jumps refer to unconditional control transfer instructions, while branches refer to conditional control transfer instructions. However, not all architectures follow this convention. For instance, in x86, all control transfer instructions are considered jumps, such as JMP, JZ, JC, and JNO.

2.2.4 Comparaisons

The processor implements only $<$ and $>$, and the assembler “creates” \leq and \geq .

Pseudoinstruction	Base Instruction(s)	Meaning
beqz rs, offset	beq rs, x0, offset	Branch if = zero
bnez rs, offset	bne rs, x0, offset	Branch if \neq zero
blez rs, offset	bge x0, rs, offset	Branch if \leq zero
bgez rs, offset	bge rs, x0, offset	Branch if \geq zero
bltz rs, offset	blt rs, x0, offset	Branch if $<$ zero
bgtz rs, offset	blt x0, rs, offset	Branch if $>$ zero
.bgt rs, rt, offset	blt rt, rs, offset	Branch if $>$
.ble rs, rt, offset	bge rt, rs, offset	Branch if \leq
.bgtu rs, rt, offset	bltu rt, rs, offset	Branch if $>$, unsigned
.bleu rs, rt, offset	bgeu rt, rs, offset	Branch if \leq , unsigned

2.2.5 Do-While

Do-while loops look like this (we obviously use control flow instructions here).

```

1 do {
2     x5 = x5 >> 1;
3     x6 = x6 + 1;
4 } while (x5 != 0);

```

```

1 .text
2 loop:
3     srl x5, x5, 1
4     add x6, x6, 1
5     bnez x5, loop

```

2.3 Functions

In higher-level programming languages, functions (routines, subroutines, procedures, methods, etc.) are used to encapsulate code and make it reusable.

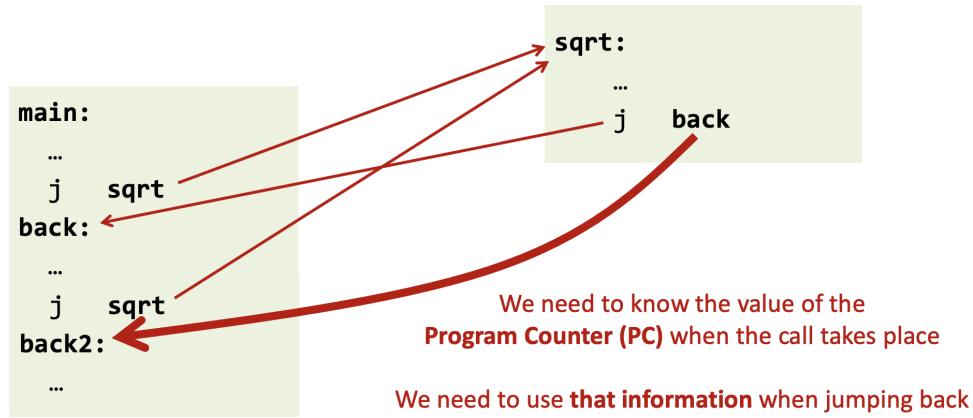
Calling a function involves these steps:

1. Place arguments where the called function can access them.
2. Jump to the function.
3. Acquire storage resources the function needs.
4. Perform the desired task of the function.
5. Communicate the result value back to the calling program.
6. Release any local storage resources.
7. Return control to the calling program.

2.3.1 Jump to the Function/Retun control to the calling program

The too simple not working approach

A simple (not working) approach for creating functions would be to do this:



With this approach the function doesn't know where to return to after being called (`back2` or `back`). For the next part, remember, the Program Counter is distinct from general-purpose registers. It is dedicated to managing the flow of instruction execution, while general registers are used for data manipulation.

The Good Approach

The right approach involves using the *Jump and Link* instruction `jal`, here loading $PC + 4$ (remember 4 bytes per Instruction) into `x1` as a way to come back from the function.

<pre> 1 main: 2 ... 3 jal x1, sqrt 4 ... 5 ... 6 jal x1, sqrt </pre>	<pre> 1 sqrt: 2 ... 3 ... 4 jr x1 </pre>
--	--

Both times `x1` was used to store the return address, and there is a reason for that (Register Conventions Sections).

2.3.2 Jump Instructions

There are only two core real jump instructions in RISCV, `jal` (jump and link) and `jalr` (jump and link register), the rest are pseudo instructions using them.

Pseudoinstr.	Base Instruction(s)	Meaning
<code>j offset</code>	<code>jal x0, offset</code>	Jump
<code>jal offset</code>	<code>jal x1, offset</code>	Jump and link
<code>jr rs</code>	<code>jalr x0, 0(rs)</code>	Jump register
<code>jalr rs</code>	<code>jalr x1, 0(rs)</code>	Jump and link register
<code>ret</code>	<code>jalr x0, 0(x1)</code>	Return from subroutine

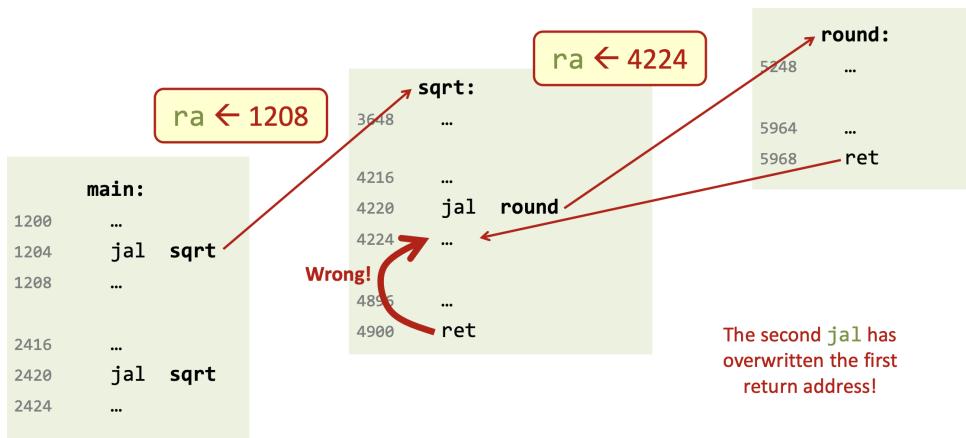
2.3.3 Register Conventions

Register conventions are rules that dictate how registers are used in a program, here are the ones we've seen for now

Register	Mnemonic	Description
x0	zero	Hard-wired zero
x1	ra	Return Address

2.3.4 Back to the good (not so good) approach

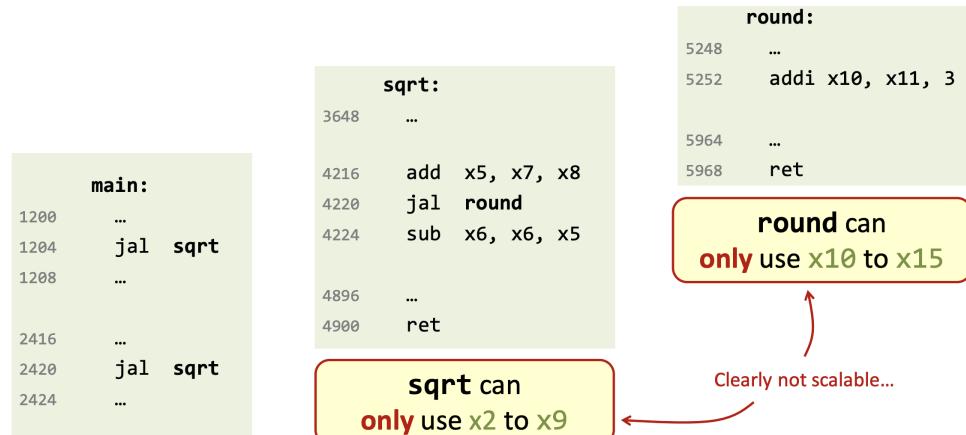
There's still a problem with the previous approach, say for example you want to call a function from another function.



Here the allocated space for the return address is overwritten by the second function call, and the first function can't return to the right place.

2.3.5 One simple solution (still not good)

One solution would be to say that a range of registers are used for certain functions and that they can't be used by other functions.



The problem here is that it's still not very scalable.

2.3.6 Acquire storage resources the function needs (still not it)

One simple solution to our problem would be to allocate memory for the function at in the data section of the program.

```

1 .data
2 sqrt_save_ra: .word 0
3 sqrt_save_x5: .word 0

```

```

1 .text
2 sqrt:
3 ...
4 add x5, x7, x8
5 sw ra, sqrt_save_ra
6 sw x5, sqrt_save_x5
7 jal round
8 lw ra, sqrt_save_ra
9 lw x5, sqrt_save_x5
10 sub x6, x6, x5
11 ...
12 ret

```

Problem: Recursive Functions

The problem here is that the return address is overwritten by the recursive call.

```

1 .data
2     find_child_save_ra: .word 0
3 .text
4     find_child:
5     ...
6     sw ra, find_child_save_ra
7     jal find_child
8     lw ra, find_child_save_ra
9     ...
10    ret

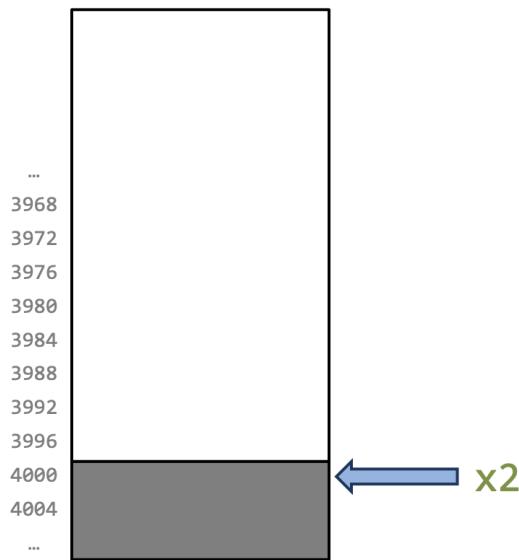
```

2.3.7 The Stack

The Solution to our Problem is this, the Stack.

The Stack is a region of memory that grows and shrinks as needed.

We may use a register (e.g x2) to point to the first used word after the end of the used region.



Dynamic Memory Allocation

The Stack, contrary to the Data Section, is dynamic and can be used to allocate memory when needed. This means that during program execution, variables or temporary data can be stored in the stack, which grows or shrinks depending on the operations performed.

The **stack pointer**, typically register x2, is used to manage the allocation and deallocation of memory.

In this instruction, for example, we allocate 12 bytes in the stack. We achieve this by decrementing the stack pointer (x2) by 12. This ensures that the new memory space is available for temporary storage.

```
1 addi x2, x2, -12
```

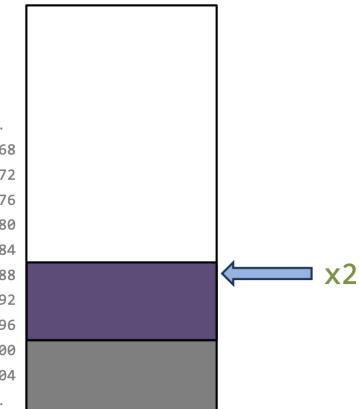


Retrieving Data from the Stack

Once memory has been allocated on the stack, we can store or retrieve data from it. In this case, we are retrieving data that was previously saved in the stack. The lw (load word) instruction is used to load the values stored at different offsets in the stack.

In this case, we retrieve three different values from the stack using the lw instruction, which loads a 4-byte value into the specified registers (ra, x5, and x6). The offsets (0, 4, and 8) refer to different positions in the 12 bytes we allocated earlier.

```
1 lw ra, 0(x2)
2 lw x5, 4(x2)
3 lw x6, 8(x2)
```



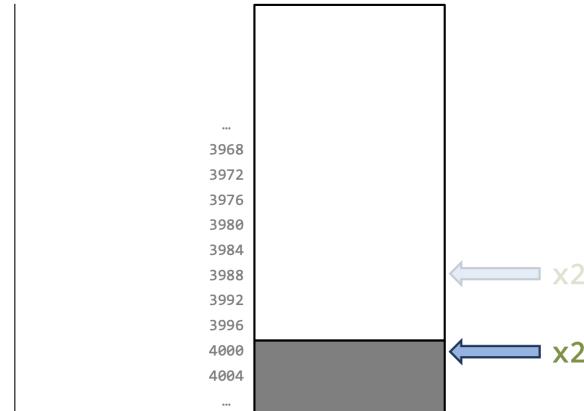
Memory Deallocation

After the data has been used or is no longer needed, it is good practice to deallocate the memory to ensure proper management of the stack. We deallocate memory by adjusting the stack pointer ($x2$) back to its original position.

In this instruction, we restore the stack to its previous state by adding 12 back to the stack pointer ($x2$).

This effectively "frees" the 12 bytes of memory we had allocated earlier.

1 addi x2, x2, 12



The Stack Pointer

The Stack Pointer is a register that points to the top of the stack, by convention it corresponds to the $x2$ register

Register	ABI Name	Description	Preserved across call?
$x2$	sp	Stack pointer	Yes

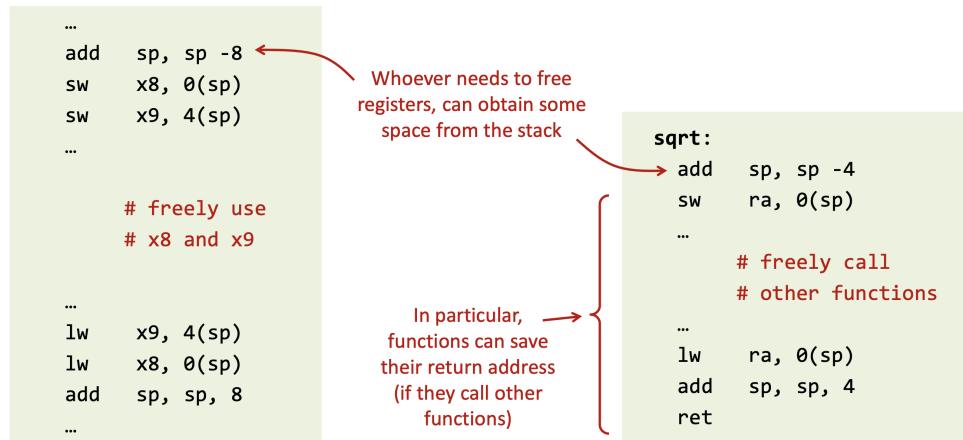
Other architectures have special instructions to place stuff on the stack (push) and to retrieve it (pop)

PUSH AX

1 add sp, sp, -4
2 sw x5, 0(sp)

2.3.8 Spilling Registers to Memory

Spilling registers to memory involves saving register values to the stack when more registers are needed or to prevent overwriting important data, allowing the registers to be reused. This technique is also used in function calls to save the return address, ensuring the program can correctly return control after the function finishes.



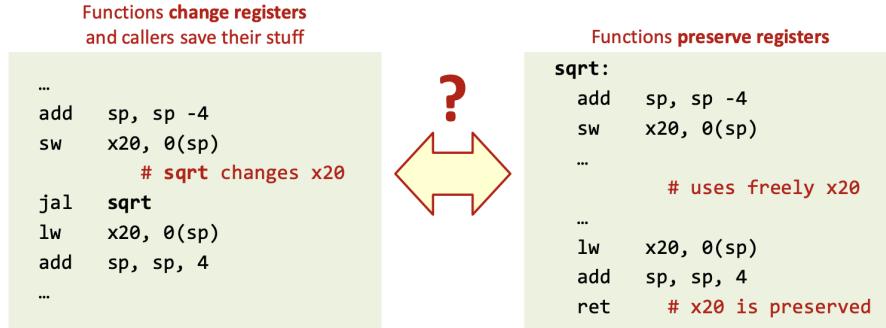
2.3.9 Register across functions

In assembly programming, handling registers across functions can be managed in two main ways: either functions **change registers** and expect the caller to save their values, or functions **preserve registers** and ensure that the register values remain the same across function calls.

- On the left, the function `sqrt` changes the value of register `x20`, requiring the caller to save and restore its value.
- On the right, the function `sqrt` preserves the value of `x20`, ensuring that the caller does not need to manage the saving and restoring.

This distinction is important, but it does not cause issues as long as there is agreement on how registers are handled.

In case it's still not clear, we're looking at the `sw` instruction



2.3.10 Preserving Registers

In RISC-V, register preservation is managed through a combination of callee-saved and caller-saved registers. Callee-saved registers (such as `s0`, `s1`, and `s2-11`) are preserved by the called function, ensuring that their values remain unchanged after the function call.

Caller-saved registers (such as `t0`, `t1-2`, and `t3-6`) are temporary and do not need to be preserved by the called function, meaning the caller must save them if their values are important.

Register	ABI Name	Description	Preserved across call?
x0	zero	Hard-wired zero	—
x1	ra	Return address	No
x2	sp	Stack pointer	Yes
x5	t0	Temporary/alternate link register	No
x6-7	t1-2	Temporaries	No
x8	s0/fp	Saved register/frame pointer	Yes
x9	s1	Saved register	Yes
x18-27	s2-11	Saved registers	Yes
x28-31	t3-6	Temporaries	No

2.4 Passing Arguments in RISC-V

In RISC-V, there are two main ways to pass arguments to functions:

2.4.1 Option 1: Using Registers

- Specific registers are used to pass arguments and return results.
- This can be done in a straightforward way, where each function uses different registers (e.g., passing an argument in `x5` and returning the result in `x6`).
- A more structured approach is to follow a convention where arguments are passed in registers `x10` to `x17`, with results returned in `x10`.
- The limitation: if there are more arguments than available registers (e.g., more than 8 arguments), this approach is insufficient.

2.4.2 Option 2: Using the Stack

- When registers are not enough, extra arguments are placed on the stack.
- The stack offers a universal solution because it has no practical limit on size.
- However, using the stack is more complex and requires additional work compared to using registers.

2.4.3 The RISC-V Approach

- RISC-V uses a combination of both methods.
- Registers x10 to x17 are used to pass arguments, with x10 and x11 also handling return values.
- If more arguments are needed beyond what these registers can handle, they are passed via the stack.

Register	ABI Name	Description	Preserved across call?
x10–11	a0–1	Function arguments/return values	No
x12–17	a2–7	Function arguments	No

Register reserved for arguments and return values in RISC-V.

2.5 Summary of RISC-V Register Conventions



Register	ABI Name	Description	Preserved across call?
x0	zero	Hard-wired zero	—
x1	ra	Return address	No
x2	sp	Stack pointer	Yes
x3	gp	Global pointer	—
x4	tp	Thread pointer	—
x5	t0	Temporary/ alternate link register	No
x6–7	t1–2	Temporaries	No
x8	s0/fp	Saved register/ frame pointer	Yes
x9	s1	Saved register	Yes
x10–11	a0–1	Function arguments/return values	No
x12–17	a2–7	Function arguments	No
x18–27	s2–11	Saved registers	Yes
x28–31	t3–6	Temporaries	No

Chapter 3

Part I(c) - ISA Memory and Addressing Modes - W 2.1

3.1 Memory

Memory is a really important component of a computing system, we store our programs in it, we store our data in it, and it's through memory that we receive and send data.

Though memory is very useful it has three main drawbacks:

- It's **slow** → Caches
- It's **finite** → Virtual Memory
- It can make an ISA **too complex** → Pipelining

no worries we'll cover each one of these in this chapter.

