

# CS 250: Computer Architecture

## Final Exam

### Spring 2025

Benjamin Lobos Lertpunyaroj

*May 8th, 10:30<sub>AM</sub> – 12:30<sub>PM</sub>*

#### Exam contents and details for referencing

- Final exam is held in Fowler Hall on May 8th (Thursday), from 10:30 <sub>AM</sub> to 12:30 <sub>PM</sub>.
- Previous cumulative book chapters
  - Chapter 1 sections 1, 2, and 3.
  - Chapter 2, sections 1, 2, 3, 4, 5, 6, 7, **8, 11, and 12**.
  - Chapter 3, sections 1, 2, and 5.
  - Chapter 4, sections 1, 2, 3, 4, 5, 6, 7, and 8.
  - Chapter 5, sections 1, 2, 3, 4, and **7**.
  - Chapter 6, sections **1, 2, 3, 4, 5, 6, 7, and 8**.
  - Chapter 8 (Appendix A), sections 1, 2, 3 (but not PLAs or ROMs), 5, 7 (lightly), and 8.
  - **Slides Week 12.**
- The order of appearance of contents in this document is arbitrary.