3.1.1 Address and Data

Data in Memory can be accessed by an address, meaning it's a Random Access (it can access a memory value without going through the preceding ones).

Professor Remark: "There's not anything random about this memory, we'd better call it and arbitrary access memory. (!not and official name)"

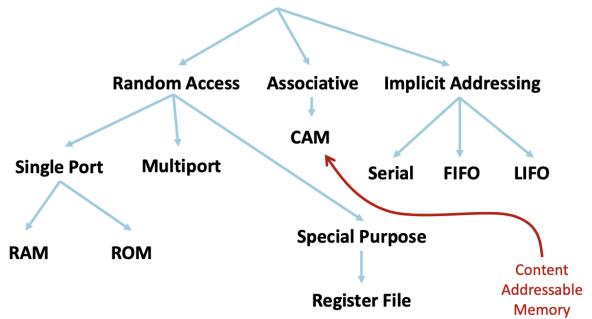
Address	Value
0	12
1	6
2	4
3	1
4	0
5	3
6	1
7	13
8	15
9	9
10	3
11	5
12	0
13	0
14	0
15	0

Write	Read
Memory[5] = 3	Memory[5]?

3.2 Many Types of Memories

We may distinguish between different types of memories based on their **technology**, such as SRAM, DRAM, EPROM, and Flash, and their **capabilities**, including **speed**, **capacity**, **density**, **writability** (whether they are writable, permanent, or reprogrammable), as well as their **size**, **volatility**, and **cost**.

3.2.1 Functional Taxonomy of Memories

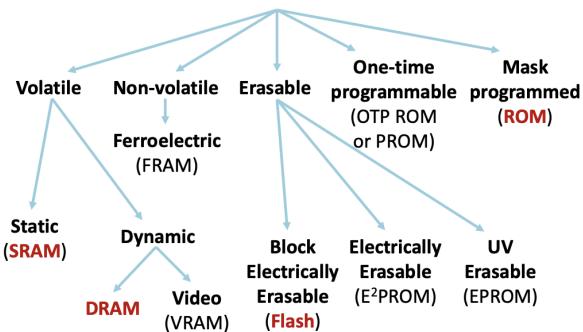


Multiport memory allows simultaneous access by multiple processors, while **single-port** memory supports only one at a time.

Non-Random Access memories

- **Associative** memories enable fast data retrieval by content rather than address, making it useful for cache memory, pattern recognition, and efficient lookups in large datasets.
 - In **Implicit addressing** the address of the data to be operated on is inferred directly by the operation code (opcode), without explicitly specifying the address in the instruction.

3.2.2 Taxonomy of Random Access Memories

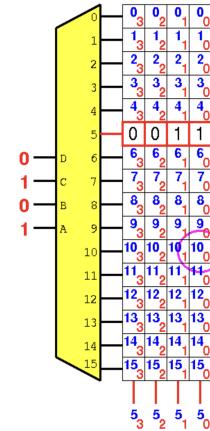


3.2.3 Basic Structure

Remember, a Data Flip Flop, stores a 1 bit value by updating the output value to the input value at the rising edge of the clock signal.

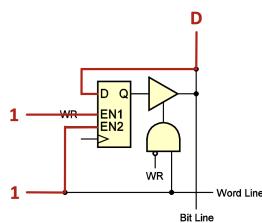
Address	Value
0	12
1	6
2	4
3	1
4	0
5	3
6	1
7	13
8	15
9	9
10	3
11	5
12	0
13	0
14	0
15	0

16 x 4 Memory Cells (Special DFFs (Data Flip-Flops))



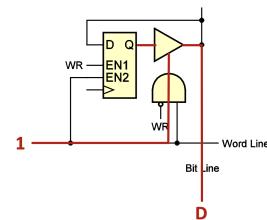
3.2.4 Write Operations

The *D* is connected to the Data outside of the system and at the rising edge it updates the value of the DFF. The AND gate ensures that the write signal is high when the clock signal is high.



3.2.5 Read Operations

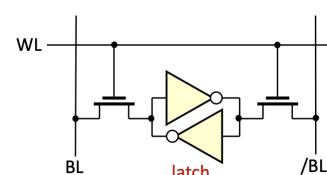
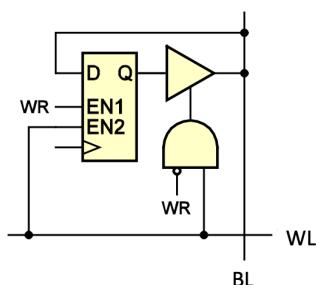
D is still connected to the Data, remember the tri-state driver is active when it's enable signal is active (so when the wr is off and the operation signal is sent.).



3.2.6 Practical SRAMs

DISCLAIMER !!: Combinational loops are prohibited as they can lead to unstable behavior, unpredictable timing, simulation and synthesis issues, excessive power consumption, and lack of a defined reset state, making them unsuitable for reliable digital circuit design.

While the type of memory we've just seen is small, and very fast, SRAM memories uses 6 transistors per cell (less than the previous design). We've also seen (in Taxonomy) that SRAM is static meaning it doesn't require periodic refresh.

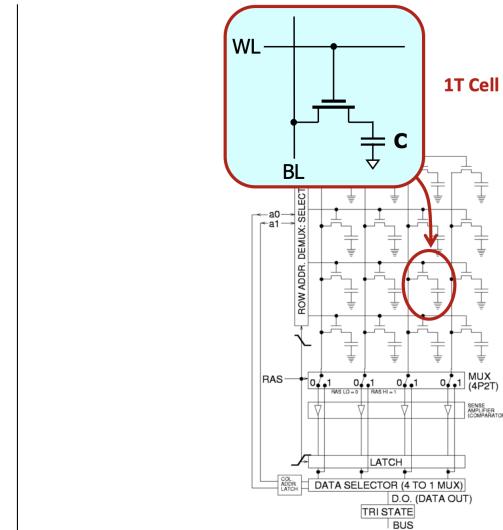


3.2.7 DRAMs

Dynamic RAMS(DRAMs) are the densest and cheapest type of RAM memory, it stores information as charge in small capacitors. This makes the DRAM need periodic refresh otherwise the charge might leak off (60ms) the capacitor due to parasitic resistances and the information lost

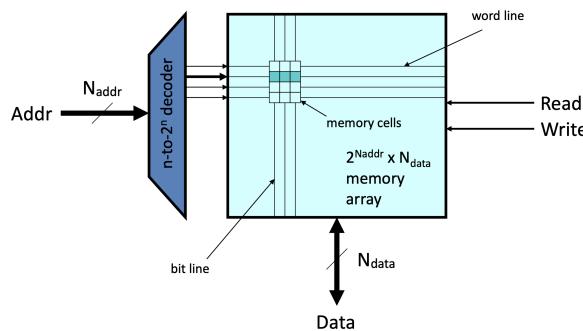
Refresh means, we come back before the end of a charge (60ms) and we rewrite the value, if there is still some charge, we add charge, if there's no charge and we keep as is.

Personal Remark: Dynamic = Bad, data disappears and needs refresh

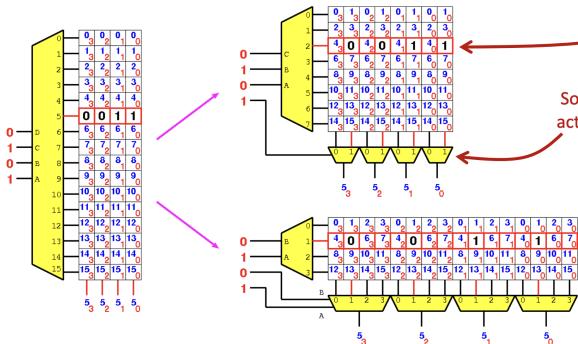


3.2.8 Ideal Random Access Memory

A memory array uses an n -to- 2^n decoder to select a word line based on the input address, enabling data to be read or written through the bit lines.



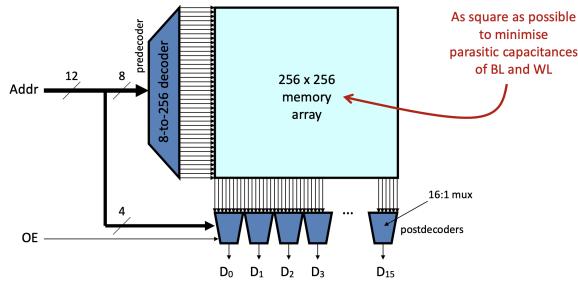
3.2.9 Physical Organisation



Out of all physical organizations, the squared one is the best one as it has the best performance. This layout facilitates faster access times and simplified wiring, resulting in improved computational efficiency and system scalability.

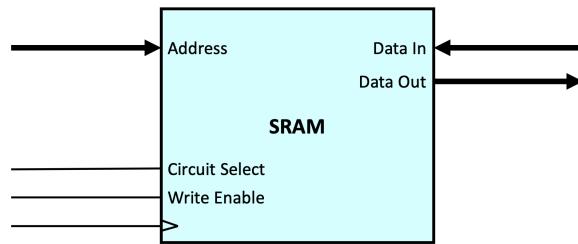
3.2.10 Realistic ROM Array

ROMs are Read-Only Memories, they are used to store the program of the computer, they are non-volatile and can't be written to.



3.2.11 Static Ram Typical Interface

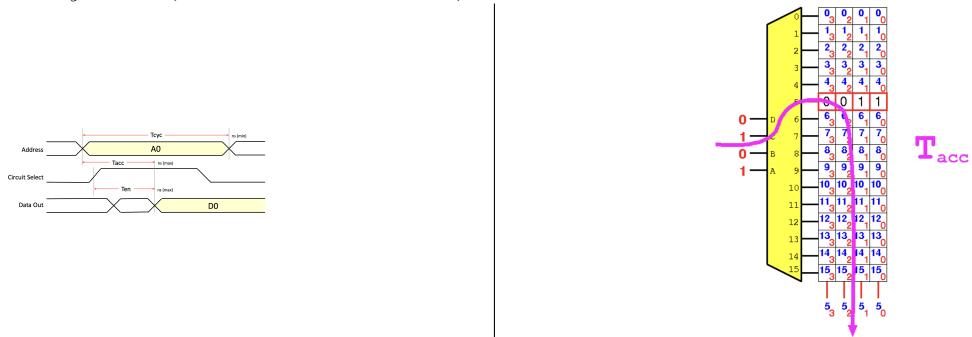
This a typical interface of a SRAM, it has a 16-bit data input/output, a 16-bit address input, a write enable signal, and a circuit select signal.



3.3 Typical Asynchronous SRAM Read Cycle

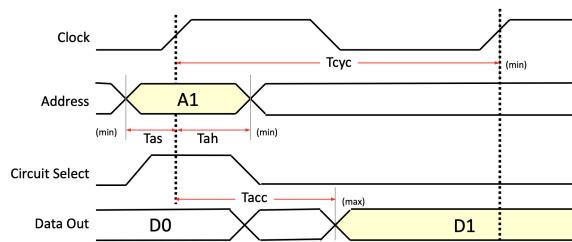
The read cycle of an asynchronous SRAM is initiated by the address input, which is decoded to select the word line, enabling the data to be read from the memory array and output to the data bus.

Here, Tcyc is the cycle time, Tacc is the access time, and Ten is the enable time.



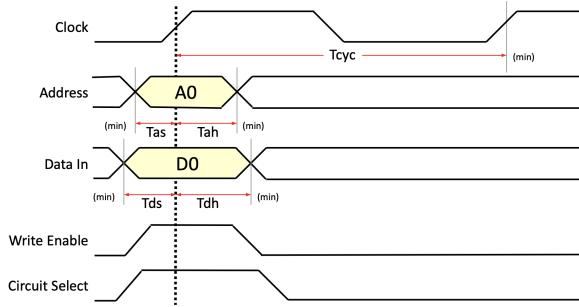
Read Cycle

Latency defined as the number of cycles between the address asserted and data available



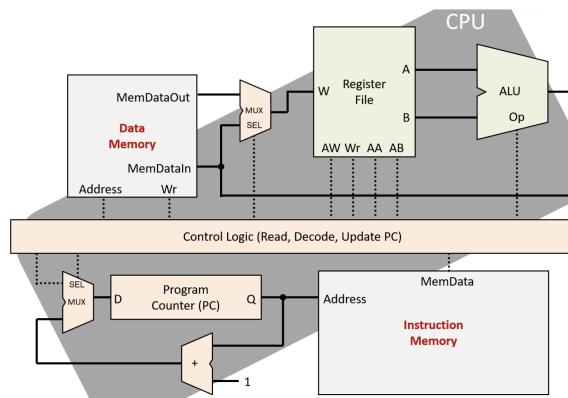
Write Cycle

Writes on the edge of the clock signal, as a DFF



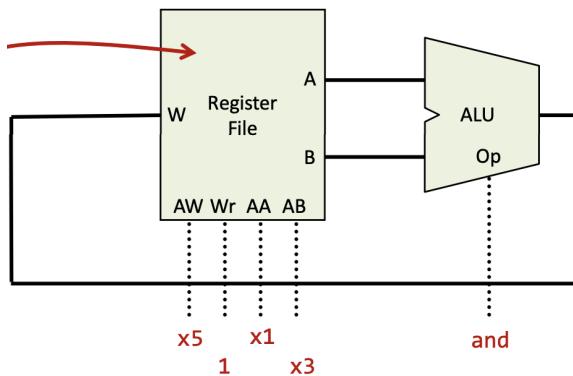
3.4 Where is Memory in the Processor?

In the processor we have memory in the Data memory component and in the Instruction memory component.

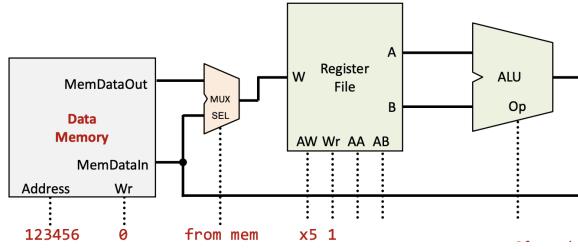


3.4.1 Arithmetic and Logic Instructions

The register file can only contain a limited number of registers making it difficult to handle more complex computations and managing data input/output efficiently.

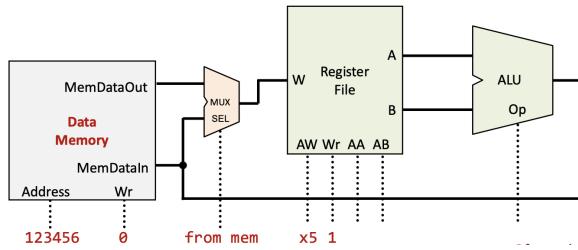


Load Instructions



Load and Store: The RISC-V Way

This instruction would never work for example because the address is too big to be sent as an immediate value : lw x5, (x7)



A Load/Store Architecture

A feature of RISC-V is that it's a Load/Store architecture, meaning that the only way to access memory is through load and store instructions. Also, instructions reading and writing in memory do exactly that and nothing else, contrary to more complex instruction set architectures (CISC), where instructions may combine memory access with other operations like arithmetic or logic. This simplicity in RISC-V's instruction set helps with streamlining the pipeline and improving performance efficiency.

Load		I	0x2	0x03
lw	rd, imm(rs1)	rd	$\leftarrow \text{mem}[\text{rs1} + \text{sext}(\text{imm})]$	
Store		S	0x2	0x23
sw	rs2, imm(rs1)	$\text{mem}[\text{rs1} + \text{sext}(\text{imm})] \leftarrow \text{rs2}$		

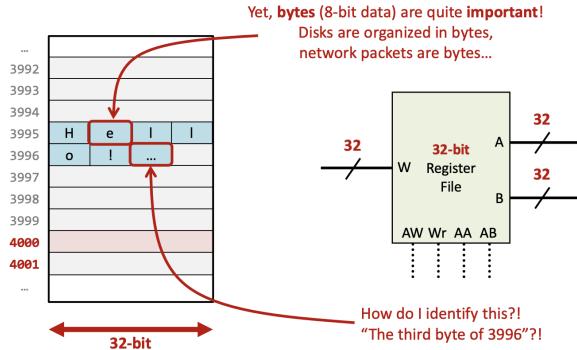
3.5 More Addressing Modes? Not in RISC-V!

Addressing Mode	Instruction	Description
Register	add x0, x1, x2	Adds the value of x1 and x2, stores the result in x0.
Immediate	add x0, x1, 123	Adds the value of x1 and the immediate constant 123, stores the result in x0.
Direct or Absolute	add x0, x1, (1234)	Adds the value of x1 and the value at memory address 1234, stores the result in x0.
Register Indirect	add x0, x1, (x2)	Adds the value of x1 and the value in memory at the address held in x2, stores in x0.
Displacement or Relative	add x0, x1, 123(x2)	Adds the value of x1 and the value in memory at x2 plus the displacement 123, stores in x0.
Base or Indexed	add x0, x1, i5(x2)	Adds the value of x1 and the value in memory at x2 plus index i5, stores in x0.
Auto-increment/-decrement	add x0, x1, (x2+)	Adds the value of x1 and the value in memory at the address in x2, then increments x2, stores in x0.
PC-Relative	add x0, x1, 123(pc)	Adds the value of x1 and the value in memory at pc plus 123, stores in x0.

Syntax here looks like RISC-V but most of these instructions do not exist in RISC-V.

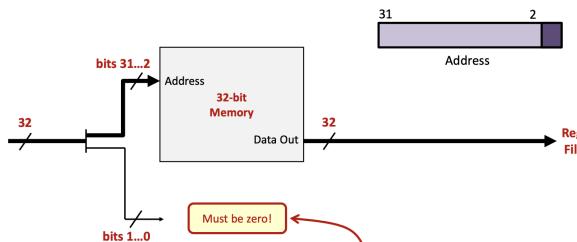
3.5.1 Word Addressed Memory

In a word addressed memory, the address is the index of the word in the memory.
The letters inside the word are identified as eg. for Hello World, H:3980, E:3981, L:3982,



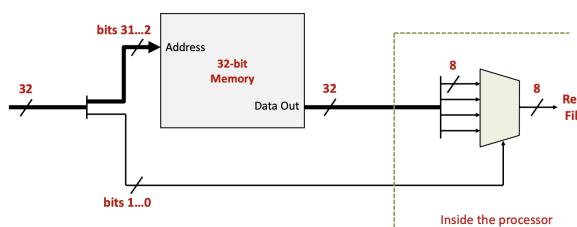
3.5.2 Loading Words (lw) and Instructions

The `lw` instruction is used to load a word from memory into a register.
The address of such words would necessarily be a multiple of 4 meaning the two least significant bits must be 0s. (to ensure the data is word aligned...)



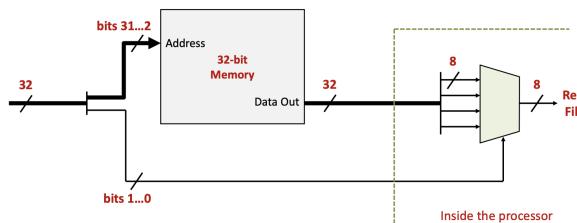
3.5.3 Loading Bytes (lb)

The `lb` (Load Byte) instruction doesn't require alignment because it only loads 1 byte (8 bits), which can be accessed at any memory address, unlike `lw` which requires word alignment to efficiently load 4 bytes (32 bits).
The `lb` instruction is used to load a byte from memory into a register.



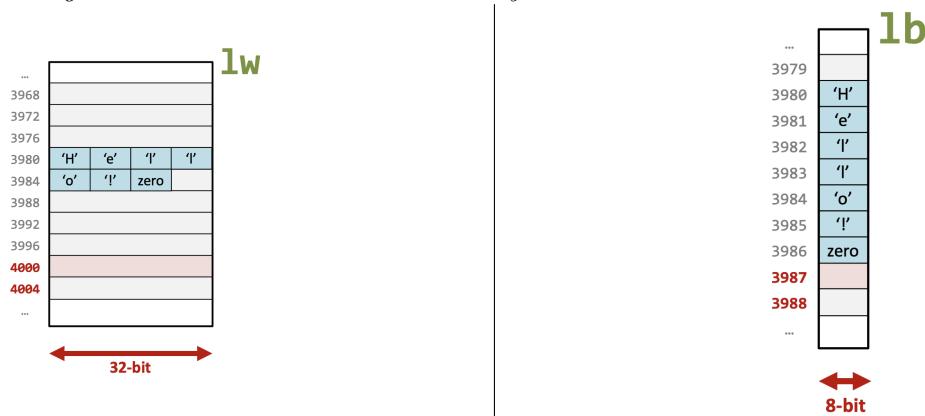
3.5.4 A Few More Load/Store Instructions

Access bytes (and half-words) as if memory were made of bytes



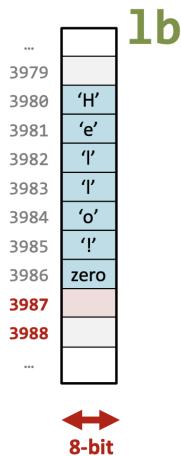
3.5.5 Access as it is more suitable

For example storing the "Hello!" zero value in the memory would like this:



Counting Characters in a String

As an example, for counting the number of characters in a string, the load byte instruction would be more suitable as seeing the string as a sequence of bytes makes use of the memory as a sort of array.



```

1  strlen:
2      mv t0, a0 # Copy the pointer (a0) into t0 to traverse the string
3      li t1, 0 # t1 will hold the length (initialized to 0)
4  loop:
5      lbu t2, 0(t0) # Load byte at address t0 into t2
6      beq t2, zero, end # If t2 is 0 (null byte), we are done
7      addi t1, t1, 1 # Increment the length counter (t1)
8      addi t0, t0, 1 # Point to the next character in the string
9  j loop # Repeat the loop
10 end:
11     mv a0, t1 # Move the length (t1) into a0 as the return value
12     ret # Return to caller

```

lbu is used here to ensure that the byte is treated as an unsigned value, which is the correct approach for processing characters in a string.

In a word addressed memory view, the code would look like such:

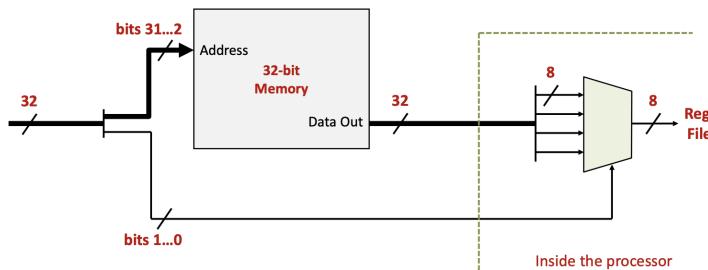
```

1  strlen:
2      li t1, 0          # t1 will hold the length (initialized to 0)
3  next_word:
4      li t2, 4          # t2 will count the bytes in a loaded word (four)
5      lw t3, 0(t0)      # Load four bytes at address t0 into t3
6  next_byte:
7      andi t4, t3, 0xff # Move the "little-end" in t4
8      beq t4, zero, end # If t4 is 0 (null byte), we are done
9      addi t1, t1, 1     # Increment the length counter (t1)
10     srli t3, t3, 8    # Prepare the next byte of the word in the "little-end" (t3)
11     addi t2, t2, -1    # One byte left in the loaded word
12     bneq t2, next_byte # If more bytes in t3, check the next
13     addi a0, a0, 4     # Else point to the next word of characters in the string
14     j next_word        # Repeat the loop
15 end:
16     mv a0, t1          # Move the length (t1) into a0 as the return value
17     ret                # Return to caller

```

3.5.6 Loading Bytes (lb)

Now, one may wonder in what ordering the bytes are stored in memory.



Which Byte Where?

Both ordering of bytes are valid the only thing we have to do is stick to one, the most generally used is little-endian as it's the RISC-V default and the Intel x86/x64 default.



Personal Remark : Mnemotecnic - Little Endian = Little End (The ending memory index takes the smallest(starting) data address), Big Endian = Big End.

Chapter 4

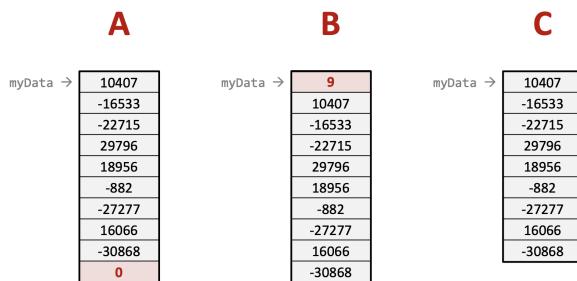
Part I(d) - ISA Arrays and Data Structures - W 2.2

4.1 Arrays

In higher level languages, are written like follows :

```
1 short[\[] myData = \{10407, -16533, -22715, 29796, 18956, \dots\}:
```

4.1.1 Different Ways to Store Arrays



- **A: Storing Arrays with a Null Terminator**

- In this method, the array is stored with each element represented using 16-bit integers.
- A null terminator (the value 0) is used at the end of the array to indicate its termination.
- This method is common when the array size is unknown in advance, and the null terminator acts as a signal to stop reading the data.

- **B: Storing Arrays with a Length Prefix**

- Here, the first element of the array contains the length of the array, stored as a 16-bit integer (in this case, the length is 9).
- The rest of the array is stored in consecutive memory locations, similar to method A.
- This method allows the array size to be known before reading all the data, making it more efficient for some use cases.

- **C: Storing Arrays without a Terminator or Length Prefix**

- In this case, the array is stored without a length prefix or a null terminator.
- The size of the array must be known externally, either through the code or an external mechanism.
- This method is the most compact but requires prior knowledge of the array's size.

4.1.2 Adding Positive Elements

Here we'll write the same program for the different ways of storing arrays.

The program will add all positive elements of an array of signed 16-bit integers. At call time, *a0* points to the array, at return time, *a0* contains the result.

A: Storing Arrays with a Null Terminator

A

myData →	10407
	-16533
	-22715
	29796
	18956
	-882
	-27277
	16066
	-30868
	0

```

1 add_pos: li t0, 0           # Initialize t0 to 0
2     lh t1, 0(a0)          # Load halfword from memory address a0 into t1
3     beqz t1, end          # Branch to 'end' if t1 equals zero
4     blez t1, donothing    # Branch to 'donothing' if t1 is less than or equal to
      zero
5     add t0, t0, t1        # Add t1 to t0 (only if t1 is positive)
6 doNothing:                 # This block does nothing for negative or zero values
7 # You can put other operations here if needed
8     j add_pos             # Jump back to the beginning of add_pos to check the next
      value
9 end:                      # Label 'end' for the program termination

```

B: Storing Arrays with a Length Prefix

B

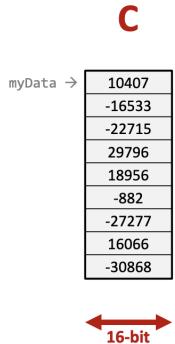
myData →	9
	10407
	-16533
	-22715
	29796
	18956
	-882
	-27277
	16066
	-30868

```

1 add_pos_b:
2     lh t2, 0(a0)          # Load the length of the array into t2
3     addi a0, a0, 2         # Move to the first element of the array (skip the length
      prefix)
4     li t0, 0               # Initialize t0 to 0 for storing the sum
5 loop_b:
6     beqz t2, end_b         # If the length (t2) is zero, branch to 'end_b'
7     lh t1, 0(a0)          # Load the current array element into t1
8     blez t1, skip_b        # If t1 is less than or equal to zero, skip the addition
9     add t0, t0, t1          # Add t1 to t0 (only if t1 is positive)
10 skip_b:
11     addi a0, a0, 2         # Move to the next element in the array
12     addi t2, t2, -1        # Decrease the length counter
13     j loop_b              # Jump back to loop_b
14 end_b:                  # End label

```

C: Storing Arrays without a Terminator or Length Prefix



```

1 add_pos_c:
2     li t0, 0           # Initialize t0 to 0 for storing the sum
3 loop_c:
4     beqz t2, end_c    # If the array size (t2) is zero, branch to 'end_c'
5     lh t1, 0(a0)       # Load the current array element into t1
6     blez t1, skip_c   # If t1 is less than or equal to zero, skip the addition
7     add t0, t0, t1     # Add t1 to t0 (only if t1 is positive)
8 skip_c:
9     addi a0, a0, 2     # Move to the next element in the array
10    addi t2, t2, -1    # Decrease the array size counter
11    j loop_c          # Jump back to loop_c
12 end_c:               # End label

```

4.1.3 Pointer to Memory vs Index in Array

Now we're wondering which one of these two ways of accessing the array is better.

Obviously the less instructions the better (not always true actually but well), Pointer to Memory.

Pointer to Memory

```

1 add_positive:
2     li t0, 0
3     mv t1, a1
4 next_short:
5     beq t1, zero, end
6     lh t2, 0(a0)
7     bltz t2, negative
8     add t0, t0, t2
9 negative:
10    addi a0, a0, 2
11    addi t1, t1, -1
12    j next_short
13 end:
14     mv a0, t0
15 ret

```

Index in array

```

1 add_positive:
2     li t0, 0
3     li t1, 0
4 next_index:
5     bge t1, a1, end
6     slli t2, t1, 1
7     add t2, a0, t2
8     lh t3, 0(t2)
9     bltz t3, negative
10    add t0, t0, t3
11 negative:
12    addi t1, t1, 1
13    j next_index
14 end:
15     mv a0, t0
16 ret

```

In C

Writing this in C for better understanding. Again, which one is better?

Obviously the less instructions the better (again not always true but ah), Index in array
Pointer to memory

```
short sum = 0;
short *ptr = myData;
short *end = myData + N;
while (ptr < end) {
    if (*ptr > 0) {
        sum += *ptr;
    }
    ptr++;
}
```

Index in array

```
1 short sum = 0;
2 int i;
3 for (i = 0; i < N; i++) {
4     if (myData[i] > 0) {
5         sum += myData[i];
6     }
7 }
```

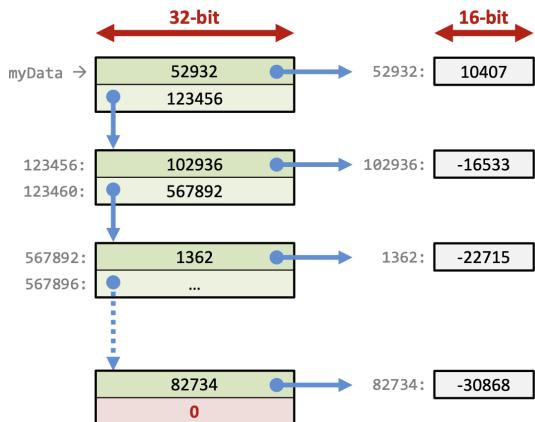
We need a good compiler

Seeing this, the idea would be to have a sufficiently good **compiler** (check I.4.3.2 if needed) such that we write our C code in *Index* in array, and we get Pointer to memory code in assembly. Thus writing better code but also getting better performance.

Another type of collection we could've used to store the data is a *Linked List*.

Linked lists are useful for efficiently inserting and deleting elements, especially in the middle of the list.

Each 32-bit element in a linked list contains 16 bits for the value and 16 bits for the address of the next element, enabling efficient insertions but slower sequential access compared to arrays.



Chapter 5

Part I(e) - ISA Arithmetic - W 3.1, 3.2

5.1 Notation

Before we start, let's define some notation:

- **Number representation (with a fixed number of digits/bits):**

$$A = A^{(n)} = A^{(m)}$$

- **Number in binary or decimal:**

$$A = A_{10} = A_2 = A_{2c}$$

With A_{2c} being the 2's complement representation.

And A_2 being the binary representation.

- **Individual digits or bits:**

$$a_{n-1}, a_{n-2}, \dots, a_2, a_1, a_0$$

- **Digit string representation:**

$$\langle a_{n-1} a_{n-2} \dots a_2 a_1 a_0 \rangle$$

5.2 Numbers

Numbers in computing can be represented in different forms, each with specific use cases.

Integers can be either signed or unsigned, representing positive and negative values, or only non-negative values. Examples include:

$$0, 1, 2, 3, 4294967295, -2147483648$$

Fixed-point numbers are essentially integers with an implicit scaling factor (e.g., 10^k or 2^k) to handle fractional values. Common in applications like signal processing. Examples include:

$$0.12, 3.14, 1073741823.75$$

Floating-point numbers represent a wide range of values using a base and exponent, providing flexibility in precision. Examples include:

$$3.14E3, -2.5E1, 1.0E0, 4.2E-2, -1.5E-3$$

5.2.1 Unsigned Integers

Unsigned integers are:

- *Weighted*: Each digit has a positional value.
- *Nonredundant*: Every number has a unique representation.
- *Based on a fixed-radix system*: Typically radix-10 (decimal) or radix-2 (binary).

- *Canonical*: Follows a standard form for representation.

Definition:

$$A = \langle a_{n-1}a_{n-2}\dots a_2a_1a_0 \rangle = \sum_{i=0}^{n-1} a_i R^i$$

where A is the unsigned integer, a_i are the digits, and R is the radix.

5.2.2 Signed Integers

We may distinguish between three methods for representing signed integers:

- **Sign-and-Magnitude (SM)**: Uses the most significant bit (MSB) to represent the sign (0 for positive, 1 for negative), with the remaining bits representing the magnitude. This method has the drawback of two zeros (+0 and -0) (Redundant).
- **Two's Complement**(Specific True-and-Complement): The most common way to represent signed integers. It avoids the two-zero problem and simplifies arithmetic operations. Negative numbers are represented by flipping the bits and adding 1.
- **Biased Representation**: Primarily used in floating-point numbers, especially for the exponent part. A fixed bias is added to the actual value to avoid negative exponents. It's rarely used for integers but is another method for handling signed numbers.

Sign and Magnitude

In the sign-and-magnitude representation, the most significant bit (MSB) is used to represent the sign of the number. The remaining bits represent the magnitude.

Definition

$$A = \langle sa_{n-2}a_{n-3}\dots a_2a_1a_0 \rangle = (-1)^s \cdot \sum_{i=0}^{n-1} a_i R^i$$

where A is the signed integer, s the most significant bit of A representing the sign of the number, a_i the digits, and R the radix.

Example (Signed 4-bit integer):

Consider the 4-bit signed binary number 1011_2 . In this case:

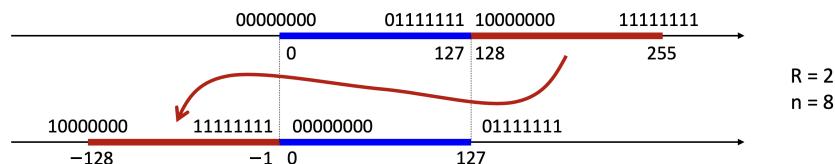
1. The MSB $s = 1$, indicating the number is negative.
2. The magnitude bits are $011_2 = 3_{10}$.
3. Therefore, the value of the number is -3 .

Thus, 1011_2 represents -3_{10} in sign-and-magnitude representation.

5.2.3 Radix's Complement

Radix's complement is a method used to represent signed numbers in different number systems.

It is a special form of *true-and-complement* where the complement $C = R^n$, with R being the radix (base) and n the number of digits.



Definition

A number A in radix's complement is represented as:

$$A = \langle a_{n-1}a_{n-2}\dots a_1a_0 \rangle = -a_{n-1}R^{n-1} + \sum_{i=0}^{n-2} a_i R^i$$

where a_{n-1} is the most significant bit, which also indicates the sign (negative for $a_{n-1} = 1$).

For binary numbers, radix's complement is known as **two's complement**, which is the most commonly used method for representing signed numbers in digital systems.

Binary (2's Complement) Representation

Two's complement uses base $R = 2$ and has a fixed word length n .

Here is an example for an 8-bit number system:

Binary	Decimal	Range
00000000	0	Positive range
01111111	127	
10000000	-128	Negative range
11111111	-1	

The two's complement system enables representation of both positive and negative numbers within a fixed bit length.

Decimal (10's Complement) Representation

In a decimal system with radix $R = 10$,

We use 10's complement to represent signed numbers. For instance:

$$5,678_{(5)}^{10c} = 05,678_{10c} = +5,678_{10}$$

This is a positive number representation in 10's complement. For a negative number:

$$9,999,999_{(7)}^{10c} = -1_{10}$$

Here, 9,999,999 in 7 digits represents -1 in decimal form.

Examples of Binary (2's Complement)

Below are several examples of numbers in binary (2's complement) and their corresponding decimal values:
This is a positive binary number.

$$0100,1101,0010_{(12)}^{2c} = 100,1101,0010_2 = +1,234_{10}$$

This is a negative binary number in 8-bit representation.

$$1111,1111_{(8)}^{2c} = -1_{10}$$

This is a negative binary number in 12-bit representation.

$$1011,0000,1110_{(12)}^{2c} = -1,234_{10}$$

5.2.4 Two's Complement Subtraction

Consider the binary subtraction using the standard paper-and-pencil method:

$$\begin{array}{r} \text{Borrow: } & -1 & -1 & -1 & & -1 \\ & 0 & 0 & 0 & 0 & 1 & 0 & 1 & 0 & (10_{10}) \\ - & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & (17_{10}) \\ \hline & 1 & 1 & 1 & 1 & 1 & 0 & 0 & 1 \end{array}$$

Since we had to borrow beyond the most significant bit, the result is negative. The binary result is:

$$-1\ 1\ 1\ 1\ 1\ 0\ 0\ 1_2$$

To find its decimal value:

$$-2^7 + 2^6 + 2^5 + 2^4 + 2^3 + 2^0 = -128 + 64 + 32 + 16 + 8 + 1 = -7$$

and

$$10_{10} - 17_{10} = -7_{10}$$

5.2.5 Addition Is Unchanged from Unsigned

In arithmetic operations, addition remains consistent whether using signed or unsigned numbers. The following instructions are available for basic arithmetic operations:

Arithmetic

add	rd, rs1, rs2	$rd \leftarrow rs1 + rs2$	R	0x00	0x0	0x33
addi	rd, rs1, imm	$rd \leftarrow rs1 + \text{sext}(imm)$	I		0x0	0x13
sub	rd, rs1, rs2	$rd \leftarrow rs1 - rs2$	R	0x20	0x0	0x33

- **add rd, rs1, rs2:** Adds the values in **rs1** and **rs2**, and stores the result in **rd**.
- **addi rd, rs1, imm:** Adds the value in **rs1** with the sign-extended immediate value **imm**, and stores the result in **rd**.
- **sub rd, rs1, rs2:** Subtracts the value in **rs2** from **rs1**, and stores the result in **rd**.

Note that older architectures (e.g., MIPS) had distinct instructions for signed (**add**) and unsigned (**addu**) addition. However, this distinction is unnecessary as the hardware handles both identically.

Sign-and-magnitude addition presents unique challenges, making **two's complement** the standard for signed integers in modern architectures.

5.2.6 Sign Extension

In digital systems, sign extension is a technique used to increase the bit width of a binary number while preserving its value and sign. It is commonly used when converting a number from a smaller to a larger bit width in a way that maintains its original meaning, whether it's unsigned or in two's complement format.

Example: 4-bit to 8-bit Conversion

Consider the 4-bit two's complement number 1110_2 , which represents -2_{10} .

When extending this number to 8 bits, we replicate the MSB (which is 1 in this case) to fill the additional bits, as shown below:

$$5_{10} = 0101_2 \quad (4 \text{ bits}) \rightarrow \quad 00000101_2 \quad (8 \text{ bits}).$$

while

$$-2_{10} = 1110_2 \quad (4 \text{ bits}) \rightarrow \quad 11111110_2 \quad (8 \text{ bits}).$$

This ensures that the number remains -2_{10} even after increasing the bit width.

Truncation is allowed when reducing bit width, but only if the truncated bits are redundant (i.e., copies of the sign bit). For example, going from 8 bits back to 4 bits would result in 1110_2 , preserving the value -2_{10} .

5.2.7 Signed and Unsigned Instructions

In RISC-V, instructions differentiate between signed (s) and unsigned (u) operations:

Shift						
srl	rd, rs1, rs2	$rd \leftarrow rs1 \gg_u rs2$	R	0x00	0x5	0x33
srli	rd, rs1, imm	$rd \leftarrow rs1 \gg_u imm$	I	0x00	0x5	0x13
sra	rd, rs1, rs2	$rd \leftarrow rs1 \gg_s rs2$	R	0x20	0x5	0x33
srai	rd, rs1, imm	$rd \leftarrow rs1 \gg_s imm$	I	0x20	0x5	0x13
Compare						
slt	rd, rs1, rs2	$rd \leftarrow rs1 <_s rs2$	R	0x00	0x2	0x33
slti	rd, rs1, imm	$rd \leftarrow rs1 <_s \text{sext}(imm)$	I		0x2	0x13
sltu	rd, rs1, rs2	$rd \leftarrow rs1 <_u rs2$	R	0x00	0x3	0x33
sltiu	rd, rs1, imm	$rd \leftarrow rs1 <_u \text{sext}(imm)$	I		0x3	0x13
Branch						
blt	rs1, rs2, imm	$pc \leftarrow pc + \text{sext}(imm \ll 1)$, if $rs1 <_s rs2$	B		0x4	0x63
bge	rs1, rs2, imm	$pc \leftarrow pc + \text{sext}(imm \ll 1)$, if $rs1 \geq_s rs2$	B		0x5	0x63
bltu	rs1, rs2, imm	$pc \leftarrow pc + \text{sext}(imm \ll 1)$, if $rs1 <_u rs2$	B		0x6	0x63
bgeu	rs1, rs2, imm	$pc \leftarrow pc + \text{sext}(imm \ll 1)$, if $rs1 \geq_u rs2$	B		0x7	0x63
Load						
lb	rd, imm(rs1)	$rd \leftarrow \text{sext}(\text{mem}[rs1 + \text{sext}(imm)][7 : 0])$	I		0x0	0x03
lbu	rd, imm(rs1)	$rd \leftarrow \text{zext}(\text{mem}[rs1 + \text{sext}(imm)][7 : 0])$	I		0x4	0x03
lh	rd, imm(rs1)	$rd \leftarrow \text{sext}(\text{mem}[rs1 + \text{sext}(imm)][15 : 0])$	I		0x1	0x03
lhu	rd, imm(rs1)	$rd \leftarrow \text{zext}(\text{mem}[rs1 + \text{sext}(imm)][15 : 0])$	I		0x5	0x03

- **Shift:** `sra, srai` (s) vs. `srl, srli` (u).
 - Signed shifts preserve the sign bit, while unsigned shifts insert zeroes.
- **Compare:** `slt, slti` (s) vs. `sltu, sltiu` (u).
 - Signed comparisons use two's complement, unsigned comparisons ignore sign.
- **Branch:** `blt, bge` (s) vs. `bltu, bgeu` (u).
 - Signed branches use two's complement; unsigned branches do not consider sign.
- **Load:** `lb, lh` (s) vs. `lbu, lhu` (u).
 - Signed loads extend the sign bit, while unsigned loads extend with zeroes.

5.3 Overflow

Overflow occurs when the result of an arithmetic operation exceeds the range of values that can be represented with a fixed number of bits. This can happen in both unsigned and signed arithmetic, though the detection method differs. In general, overflow results in an incorrect outcome that needs to be detected and handled.

5.3.1 Overflow in 2's Complement

In 2's complement arithmetic, overflow occurs when the result of an addition or subtraction operation falls outside the representable range for the number of bits. For an n -bit 2's complement system, the representable range is -2^{n-1} to $2^{n-1} - 1$.

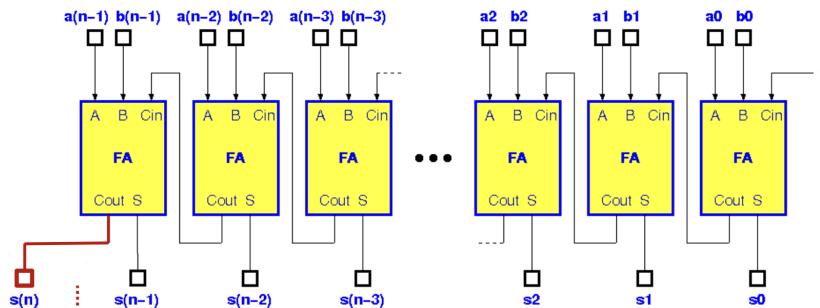
Overflow is detected by examining the carry into and out of the most significant bit (MSB). Specifically, overflow occurs if:

$$\text{Overflow} = \text{Cout}_{n-1} \oplus \text{Cout}_n$$

Where:

- Cout_{n-1} is the carry into the MSB.
- Cout_n is the carry out of the MSB.

An overflow occurs when these two carry bits differ. This is because the sign of the result is incorrect if there is a mismatch, leading to an incorrect outcome.

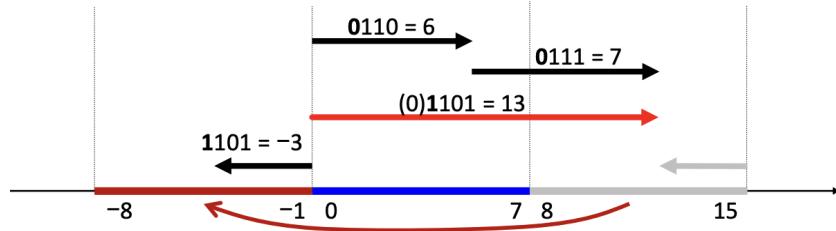


For example, if two large positive numbers are added and result in a negative value (or two negative numbers added result in a positive value), this indicates an overflow in 2's complement addition.

5.3.2 Overflow in Software

In many architectures, detecting overflow during arithmetic operations is a critical aspect of software implementation. Overflow occurs when the result of an addition or subtraction exceeds the capacity of the register used to store it. Detection methods vary depending on the type of architecture:

- **Traditional architectures (e.g., x86):** These systems provide a *carry bit* in a special register, known as a flag, that is set when an overflow occurs. Thus, overflow detection operates similarly to hardware-based overflow detection.
- **Modern architectures (e.g., RISC-V):** These architectures typically provide only the result of the addition or subtraction without a carry bit. Overflow detection must be handled in software, based on analyzing the sign and magnitude of the result.



Overflow detection can be based on the following observations:

- **Addition of opposite sign numbers:** The magnitude of the result decreases, making overflow impossible.
- **Addition of same sign numbers:** Overflow is possible if the result exceeds the range representable by the register, leading to an incorrect sign in the result.

5.3.3 Detect Addition Overflow in Software

- Add two 32-bit signed integers and detect overflow
 - At call time, `a0` and `a1` contain the two integers.
 - On return, `a0` contains the result and `a1` must be nonzero in case of overflow.

```

1 srai a2, a0, 31      # a2 = sign of a0 (0 or -1)
2 srai a3, a1, 31      # a3 = sign of a1 (0 or -1)
3 xor a4, a2, a3       # a4 = 0 if signs are same, -1 if different
4 add a0, a0, a1        # compute sum in a0
5 srai a5, a0, 31       # a5 = sign of sum (0 or -1)
6 xor a6, a2, a5       # a6 = 0 if sign of sum same as a0, -1 if different
7 and a1, a4, a6        # a1 = -1 if overflow occurred, else 0
8 srli a1, a1, 31       # a1 = 1 if overflow occurred, else 0

```

5.4 A Strange but Useful Property

Personal Remark: don't mistake A and \bar{A} as sets of elements which might confuse you. They are binary numbers.

In binary arithmetic, there is a particularly useful property that can be expressed as follows:

$$A + \bar{A} = -1$$

or equivalently,

$$-A = \bar{A} + 1$$

Proof: Consider a binary number $A = a_{n-1}2^{n-1} + \sum_{i=0}^{n-2} a_i 2^i$, where $a_i \in \{0, 1\}$ represents the binary digits of A . The complement of A , denoted \bar{A} , is given by replacing each a_i with its complement \bar{a}_i .

$$\begin{aligned}
 A + \overline{A} &= \left(-a_{n-1}2^{n-1} + \sum_{i=0}^{n-2} a_i 2^i \right) + \left(-\overline{a}_{n-1}2^{n-1} + \sum_{i=0}^{n-2} \overline{a}_i 2^i \right) \\
 &= -(a_{n-1} + \overline{a}_{n-1}) \cdot 2^{n-1} + \sum_{i=0}^{n-2} (a_i + \overline{a}_i) \cdot 2^i \\
 &= -2^{n-1} + \sum_{i=0}^{n-2} 2^i = -1
 \end{aligned}$$

Where \overline{A} is the two's complement of A .

Intuition: For each binary digit, adding a_i and its complement \overline{a}_i results in 1. Therefore, $A + \overline{A}$ consists entirely of 1s, representing -1 in two's complement.

5.4.1 Two's Complement Subtractor

Using the property of two's complement, we can create a subtractor circuit. The subtractor is implemented using an adder, where the number to be subtracted is inverted and incremented by 1.

- **Step 1: Inversion of Subtrahend (B)**

The subtrahend B is inverted using NOT gates, as shown in the diagram. This converts B into its one's complement.

- **Step 2: Addition of A and Inverted B**

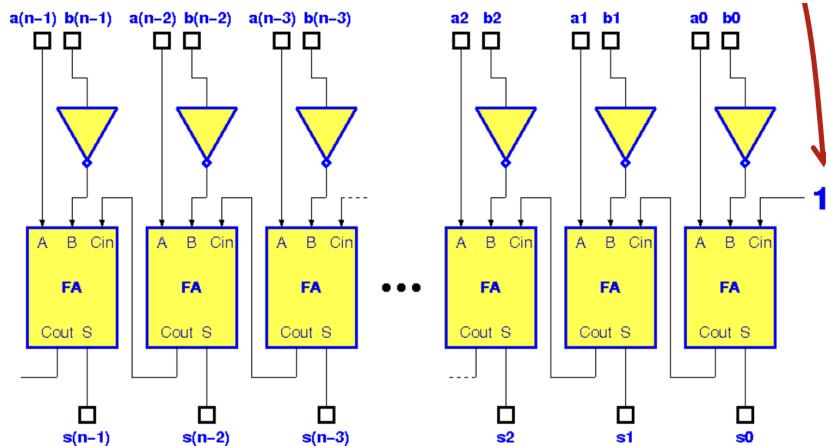
The full adders (FA) add each bit of the minuend A to the inverted bits of B . The full adders also handle any carry-over from the previous addition.

- **Step 3: Add 1 (Two's Complement)**

To complete the two's complement operation, a carry-in of 1 is added to the least significant bit (LSB), which effectively adds 1 to the inverted B .

- **Output:**

The sum outputs S (s_0, s_1, s_2, \dots) represent the result of the subtraction $A - B$, while the final carry-out can be used to detect overflow.

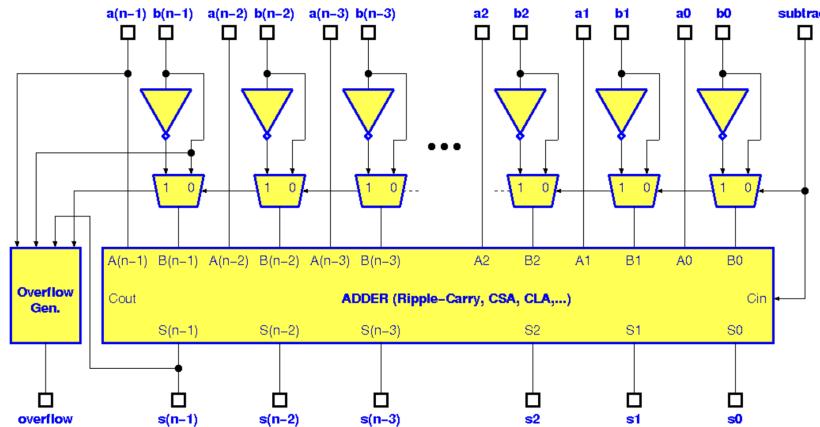


5.4.2 Two's Complement Add/Subtract Unit

This circuit performs both addition and subtraction using two's complement arithmetic. The operation is selected based on the control input signal for subtraction. The unit consists of several key components:

- **Input Inversion:** Each bit of the subtrahend B is passed through a XOR gate controlled by the 'subtract' signal. When the 'subtract' signal is high (logic 1), the bits of B are inverted to form the two's complement of B , effectively switching the operation to subtraction.

- **Addition:** The ripple-carry adder, represented by the ADDER block, performs binary addition of the bits from A and B . The carry-in (Cin) to the least significant bit is used to add 1 when performing subtraction, completing the two's complement process.
- **Overflow Detection:** The overflow generator block detects if the result of the addition/subtraction operation has exceeded the range representable by the fixed number of bits. The ‘overflow’ output is asserted in such cases.
- **Output:** The result of the operation is provided as the sum output (S), representing either the sum $A + B$ or the result of $A - B$, depending on the control signal.



5.5 Bounds Check Optimization

Very, very, very useful. When working with signed integers (e.g., array indices), a common task is to ensure that the index remains within a valid range, typically $0 \leq t0 < N$, where N is some predefined boundary. This can be achieved efficiently using a single branch check that combines both lower and upper bound constraints.

Single Branch Bound Check

The instruction `bgeu` (branch if greater than or equal, unsigned) can perform two checks at once:

```
bgeu t0, t1, out_of_bound
```

Here, $t0$ is the signed number to be checked, and $t1 = N$ is the boundary.

Explanation

- If $t0 \geq 0$, the behavior of `bgeu` mimics that of `bge` (branch if greater than or equal) for signed integers, thus effectively performing an upper bound check.
- If $t0 < 0$, since the comparison is unsigned, $t0$ will appear as a very large positive value, hence automatically triggering the out-of-bound case.

This approach efficiently checks both the lower and upper bounds in one instruction, streamlining the bounds checking process.

5.6 Floating Point Representation

Floating point numbers are widely used in computing to represent real numbers in a way that supports a wide dynamic range.

They are composed of a *significand* (or *mantissa*) and an *exponent* of the base. This representation allows for the approximation of very large and very small values, similar to the way scientific notation is used in everyday practices.

Such as

$$\begin{aligned} 0.18 \mu\text{m} &\rightarrow 0.18 \cdot 10^{-6} \text{ m} \rightarrow 1.8 \cdot 10^{-7} \text{ m} \\ 75 \text{ km} &\rightarrow 75 \cdot 10^3 \text{ m} \rightarrow 7.5 \cdot 10^4 \text{ m} \end{aligned}$$

In floating point representation, a number X is expressed as:

$$X = (-1)^s \cdot \left(\sum_{i=0}^{n-1} a_i \cdot 2^i \right) \cdot 2^{\left(-e_{m-1} 2^{m-1} + \sum_{j=0}^{m-2} e_j 2^j \right)}$$

where:

- s is the sign bit,
- a_i represents the bits of the significand (in sign-and-magnitude form),
- e_j represents the bits of the exponent (in 2's complement form).

Properties of Floating Point Numbers

- Large dynamic range, but *variable accuracy*.
- Numbers are **redundant** unless *normalized*.
- Floating point operations are **not associative**, unlike real numbers.
- Exponents are typically stored in a **biased signed representation**, making zero easier to represent and simplifying comparisons in hardware.
- The **mantissa** (significand) is usually normalized such that $1 \leq m < 2$, with a *hidden bit* to store the leading 1.

Standardization and Hardware Support

Floating point representation is standardized by the IEEE 754 standard, which is widely adopted in modern computing systems:

- **x86/x64** architectures have supported floating point operations through SSE/AVX extensions since 1999.
- **RISC-V** also includes support for floating point through ISA extensions.

Example: Decimal to IEEE 754 Simple Precision (32 Bits) Conversion

Convert -7.75 to IEEE 754 single-precision:

Step 1: Sign Bit (1 Bit)

$s = 1$ (negative number).

Step 2: Binary Conversion

$7_{10} = 111_2$, $0.75_{10} = 0.11_2$, so $7.75_{10} = 111.11_2$.

Step 3: Normalize

$111.11_2 = 1.1111_2 \times 2^2$.

Step 4: Exponent (8 Bits)

$E = 2 + 127 = 129$, $129_{10} = 10000001_2$.

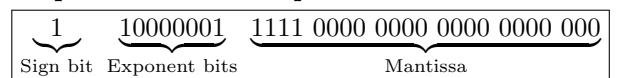
Step 5: Mantissa (23 Bits)

Take the fractional part after the leading 1 and pad with zeros to make 23 bits:

1111 0000 0000 0000 0000 000

(fractional part after the leading 1).

Step 6: IEEE 754 Representation



5.6.1 Sign-and-Magnitude Addition

(Assembly)

In this exercise, we aim to write a function in RISC-V assembler to sum two 32-bit signed numbers represented in sign-and-magnitude (S&M) format. The result should also be produced in the sign-and-magnitude format.

- The two operands are stored in registers `a0` and `a1` on entry.
- The result should be placed in register `a0`.
- Overflow cases should be ignored.

Solution 1

Basic Algorithm:

- If the operands have the same sign:
 - Add the absolute values
 - Attach to the result the same sign as the operands
- If the operands have different signs:
 - Subtract the smallest absolute value from the largest one
 - Attach to the result the sign of the largest value

```

1 add_sandm:
2   lui      t1, 0x80000      # mask for sign bit
3   and     t0, a0, t1        # check a0 sign
4   beqz    t0, a0_positive  # if positive, skip
5   xor     a0, a0, t1        # flip sign bit
6   neg     a0, a0          # negate a0
7
8 a0_positive:
9   and     t0, a1, t1        # check a1 sign
10  beqz   t0, a1_positive  # if positive, skip
11  xor     a1, a1, t1        # flip sign bit
12  neg     a1, a1          # negate a1
13
14 a1_positive:
15  add     a0, a0, a1        # add values
16  and     t0, a0, t1        # check result sign
17  beqz   t0, sum_positive # if positive, skip
18  neg     a1, a1          # negate a1
19  xor     a0, a0, t1        # flip sign bit
20
21 sum_positive:
22   ret                  # return result

```

Solution 2**Basic Algorithm:**

- Convert the two operands from sign-and-magnitude to 2's complement
- Add the two operands
- Convert the result from 2's complement back to sign-and-magnitude

```

1 convert_to_twos_comp:
2   lui      t1, 0x80000      # mask for sign bit
3   and     t0, a0, t1        # check a0 sign
4   beqz    t0, a0_twos_comp # if positive, skip
5   xor     a0, a0, t1        # flip sign bit
6   neg     a0, a0            # negate a0
7
8 a0_twos_comp:
9   and     t0, a1, t1        # check a1 sign
10  beqz   t0, a1_twos_comp # if positive, skip
11  xor     a1, a1, t1        # flip sign bit
12  neg     a1, a1            # negate a1
13
14 add_twos_comp:
15   add    a0, a0, a1        # add values
16
17 convert_to_sign_mag:
18   lui      t1, 0x80000      # mask for sign bit
19   and     t0, a0, t1        # check result sign
20   beqz    t0, result_positive # if positive, skip
21   xor     a0, a0, t1        # flip sign bit
22   neg     a0, a0            # negate result
23
24 result_positive:
25   ret                  # return result

```

Chapter 6

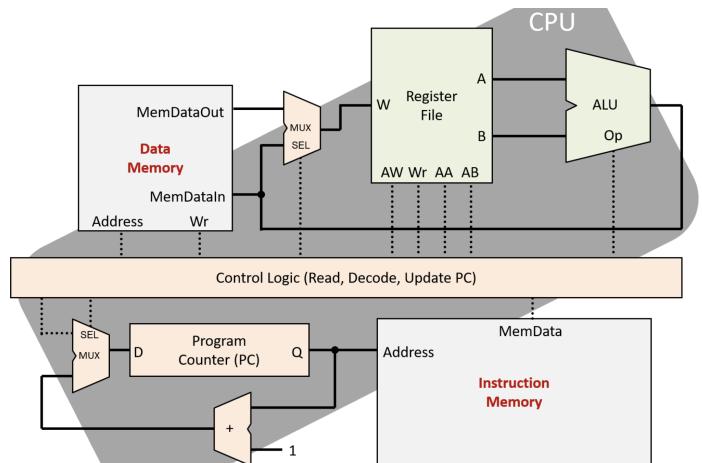
Part II(a) - I/O - Exceptions Multicycle Processor W - 3.2, 4.1

In this chapter we will be discussing how we can actually design a processor (subject of our LAB B).

6.1 Processor

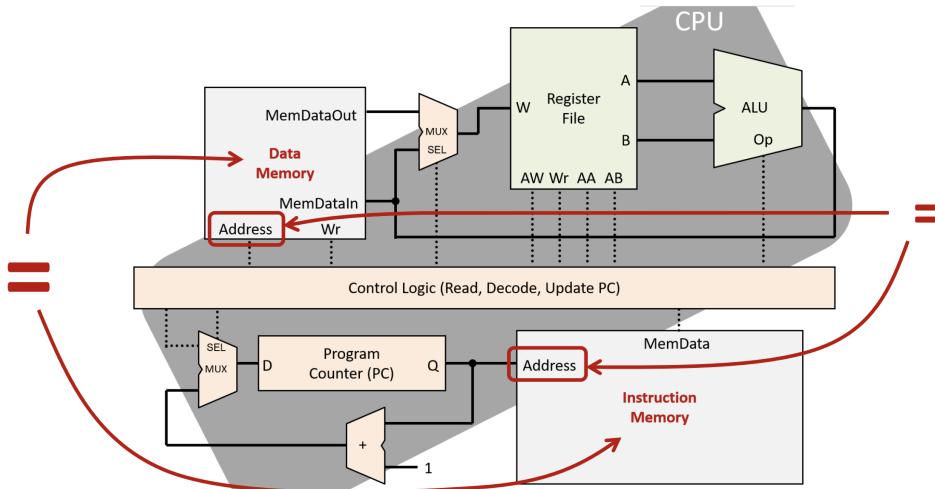
(yes, one more time) A processor is composed of several fundamental components that work together to perform computations.

- **Program Counter (PC):** Holds the address of the next instruction to be executed from the instruction memory. It increments after each instruction fetch or is updated based on control logic.
- **Instruction Memory:** Stores the instructions that the processor fetches and executes. Instructions are read sequentially unless altered by control logic.
- **Control Logic:** Manages the flow of data and the sequence of operations, including reading instructions, decoding them, and updating the program counter.
- **Register File:** A set of registers where data is temporarily stored. It allows the processor to access and manipulate values quickly. Each register has read/write capabilities.
- **Arithmetic Logic Unit (ALU):** Performs arithmetic and logical operations. The inputs are provided by the register file, and the result is stored back into the registers or data memory.
- **Data Memory:** Stores data that can be written to or read from during program execution. It interacts with both the register file and the ALU for storing operands and results.



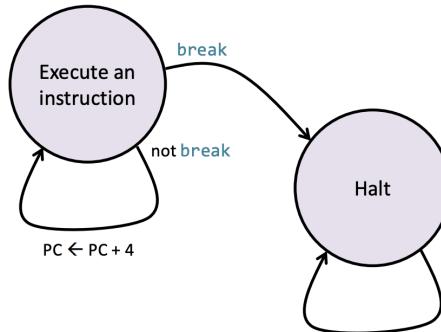
6.1.1 Unified Memory

In the image above, we see that the data memory and instruction memory are separate. However, a choice that is often made is to have a unified memory.



6.1.2 Single-Cycle Processor

At the end, like most circuits, a processor is just another Finite State Machine. The simplified state diagram of a single-cycle processor would like this:



Execute an instruction, move to the next, repeat.

This simplified view doesn't reflect actual CPU design. In reality, instructions take different amounts of time due to complexity and **Propagation Time**—the delay in signal travel through the processor.

6.2 Propagation Time

Remember the difference (**this is absolutely critical to understand the rest of the course**) between combinational circuits and sequential circuits.

As the name suggests, sequential circuits are built like a *sequence*(mnemotechnic), meaning the current output depends on both the current input and the previous state.

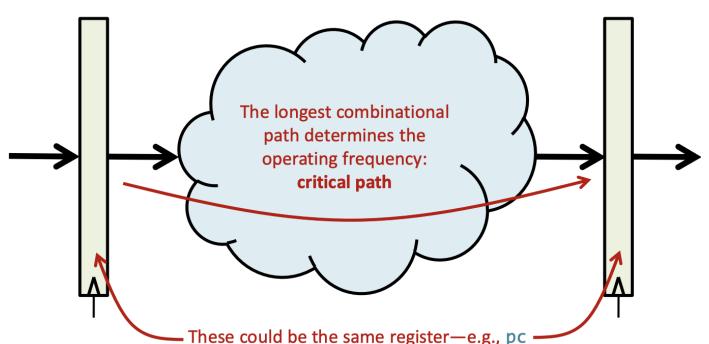
While combinational circuits, don't have a memory, they just take an input and give out an output.

The main thing to understand here is that, for our circuits to function as intended, the **propagation time** must allow the combinational circuits to complete before the next clock cycle (otherwise, it would lead to *obvious bugs*).

This implies that we need to observe the **longest combinational path** and account for it when designing our circuits.

While this is the *efficient approach*, one could, in theory, design a propagation time that is longer than the longest path. However, this would result in a *waste of both time and resources*.

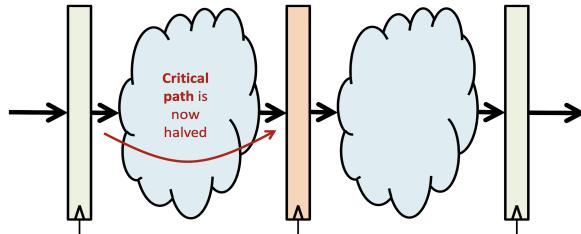
Remember, lower propagation time means higher clock frequency, which means faster processing.



6.2.1 Increasing the Frequency

To increase the frequency, we need to decrease the propagation time. This can be achieved by breaking down the combinational path into smaller parts.

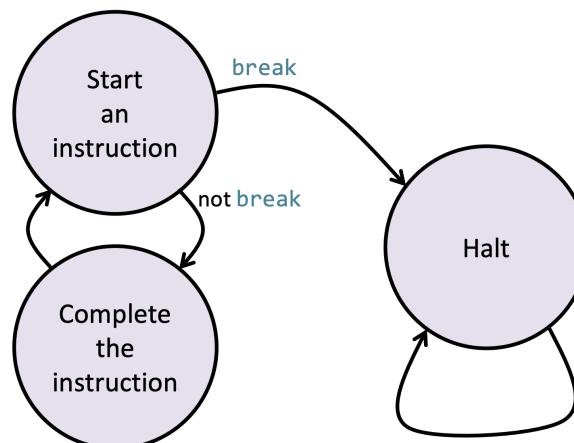
For example, consider the ‘lw’ instruction. This requires adding the offset to the base address (which involves addition, not completely trivial), and then reading the data from memory. This process can be broken down into two stages: first, the addition, and then the memory read.



By doing this, we can operate at twice the “”speed””(we’ll see why this is wrong in a moment).

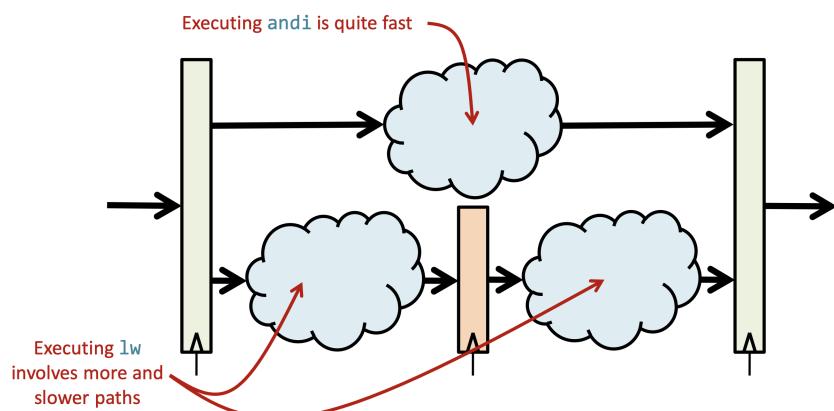
6.2.2 Two-Cycle Processor

However, what we quickly realize is that this approach doesn’t result in a real performance gain. While the processor runs at twice the frequency, it also takes twice as long to complete the instruction, leading to no overall improvement. *Historically, Intel often used this strategy to persuade uninformed consumers that their processors were getting faster.*



6.2.3 Not All Paths Are Born Equal

The reason we’re discussing this is that not all paths are equal. Some instructions are faster to compute than others. For example, the `andi` instruction is much faster than the `lw` instruction.

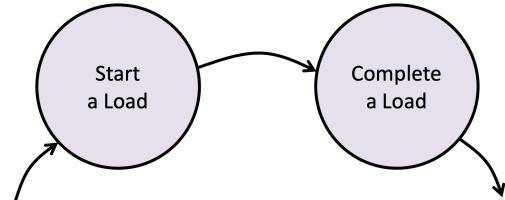
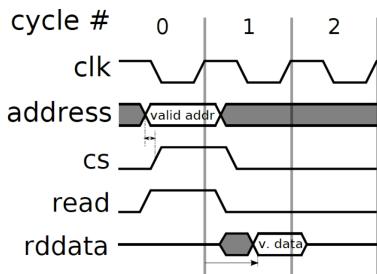


6.2.4 Asynchronous/Synchronous Memories

Another reason why breaking down the combinational path could be beneficial is that certain memories are **Synchronous**, meaning they only read data from a valid memory address on the rising edge of the clock cycle.

On the other hand, **Asynchronous** memories read data as soon as a valid memory address is available, without waiting for the clock cycle.

So, for **Synchronous** memories, breaking down combinational paths into smaller segments allows us to increase the clock frequency, making memory updates faster.



6.3 Multicycle Processor

Now let's try to construct a more convincing representation for our processor.

The processor operates in two cycles: a faster path for simple instructions and a slower path for more complex ones.

- Fetch1/Fetch2:

Simple: Uses only Fetch1 for single-word instructions.

Complex: Uses Fetch2 to fetch additional data when needed (e.g., multi-word instructions).

- Decode:

Simple: Quick decoding with fewer control signals.

Complex: More control signals and operands, requiring extra decoding time (and extra Optimization could be to introduce two Decoding stages for simple/complex instructions).

- Execute:

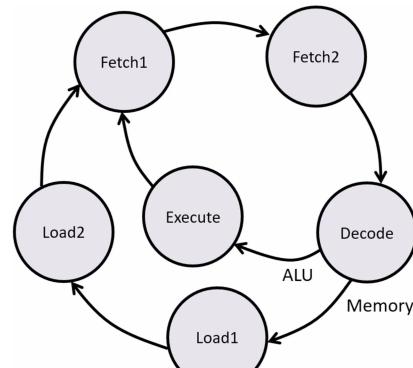
Simple: Fast ALU operations like additions.

Complex: Involves branches or complex ALU operations.

- Load1/Load2:

Simple: Skips Load stages if no memory access.

Complex: Memory operations use Load1 and Load2 to fetch and process data.



While this is an efficient design, it is not unique. The two things to keep in mind when designing a processor are:

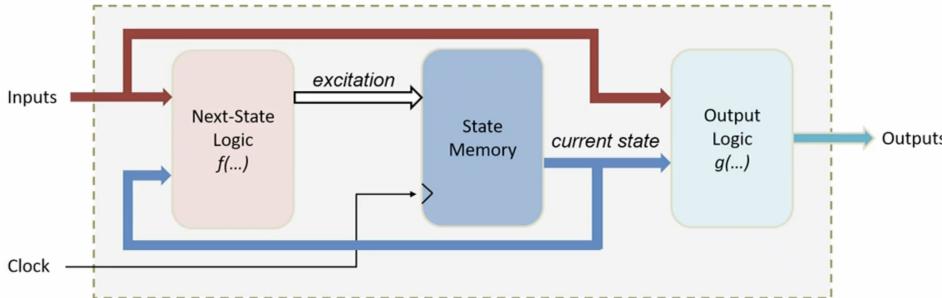
not to have too many stages, *meaning that having an excessive number of stages could increase the complexity and latency of the processor (this we will see later in the course)*.

to have paths as balanced as possible, *meaning that the duration of each stage should be similar to avoid bottlenecks that would slow down the overall process. The more balance we have the more we can profit from fast cases*.

6.4 Mealy or Moore?

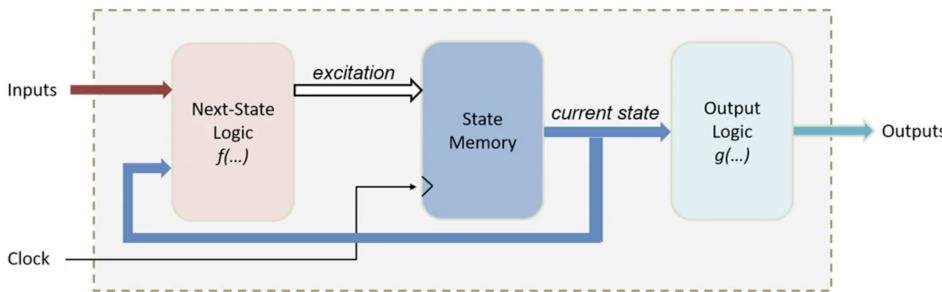
Personal Remark (mnemotechnic)

Moore - Output Only depends on state (double O like in Moore),
Mealy - Output depends on state and input



Mealy FSM

Output depends on
input and state



Moore FSM

Output depends on
state only

It is generally preferable to use Moore state machines because their outputs depend only on the current state, making them simpler to design, debug, and predict, whereas Mealy machines depend on both state and input, introducing complexity and potential glitches. So unless the specifications requires us to do otherwise, we'll generally tend to represent our state machines as Moore machines.

6.5 Processor - Building the Circuit

In this part, we will be incrementally adding the components needed to build our processor circuit.

For now, we've added two components to our CPU:

Controller: This component, although empty for now, will eventually manage the flow of data and sequence of operations within the CPU. It will control how data moves and instructions are executed.

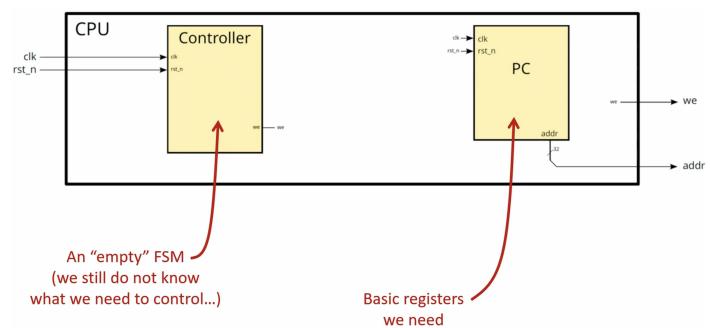
PC (Program Counter): The PC holds the address of the next instruction to be executed from the instruction memory. It increments after each instruction fetch or is updated based on control logic.

Inputs

- **clk:** The clock input ensures the program counter updates synchronously with the system clock.
- **rst_n:** An active-low reset signal that resets the program counter to a default value when low (0).
- **en:** The enable signal controls whether the PC updates its value (controlled by the Controller's FSM).

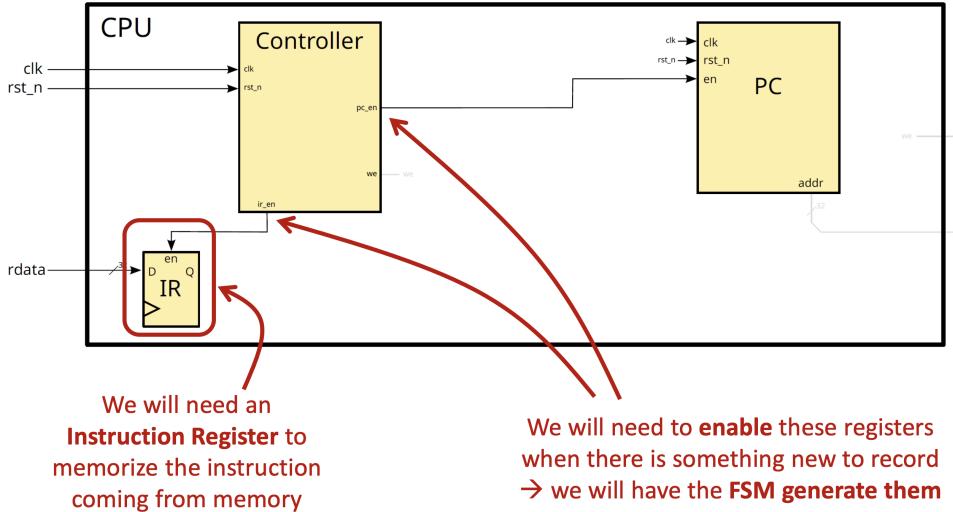
Outputs

- **addr:** The address output representing the next instruction to be fetched from memory.



6.5.1 Adding the Instruction Register

Now we're adding an Instruction Register.



In this step, we introduce the Instruction Register (IR) to our CPU:

Instruction Register (IR): The IR is responsible for storing the instruction fetched from memory. It captures the instruction ready to be decoded and executed. The Controller generates enable signals to control when the PC and IR should update their contents.

PC (Program Counter)

Inputs

- **clk:** The clock input ensures the program counter updates synchronously with the system clock.
- **rst_n:** An active-low reset signal resets the program counter to a default value when low (0).
- **en:** The enable signal controls whether the PC updates its value. It is driven by the FSM in the Controller.

Outputs

- **addr:** The address output representing the next instruction to be fetched from memory.

IR (Instruction Register)

Inputs

- **clk:** Ensures the instruction register captures the instruction at the correct clock edge.
- **rst_n:** Active-low reset to reset the IR to its default state.
- **en:** The enable signal controls whether the IR updates its contents. It is activated when a new instruction is fetched from memory.
- **D:** The data input, which represents the instruction fetched from memory (**rdata**).

Outputs

- **Q:** The output of the instruction register, representing the stored instruction that will be decoded and executed.

Controller

Inputs

- **clk:** The clock signal to ensure synchronization with other components.
- **rst_n:** The active-low reset signal to reset the controller to its initial state.

Outputs

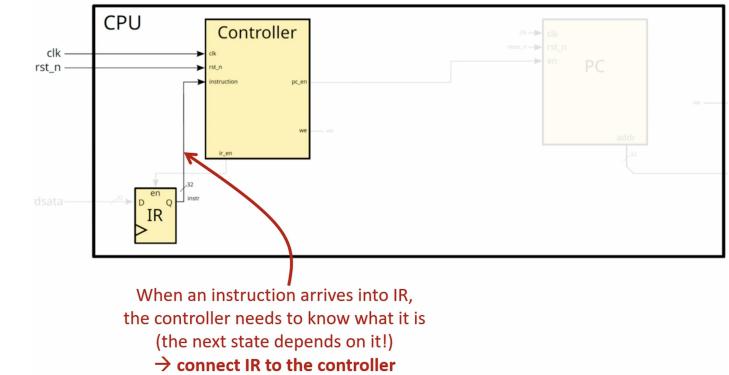
- **pc_en:** The enable signal sent to the Program Counter (PC) to control when it should update its value.
- **ir_en:** The enable signal sent to the Instruction Register (IR) to control when it should store a new instruction.

6.5.2 Adding functionality

Once an instruction is fetched from memory and stored in the Instruction Register (IR), it is crucial for the Controller to receive this instruction. The Controller needs the instruction to determine the next sequence of operations, as the next state of the system is dependent on the instruction being executed.

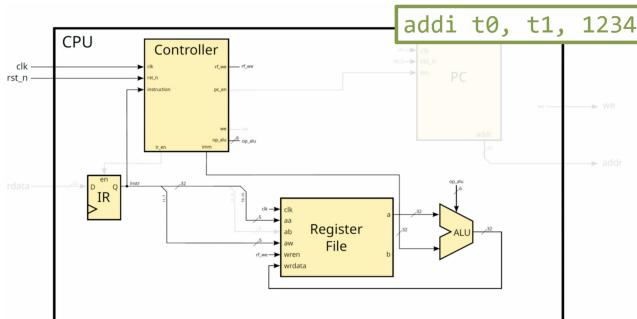
The Q output of the IR, which holds the stored instruction, is fed directly to the Controller. This connection allows the Controller to decode the instruction and control the subsequent operations of the CPU.

Specifically, the Controller will enable or disable other components, such as the Program Counter (PC), based on the instruction being processed.



6.5.3 I-Type Instructions Need RF and ALU

I-Type instructions such as `addi t0, t1, 1234` require both the register file (RF) and the Arithmetic Logic Unit (ALU) for execution. The operation consists of an addition between a register value and an immediate value, and the result is stored back into a register.



ALU (Arithmetic Logic Unit)

Inputs

- **a:** First operand input from the register file (e.g., `t1`).
- **b:** Second operand input, typically the immediate value for I-type instructions (e.g., `1234`).
- **op_alu:** Control signal from the controller specifying the operation to perform (e.g., addition for the `addi` instruction).

Outputs

- **alu_out:** The result of the operation performed by the ALU (e.g., the sum of `t1` and the immediate value).

Register File

Inputs

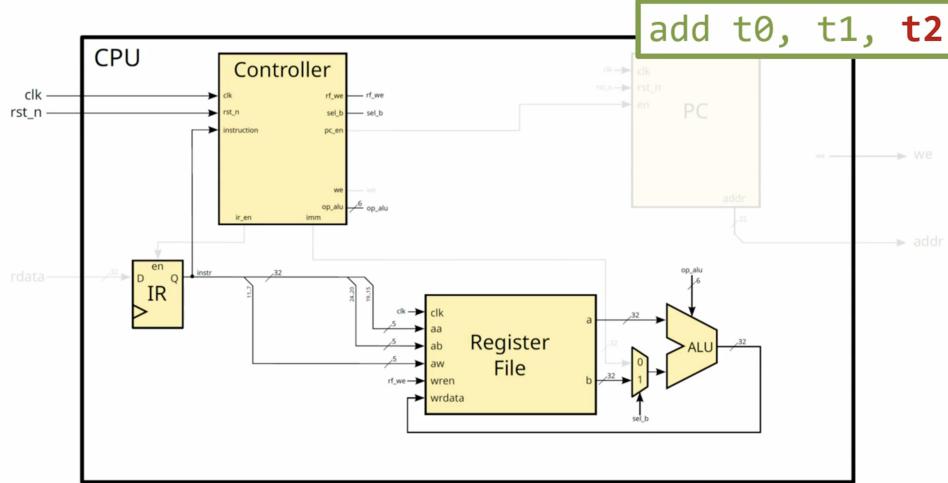
- **clk:** The clock input that ensures register updates are synchronous with the system clock.
- **aa:** The address of the first register (e.g., `t1`) from which data will be read.
- **ab:** The address of the second register (for other instruction types).
- **aw:** The address of the destination register (e.g., `t0`) to which the result will be written.
- **wren:** Write enable signal that allows data to be written into the destination register.
- **wrdata:** The data to be written into the destination register (e.g., the result from the ALU).

Outputs

- **a:** The data from the first register (e.g., the value stored in `t1`).
- **b:** The data from the second register (for other instruction types).

6.5.4 R-Type Instructions and Second Operand Selection

R-Type instructions, such as `add t0, t1, t2`, require two register operands and involve several components for execution. The instruction specifies two source registers (`t1` and `t2`) and a destination register (`t0`), with the second operand selected from the register file rather than using an immediate value. The multiplexer plays a key role in selecting the correct second operand based on the instruction type.



Register File

Inputs

- **clk:** The clock input that ensures register updates are synchronous with the system clock.
- **aa:** The address of the first register (e.g., `t1`) from which data will be read.
- **ab:** The address of the second register (e.g., `t2`) from which data will be read.
- **aw:** The address of the destination register (e.g., `t0`) where the result will be written.
- **wren:** Write enable signal that allows data to be written into the destination register.
- **wrdata:** Data to be written into the destination register (e.g., the result from the ALU).

Outputs

- **a:** The data from the first register (e.g., the value stored in `t1`).
- **b:** The data from the second register (e.g., the value stored in `t2`).

Multiplexer (sel_b)

Inputs

- **b:** The second operand, which can either be the register value (`t2`) or an immediate value, depending on the instruction type.
- **sel_b:** The select signal from the controller, determining whether the second operand is a register value (`t2`) or an immediate value.

Outputs

- **selected_b:** The selected operand output, which forwards either the register value (`t2`) or the immediate value to the ALU as the second operand.

ALU (Arithmetic Logic Unit)

Inputs

- **a:** First operand input from the register file (e.g., `t1`).
- **b:** Second operand input, selected by the multiplexer, from the register file (e.g., `t2`).
- **op_alu:** Control signal from the controller specifying the operation to perform (e.g., addition for the `add` instruction).

Outputs

- **alu_out:** The result of the operation performed by the ALU (e.g., the sum of `t1` and `t2`), which is written back into the destination register.

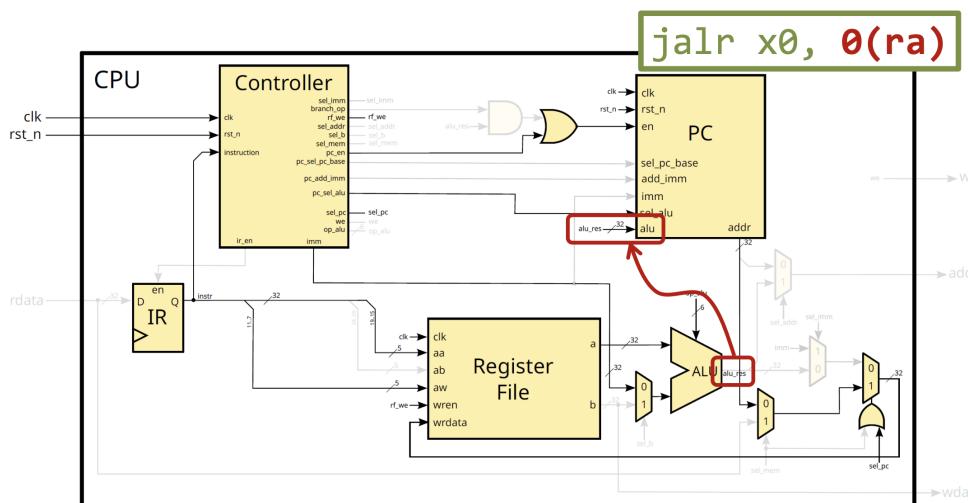
6.5.5 And More, and More...

After these few additions, you basically get the point, we keep adding, block by block, the components we need for the full use of our processor, professor also goes pretty quickly over this (you'll also see the full implementation of this in LAB B.)

The rest of the additions being :

- U-Type Instructions Write an Immediate
- Load and Stores Produce a Memory Address
- Loads Write the Read Data into the RF
- Stores Send an Operand to Memory
- Branches Need to Write an Offset to the PC
- jal Needs to Store PC + 4 in the RF
- Jumps Need to Write an Address to the PC

The processor after all of this looks like this:



6.5.6 Guidelines for Writing Verilog

Before beginning to write Verilog code, it is crucial to follow certain guidelines to ensure clarity and correctness in your hardware design. Verilog and VHDL are Hardware Description Languages (HDLs) that require a clear and structured approach.

Anything that's complicated is a Module, anything that is trivial, we need to know if it's sequential or combinational.

- **Clarity and Preparation:**

- Ensure that you have drawn a diagram, as demonstrated in previous examples.
- Clearly distinguish between **combinational** and **sequential** blocks in your design.

- **Decomposition of Complex Sequential Blocks:**

- Break down complex sequential blocks into simpler, well-defined elements. For instance, sequential blocks should primarily consist of simple registers (e.g., Instruction Register - IR).
- Continue refining your hierarchical diagrams until all sequential blocks become trivial to implement.

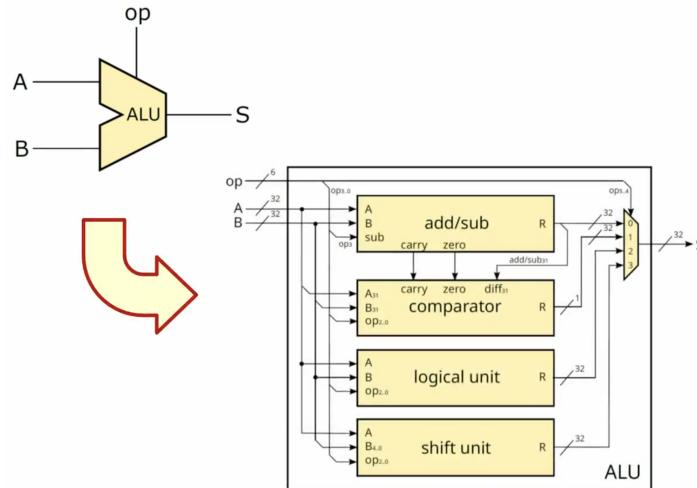
- **Adopt a Hierarchical Approach:**

- Use a hierarchical approach, similar to programming practices, and employ your diagrams to guide the creation of modules, such as the Program Counter (PC).

For example, for our processor, identifying that a register file is sequential while a PC is combinational is crucial before starting to write Verilog.

6.5.7 Detailing Complex Combinational Modules (ALU)

When designing complex combinational modules, it is essential to clearly define and break down each component to ensure accurate and efficient implementation. The following steps outline the process of detailing these modules:



- **ALU (Arithmetic Logic Unit) Overview:**

- The ALU receives inputs A , B , and an operation code (op), and produces an output S .
- It contains multiple submodules, such as add/subtract, comparator, logical unit, and shift unit.

- **Add/Subtract Unit:**

- The add/subtract unit performs addition and subtraction operations based on the control signal sub .
- It includes circuitry to handle carry and zero detection, essential for arithmetic operations.

- **Hierarchical Design:**

- The ALU is composed hierarchically, where each submodule (e.g., add/subtract, comparator) performs specific functions and connects to the overall ALU structure.
- Such a design allows for easier debugging, maintenance, and understanding of each module's role within the ALU.

6.5.8 Verilog - Sticking to Basic Patterns

When writing Verilog, it is essential to adhere to basic patterns for describing combinational and sequential logic. This section provides guidelines on structuring Verilog code efficiently.

Combinational Logic

Combinational logic blocks should be described using the `always @(*) begin` construct. This approach ensures that outputs are updated whenever the inputs change.

```

1  always @(*) begin
2    if (a)
3      y = \~b;
4    else
5      y = b;
6  end

```

Complex combinational blocks, such as the next state in a finite state machine (FSM), can also be described using this pattern.

For detailed guidelines, refer to the Verilog guidelines provided in Moodle.

Sequential Logic

Sequential logic blocks should be described using the `always @(posedge clk) begin` construct. This pattern is suitable for describing registers and counters.

```

1  always @(posedge clk) begin
2    if (reset == 1)
3      q <= 0;
4    else if ((enable1 == 1) && (
5      enable2 == 1))
6      q <= d;

```

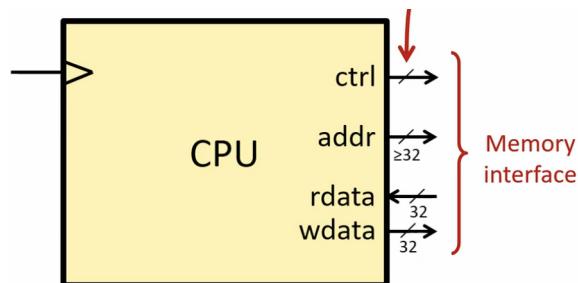
Use `posedge clk` to trigger updates on the rising clock edge.

Chapter 7

Part II(b) - Processor, I/Os, and Exceptions W - 4.1 - 4.2

7.1 The CPU

The CPU is a very sequential component responsible for executing instructions in a controlled manner. The CPU interacts with the memory through a defined memory interface, which includes various control signals and data pathways.



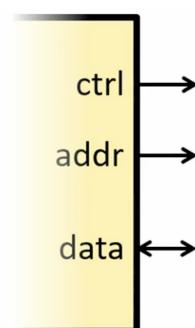
- **Control Signals (ctrl):** These signals manage the behavior of memory access, indicating whether to read or write data.
- **Address (addr):** Specifies the memory address where the CPU wants to read or write data. The width of the address bus is typically 32 bits or more.
- **Read Data (rdata):** A 32-bit pathway through which the CPU receives data from the memory.
- **Write Data (wdata):** A 32-bit pathway through which the CPU sends data to be stored in memory.

The memory interface is also controlled by two important signals:

- **Circuit Enable (CE):** Validates the address, indicating that the address provided is active and the operation should proceed.
- **Write Enable (WE):** Indicates that the current access is a store operation, allowing data to be written into memory.

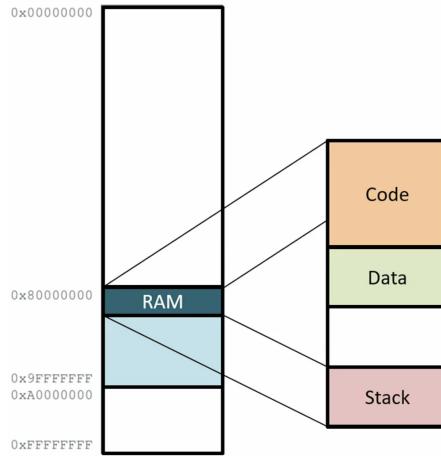
From now on, the clock signal, which drives the sequential behavior, may be omitted for simplicity.
This interface design allows for a clear and structured method of communication between the CPU and memory, ensuring reliable execution of instructions and data management.

Some processors, instead of having two separate buses, have a single data bus, that can be used for both reading and writing. This is known as a **bidirectional data bus**. This means, this kind of system uses a tri-state buffer to control the direction of the data flow.

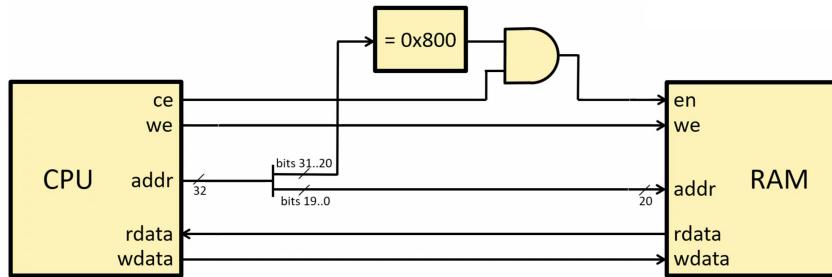


7.2 Physical Memory Map

To connect the CPU to memory, we need to define a *physical memory map*. As we had 32 bits of address, we can address 2^{32} bytes of memory (approx. 4GB of memory at least in a CPU).



7.2.1 Connecting CPU and Memory



Address Mapping: The upper segment of the address (bits 31..20) is compared with the hexadecimal value 0x800 to determine whether the target memory location resides within the RAM range. If the comparison (XNOR gate) yields a match, the **en** (enable) signal for the RAM is activated, allowing the subsequent read or write operation.

Control Signals: Two primary control signals are involved in coordinating memory operations:

- **ce** (chip enable): Indicates whether the CPU is ready to perform an operation (read or write) on a specific memory chip.
- **we** (write enable): Indicates whether the operation is a write operation, allowing data to be written into RAM.

The lower segment of the address (bits 19..0) is passed directly to the RAM as the 20-bit address input, specifying the exact memory location within the enabled RAM region.

This mapping allows the CPU to access specific regions of RAM by comparing higher address bits with predefined values and appropriately enabling or disabling memory segments.

7.3 Input/Output (I/O) Devices

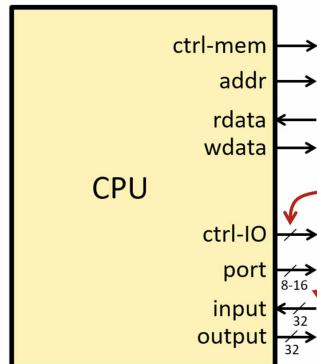
Input/Output (I/O) devices serve as crucial interfaces for various types of communication between a computer system and the external environment. I/O device data rates vary widely based on their type and purpose. Low-bandwidth devices like keyboards and mice handle simple inputs, while modern storage and networking technologies, such as PCIe 4.0 and USB 4.0, achieve much higher data rates for demanding tasks.

Type	Peripheral	Data Rate
Human Interaction	Keyboard	~ kbps
Human Interaction	Mouse	~ kbps
Generic	Serial Port (RS-232)	115.2 kbps (max)
Generic	Parallel Port (LPT)	150 kbps
Generic	USB 4.0	20-40 Gbps
Generic	Bluetooth 5.0	2 Mbps
Generic	PCIe 4.0	16 Gbps per lane
Storage	SATA III (HDD/SSD)	6.0 Gbps
Storage	NVMe (PCIe 4.0)	64 Gbps (4-lane)
Networking	Ethernet (10BASE-T)	10 Mbps
Networking	10 Gigabit Ethernet (10GBASE-T)	10 Gbps
Networking	Wi-Fi 6 (802.11ax)	Up to 9.6 Gbps
Displays	VGA (analog video)	0.6-1.5 Gbps (approx.)
Displays	HDMI 2.1	48 Gbps
Optical Discs	CD-ROM	150 KB/s (1x) - 7.68 MB/s (52x)
Optical Discs	DVD-ROM	1.32 MB/s (1x) - 21.1 MB/s (16x)
Optical Discs	Blu-ray	4.5 MB/s (1x) - 54 MB/s (12x)

7.3.1 Accessing I/Os: Port-Mapped I/O (PMIO)

Port-Mapped I/O (PMIO) is a technique used to create a separate interface for Input/Output (I/O) operations, which is distinct from the memory interface. This method allows the CPU to access peripheral devices using dedicated I/O ports. In PMIO, specific control signals and new instructions are introduced to facilitate I/O operations.

Similar to the Memory Interface.



- **New Interface:** The CPU has a control interface for both memory (**ctrl-mem**) and I/O (**ctrl-IO**), along with an address bus (**addr**), read data bus (**rdata**), and write data bus (**wdata**).
- **Port Numbering and Control Signals:** The I/O ports are addressed separately using a dedicated **port** line (typically 8-16 bits wide), and additional control signals are introduced:
 - **CE** (Circuit Enable): Indicates that a valid port number is provided.
 - **OE** (Output Enable): Indicates that the I/O access is an output operation.
- **Data Buses:** The CPU communicates with I/O devices using dedicated **input** and **output** buses, which may not necessarily be 32 bits wide due to the limited number of peripheral devices.

Accessing I/Os: Memory Mapped I/O(MMIO)

Memory-Mapped I/O is a technique that leverages the same address space for both memory and I/O operations. This allows devices to be accessed using standard memory instructions, eliminating the need for special hardware or dedicated I/O instructions.

Address Space Allocation: In this configuration, specific memory addresses are allocated for different peripherals:

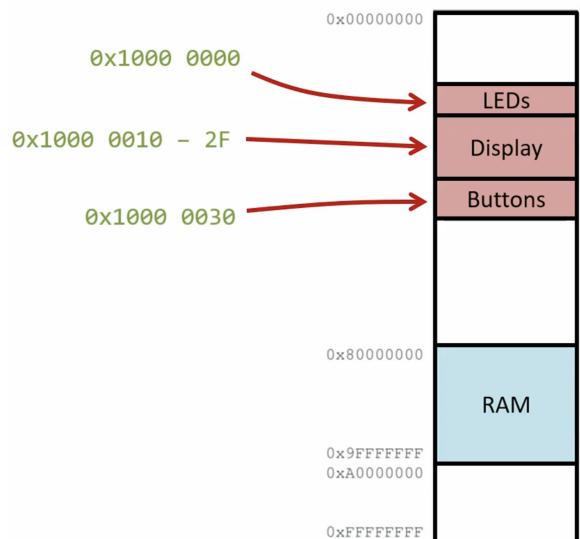
- 0x1000 0000: Base address for controlling LEDs.
- 0x1000 0010 to 0x1000 002F: Address range for controlling the display.
- 0x1000 0030: Base address for reading button states.

Standard Instructions: Since I/O devices are accessed as part of the memory space, standard load and store instructions can be used to interact with them. For example:

```

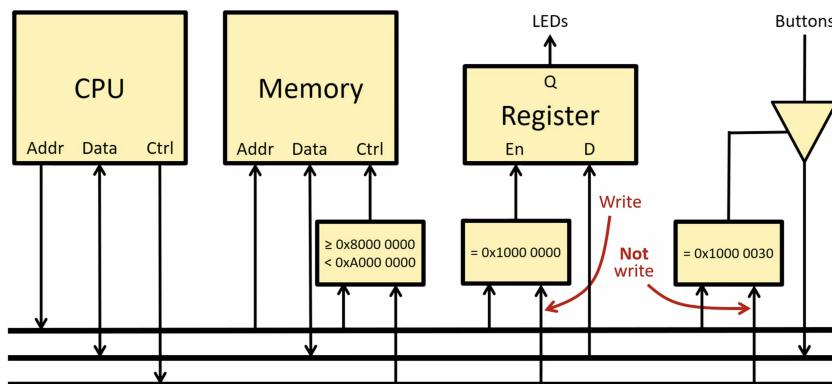
1 # Load upper immediate to set pointer to I/Os
2 lui t0, 0x10000
3 # Write value in t1 to the address of LEDs
4 sw t1, 0(t0)
5 # Read button states into t2
6 lw t2, 0x30(t0)

```



7.3.2 Memory Mapped I/O (MMIO)

In computer architecture, Memory-Mapped I/O (MMIO) is a technique used to access input and output devices. Instead of using separate I/O instructions, devices are assigned specific memory addresses. The CPU interacts with these devices by reading or writing data to these memory locations as if they were normal memory addresses.



Components

- **CPU:** Acts as the central processing unit that interacts with memory and registers. It uses address, data, and control buses to communicate.
- **Memory:** Represents the conventional memory address space accessible by the CPU. Specific memory ranges (e.g., between 0x8000 0000 and 0xA000 0000) are reserved for memory-mapped devices.
- **Register:** The CPU communicates with the register through a dedicated memory-mapped address (e.g., 0x1000 0000 for writing operations). This register drives outputs such as LEDs or accepts inputs from buttons.

Operation

The CPU interacts with memory and I/O devices through common buses. Depending on the address, data is either directed to regular memory or to an I/O device register. When the address matches a specific register (e.g., 0x1000 0000 for a write operation), the corresponding action is triggered, such as updating an LED state. In contrast, accessing 0x1000 0030 might perform a read operation, retrieving button states.

7.4 Example - A/D Converter

This example describes an Analog-to-Digital (A/D) Converter or (ADC) and its associated signals. The A/D Converter converts an analog input signal into a digital representation. The conversion process and signal behaviors are described below.

Signals

- **Start (START)**: This input signal, when active a new conversion begins.
- **Data Valid (/DV)**: This output signal indicates the validity of the data. When active, the output data bits (D7–D0) contain the converted digital value and are Valid.
- **Data (D7–D0)**: The output data bits representing the last conversion result in digital form.

7.4.1 Bus Interface

The A/D Converter interacts with an 8-bit processor using a simple bus interface. This bus interface allows data exchange and control signals to flow between the A/D Converter and the processor. The following signals are used:

- **Address (A23–A0)**: Output - Serves as the address bus to select a specific device or memory location.
- **Data (D7–D0)**: Bi-directional - Represents the data bus used for data exchange between the processor and the A/D Converter.
- **Address Strobe (/AS)**: Output - Indicates the presence of a **valid** address on the address bus during a memory access cycle.
- **Read/Write (R//W)**: Output - Determines the direction of the data flow (read from or write to the A/D Converter).
- **Data Acknowledge (/DTACK)**: Input - Signals the completion of a memory access by the A/D Converter when activated, indicating that the data is ready or has been latched.

The bus interface provides a simple mechanism for connecting the A/D Converter to the system bus, allowing the processor to initiate conversions and read the results.

7.4.2 Memory Mapping

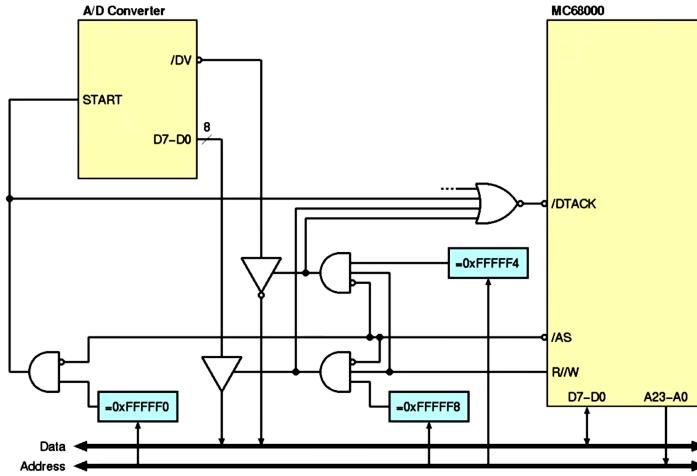
The A/D Converter is connected to the processor using a memory-mapped interface. Specific memory addresses are reserved for starting conversions, reading the data valid signal, and accessing the conversion result. The following address mapping is used:

- **0xFFFFF0**: Any access (read or write) to this address initiates a new conversion by the A/D Converter.
- **0xFFFFF4**: The processor reads the data valid signal from this address. Bit 0 of this location indicates whether the conversion result is ready.
- **0xFFFFF8**: The processor reads the conversion result from this address. The value stored here represents the digital output of the A/D Converter.

This memory-mapped interface simplifies the interaction between the processor and the A/D Converter by using standard read and write instructions to control the conversion process and retrieve the results.

7.4.3 Assembling everything

To get to this point, it is highly advised to first draw a timing diagram of the expected signals, and then start drawing connections, also, be careful, the notation "/AS" for example, means that the signal is active low.



Software Implementation

```

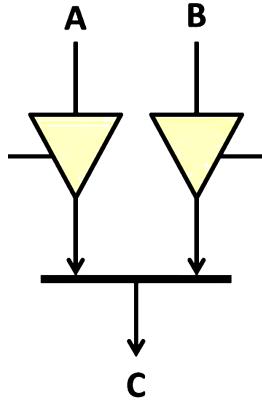
1 .section .data
2 START_ADDR:    .word 0xFFFFF0      # Address to initiate conversion
3 DATA_VALID_ADDR: .word 0xFFFFF4    # Address to check if data is valid
4 RESULT_ADDR:   .word 0xFFFFF8    # Address to read conversion result
5
6 .section .text
7 .globl _start
8
9 # Start of the main program
10 _start:
11     # Initiate A/D Conversion
12     lui t0, %hi(START_ADDR)      # Load upper 20 bits of START_ADDR
13     lw t1, %lo(START_ADDR)(t0)    # Load lower 12 bits into t1
14     # Write zero to start the conversion (writing to address 0xFFFFF0)
15     sw zero, 0(t1)
16
17 wait_for_data:
18     # Check if the data is valid
19     lui t0, %hi(DATA_VALID_ADDR) # Load upper 20 bits of DATA_VALID_ADDR
20     lw t1, %lo(DATA_VALID_ADDR)(t0) # Load lower 12 bits into t1
21     lw t2, 0(t1)                  # Load the data valid status into t2
22     andi t2, t2, 0x1             # Mask bit 0 (check if data is ready)
23     beq t2, zero, wait_for_data # If bit 0 is zero, data is not valid, wait
24
25     # Read the conversion result
26     lui t0, %hi(RESULT_ADDR)    # Load upper 20 bits of RESULT_ADDR
27     lw t1, %lo(RESULT_ADDR)(t0)  # Load lower 12 bits into t1
28     lw t3, 0(t1)                # Read conversion result into t3
29
30     # End of program (infinite loop)
31 end:
32     j end                      # Loop indefinitely

```

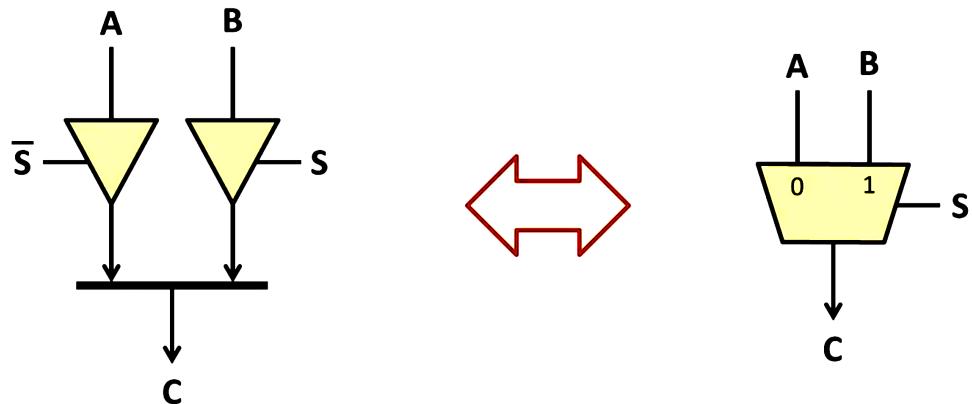
7.5 What do these tri-state buffers do?

Tri-state buffers are crucial components used to control data flow on a bus. They can exist in one of three states: high impedance (effectively disconnected), logic high, or logic low.

If not controlled properly, a tri-state buffer can cause bus contention, where multiple devices attempt to drive the bus simultaneously, leading to data corruption or damage.



In essence, a tri-state buffer acts like a decentralized multiplexer, functioning similarly to a multiplexer with a select line to manage data transmission.



7.5.1 A Classic UART

UART (Universal Asynchronous Receiver-Transmitter) is one of the simplest and most common communication peripherals, typically used to connect terminals to embedded devices.

Our UART employs a simple programmed I/O interface, consisting of four key registers:

- **Control register:** Configures the UART. Bit 7 must be set to 1 to enable the UART, while bits 2 to 0 determine the communication speed (e.g., 0b001 for 9600 baud).
- **Status register:** Provides the current status of the UART. Bit 1 indicates if data is available, and bit 0 signals if the UART is ready to transmit data.
- **Data input register:** Holds the received data available to the processor.
- **Data output register:** Contains data placed by the processor for transmission.

```

1 # Constants
2 UART_CTRL_ADDR      = 0x10000000 # UART control register address
3 UART_ENABLE_BIT     = 0x80       # Enable bit (bit 7)
4 UART_SPEED_9600     = 0x01       # Speed setting for 9600 baud (4 bits, [3:0])
5 UART_STATUS_ADDR    = 0x10000004 # UART status register address
6 TX_READY_BIT        = 0x01       # Transmitter ready bit (bit 0)
7 UART_DATAIN_ADDR   = 0x10000008 # UART data input (receive) register address
8 UART_DATAOUT_ADDR  = 0x10000008 # UART data output (send) register address
9
10 # Send a string using UART
11 send_string:
12     li t0, UART_CTRL_ADDR      # Load UART control register address into t0
13     li t1, UART_STATUS_ADDR    # Load UART status register address into t1
14     li t2, UART_DATAOUT_ADDR   # Load UART data output register address into t2
15     li t3, UART_ENABLE_BIT     # Load enable bit (0x80) into t3
16     li t4, UART_SPEED_9600     # Load speed setting (0x01) into t4
17     or t3, t3, t4             # Combine enable and speed bits into t3
18     sw t3, 0(t0)              # Configure UART by storing combined value into
19         control register
20
21 next_char:
22     lb t5, 0(a0)               # Load the current byte of the string into t5
23     beqz t5, finish            # If the byte is zero (null terminator), jump to
24         finish
25
26 check_tx_ready:
27     lw t6, 0(t1)               # Load the value of the UART status register into
28         t6
29     andi t6, t6, TX_READY_BIT  # Check if the TX ready bit is set
30     beqz t6, check_tx_ready    # If not ready, loop back to check again
31
32     sw t5, 4(t2)               # Store the character in UART data register
33     addi a0, a0, 1              # Move to the next character in the string
34     j next_char                # Jump to send the next character
35
36 finish:
37     ret                       # Return from function

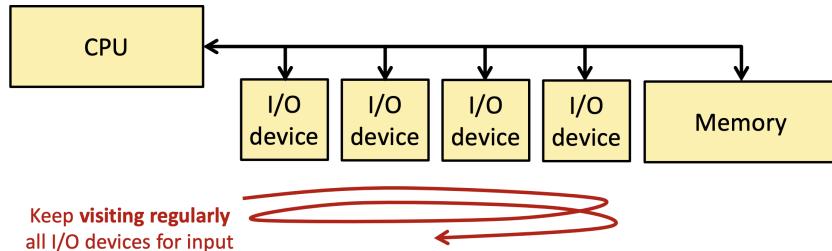
```

Chapter 8

Part II(c) - Interrupts

8.1 I/O Polling

I/O Polling is a method used by the CPU to check if any peripheral devices, such as a keyboard or network interface, have data to provide. The CPU continuously monitors each connected I/O device at regular intervals to see if they need attention.



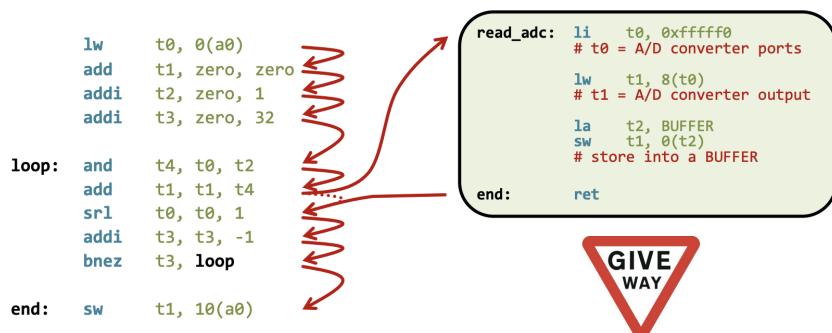
How It Works: The CPU keeps visiting each I/O device in a loop to check for input or status changes. This is known as "polling" the devices.

Drawbacks: This approach can be very resource-intensive. If a device operates at high speed and requires immediate handling, the CPU must check it frequently, which can consume significant processing time.

8.2 I/O Interrupts

Instead of continuously checking the status of peripherals, it is more efficient to have them request attention when needed. This approach minimizes CPU usage by eliminating the need for constant polling.

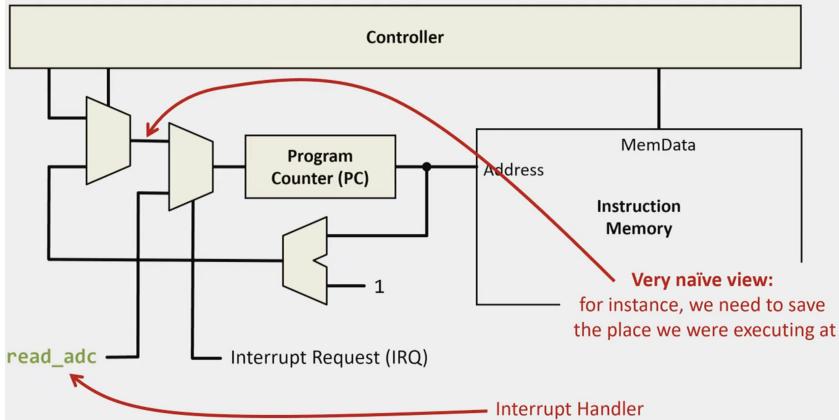
- **Polling Method:** The CPU checks the status of a peripheral device by repeatedly executing a loop to monitor the peripheral register. This approach requires continuous CPU attention, which can be inefficient in systems with multiple peripherals.
- **Interrupt Method:** In an interrupt-driven approach, peripherals alert the CPU only when they need attention. The CPU executes an interrupt service routine (ISR) to handle the request. This method allows the CPU to focus on other tasks until interrupted, improving efficiency.



8.2.1 The Basic Concept of I/O Interrupts

I/O interrupts provide a mechanism for a controller to handle external requests efficiently by temporarily diverting program execution.

This is not the actual implementation, but basic concept for you to help you understand what we're aiming for.



- **Interrupt Request (IRQ):** An interrupt signal is triggered, typically from an I/O device, to request immediate attention from the controller.
- **Program Counter (PC) Preservation:** The current value of the Program Counter (PC), which holds the address of the next instruction, is saved to allow resumption of normal execution after the interrupt is handled.
- **Interrupt Service Routine (ISR):** The controller redirects the PC to the address of the interrupt handler function, denoted here as `read_adc`. This function processes the interrupt by executing specific instructions related to the I/O request.
- **Instruction Memory Access:** The Instruction Memory is accessed to fetch instructions at the new PC address, executing the ISR for the interrupt.
- **Resuming Program Execution:** Once the interrupt has been serviced, the controller restores the saved PC value, allowing the program to continue from the point it was interrupted.

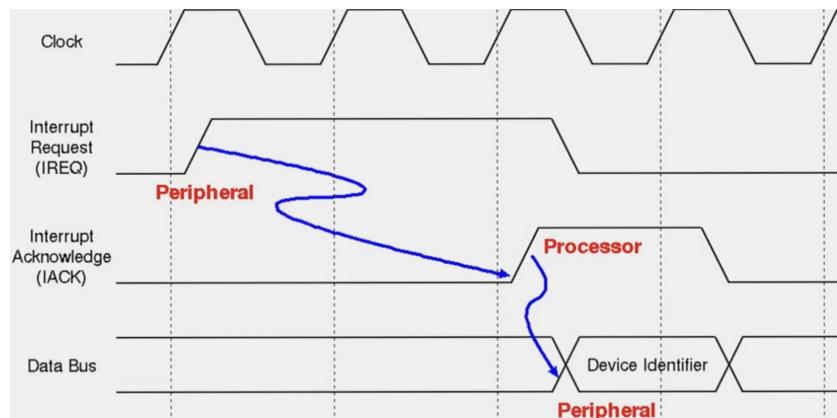
Considerations for I/O Interrupt Handling

When managing multiple I/O interrupts, several issues must be addressed:

- **Identifying the Source of the Interrupt:** In systems with multiple peripherals, it is essential to determine which device triggered the interrupt. This can be achieved through:
 - *Polling:* After an interrupt, the software checks each peripheral sequentially.
 - *Identification by the Peripheral:* The I/O peripheral itself sends an identification signal.
- **Handling Different Priorities:** Some interrupts may require immediate attention, while others can be delayed. Assigning priorities ensures critical interrupts are serviced promptly, while less urgent ones may wait.
- **Impact on Current Execution:** The system must decide whether to allow the current instruction(s) to complete before handling the interrupt or to pause immediately. This decision impacts program flow and execution timing.

8.2.2 Interrupt Cycle Description

The interrupt cycle is a sequence where a peripheral device signals an interrupt to the processor, which responds by acknowledging the interrupt and reading the device identifier from the data bus. The following signals are involved in this process:



- **Clock Signal**: The clock signal provides the timing for synchronization between the processor and peripherals.
- **Interrupt Request (IREQ)**: A peripheral asserts this signal to request service from the processor. When IREQ goes high, the processor detects an interrupt request.
- **Interrupt Acknowledge (IACK)**: In response to IREQ, the processor sends an acknowledgment signal (IACK) to the peripheral. This signal indicates that the processor is ready to handle the interrupt.
- **Data Bus**: After the IACK signal is asserted, the peripheral places its device identifier on the data bus, allowing the processor to identify the source of the interrupt.

The interrupt cycle proceeds as follows:

1. The peripheral raises the **IREQ** line to signal the interrupt request.
2. The processor detects the interrupt and, after some clock cycles, responds by asserting the **IACK** line.
3. The peripheral then places its **Device Identifier** on the data bus.
4. The processor reads the device identifier to determine the source of the interrupt and proceeds with the appropriate interrupt service routine.

This cycle ensures that the processor can handle asynchronous requests from peripheral devices in an organized and timely manner.

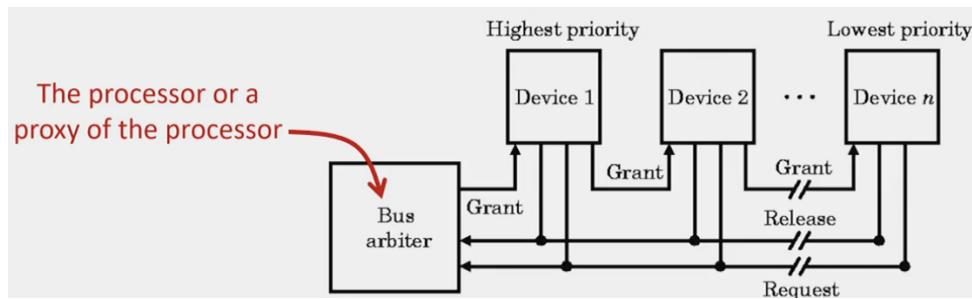
8.2.3 I/O Interrupt Priorities: Daisy Chain Arbitration

Daisy Chain Arbitration is a basic method used to manage I/O interrupt priorities. The process operates as follows:

- **Request Placement:** Any device can initiate a request to access the bus, indicated by signals such as **IREQ** (Interrupt Request).
- **Acknowledgment Line:** An acknowledgment signal, referred to as **IACK** or **Grant**, is sequentially passed from one device to the next.
- **Signal Interception:** The first device that requires access intercepts the acknowledgment signal, preventing it from being passed to devices further down the chain.

This method, while simple and easy to implement, has some limitations:

- **Slow Performance:** Due to the sequential nature of signal passing, response times can be slower as the chain length increases.
- **Fixed Priorities:** Devices closer to the bus arbiter have higher priority by design, leading to a rigid priority structure.



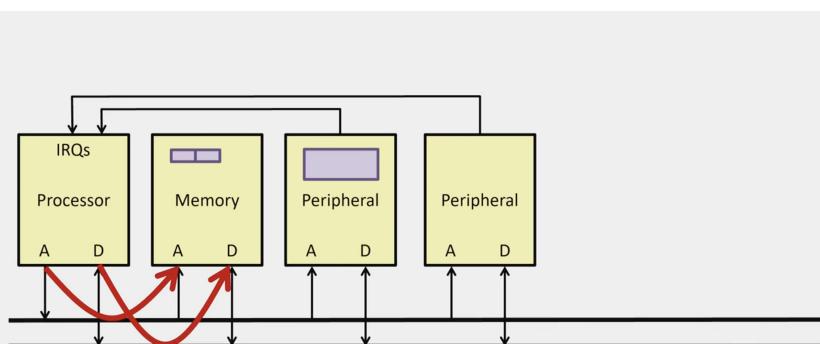
In this setup, the bus arbiter, which acts as the processor or a proxy for the processor, grants access in a priority chain from the highest-priority device to the lowest. This method is suitable for systems where simplicity is valued over flexibility and speed.

8.3 Direct Memory Access (DMA)

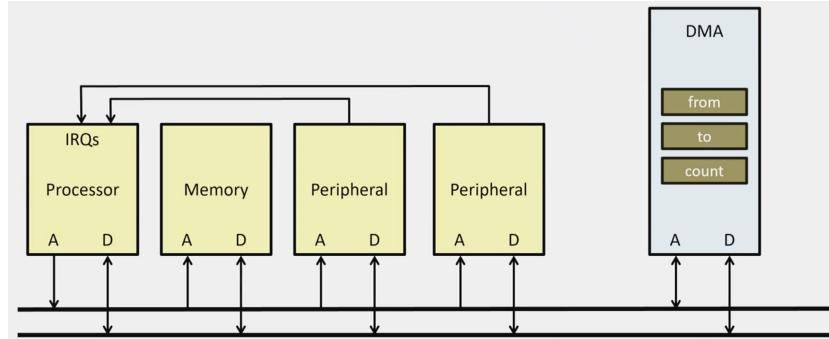
Direct Memory Access (*DMA*) is an efficient mechanism designed to offload the processor from managing repetitive and resource-intensive data transfers. Key considerations include:

- **Interrupts Efficiency:** *Interrupts* save the processor from continuously polling Input/Output (I/O) devices, allowing it to focus on computation.
- **Large Data Transfers:** Despite the use of interrupts, the processor may still spend considerable time transferring large chunks of data to and from high-throughput peripherals (e.g., disks, networks).
- **Solution:** A dedicated peripheral, known as the *DMA Controller*, is introduced. This controller autonomously handles data transfers between memory and peripherals (read/write operations), freeing the processor to focus on more critical tasks.

DMA significantly enhances system performance by reducing processor overhead during data transfer operations.

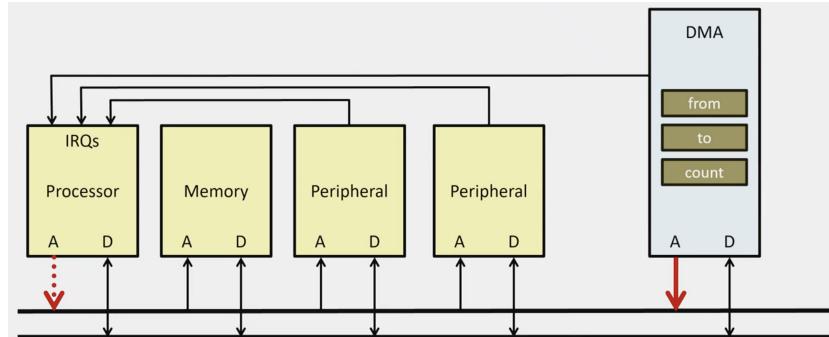


The diagram above illustrates a key inefficiency: *using the processor—a complex and expensive machine—to handle simple data transfer operations*. This inefficiency forms the basis for introducing DMA.



When initiating a data transfer, the **CPU** communicates with the **DMA Controller** to start the operation with a specific peripheral. The **DMA Controller** then handles communication with the **I/O device**, transferring the data to or from memory.

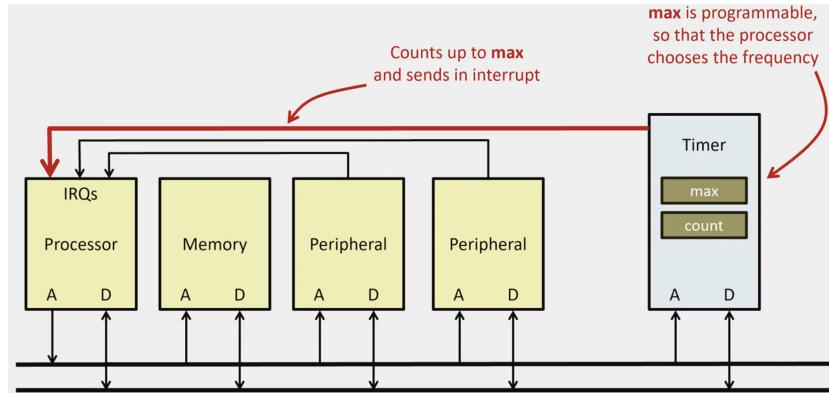
However, this introduces a new challenge. Previously, the *CPU* was the sole **master of the BUS**. With the **DMA Controller** capable of two-way communication with the BUS (on its A input), it also becomes a master of the BUS. This requires a mechanism to manage BUS access between the processor and the DMA Controller.



This is achieved using a **Tri-State Buffer**. However, during the data transfer, the processor is temporarily *disconnected from the BUS*, meaning it cannot track the progress of the transfer. Once the transfer is complete, an **interrupt** is required to notify the processor, ensuring it can resume operations with the updated data.

8.3.1 Timer and Interrupt Mechanism

A **timer** is a critical hardware component used to manage periodic tasks in embedded systems. The timer operates by incrementing a **count** register until it reaches a programmable **max** value, at which point it generates an **interrupt request (IRQ)**. The key features of this mechanism are outlined below:



- **Programmable Frequency:** The **max** value is configurable, allowing the processor to adjust the interrupt frequency based on system requirements.
- **Interrupt Handling:** Upon reaching the **max** value, the timer sends an IRQ signal to the processor. This allows the processor to execute specific tasks at regular intervals.
- **System Integration:** The timer interacts with the processor, memory, and peripherals via the system bus, ensuring synchronized operation.
- **Task Management:** Without a timer, it would be impossible for a processor to manage multiple tasks simultaneously, as there would be no mechanism to divide time between different operations. The timer enables multitasking by providing precise time slicing for task scheduling.

This mechanism is essential for time-sensitive operations such as task scheduling, event triggering, and real-time control in embedded systems, enabling efficient multitasking and coordination between components.

Chapter 9

Part II(d) - Processor, I/Os, and Exceptions

9.1 Exceptions, Interrupts, Faults, Traps, and Checks

Control Flow Under normal circumstances, the *control flow*—the sequence of instructions executed by a program—is fully determined by the programmer. This includes the use of jumps, branches, and procedure calls.

Exceptions Exceptions represent a deviation from the normal control flow. They are triggered by **special conditions** that are not explicitly defined in the program. When an exception occurs, the control flow changes unexpectedly, and the program must respond accordingly.

Exception Handlers To manage exceptions, *exception handlers* are invoked. These are specialized functions designed to take appropriate actions when an exception arises. An example of this is **I/O interrupts**, which signal specific events related to input/output operations.

Naming Conventions The terminology for exceptions and related events varies widely across systems. For clarity, we adopt the following convention based on RISC-V and the COD:

- **Exceptions:** A general term encompassing all control flow deviations.
- **Interrupts:** A specific type of exception generated outside the processor.

Thus far, interrupts are the only form of exception encountered.

9.1.1 Undefined Instruction

Undefined instructions are instructions that the controller does not recognize, as they do not correspond to any valid operation in the Instruction Register (IR). These scenarios require special handling to ensure system stability and proper exception processing.

- Detection: When an undefined instruction is detected in the IR, the controller generates a signal (`undef`) indicating the presence of an invalid operation.

- Exception Handling: The Program Counter (PC) is updated to the address of the Exception Handler to manage the undefined instruction. This involves:

- Saving the current PC for potential recovery.
- Redirecting the control flow to the exception handler's address using multiplexer logic.

- Control Logic: The system leverages the Next PC Logic to determine whether the next instruction comes from the regular PC logic or the exception handler, based on the `undef` signal or an external interrupt (IRQ).

- Synchronous Nature: These exceptions occur at a specific point in the program, precisely where the undefined instruction resides. This predictable behavior ensures that if the program is re-executed from the same initial state, the exception will occur at the exact same point, making debugging more straightforward.

- Immediate Handling: Serving the exception before executing the next instruction allows advanced features, such as efficient error recovery and the potential to extend system capabilities. This mechanism ensures that undefined instructions do not disrupt the execution flow and are handled systematically, enabling robust error recovery and system stability.

9.1.2 Optional fadd.s Instruction

Suppose we want to include a floating-point addition instruction, denoted as:

1 `fadd.s rd, rs1, rs2`

- Some processors might include a specialized ALU to support this instruction, whereas **cheaper processors do not**.
- For processors that lack support for this instruction, its execution would trigger an *undefined instruction exception*, which invokes a handler.
- The handler can **emulate** the behavior of the `fadd.s` instruction, ensuring compatibility across processors.

9.1.3 Outline of an Undefined Instruction Handler

To handle an undefined instruction, such as `fadd.s`, the following steps wouls be executed:

Save all registers on the stack that the handler or its callees might modify.

- Note: Standard calling conventions do not apply.

Retrieve the problematic instruction:

- If the program counter (PC) is saved, load the instruction from the corresponding address.

Decode the instruction in software and identify it as `fadd.s`.

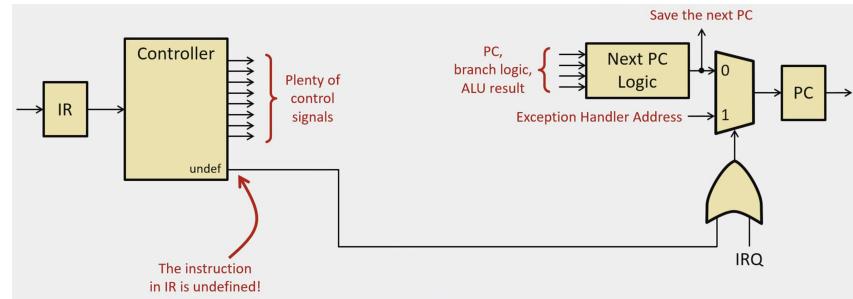
Read the source registers (operands) and either:

- Call a library function, or
- Implement the floating-point addition in software.

Store the result in the destination register.

Update the program counter (PC) to point to the next instruction.

Jump to the updated PC to resume execution.



9.2 Exceptions and Interrupts

Exceptions, interrupts, and related mechanisms handle critical events during execution. Key use cases include:

I/O Requests: Processing data or new inputs.

Timer Interrupts: Handling time-based events.

Undefined Instructions: E.g., unsupported floating-point operations.

Arithmetic Faults: Errors like division by zero.

Memory Violations: Unauthorized access to restricted memory.

Debugging: Breakpoints and execution control.

Hardware Failures: Malfunctions such as power loss.

9.2.1 A Possible Classification of Exceptions

Type	Synchronous?	Coerced?	Resume?
I/O request	Asynchronous	Coerced	Resume
Invoke OS	Synchronous	User requested	Resume
Trace instruction	Synchronous	User requested	Resume
Breakpoint	Synchronous	User requested	Resume
Page fault	Synchronous	Coerced	Resume
Misaligned access	Synchronous	Coerced	Resume
Memory protection violation	Synchronous	Coerced	Terminate
Bus error	Synchronous	Coerced	Terminate
Arithmetic fault	Synchronous	Coerced	Terminate
Undefined instruction	Synchronous	Coerced	Terminate
Hardware malfunction	Asynchronous	Coerced	Terminate
Power failure	Asynchronous	Coerced	Terminate

- **Synchronous?** Indicates whether the exception occurs as a direct result of the execution flow (synchronous) or independently of it (asynchronous).
- **Coerced?** Specifies whether the exception is forced by the system (coerced) or triggered by a user request.
- **Resume?** Denotes whether the system can continue executing after handling the exception (resume) or must terminate.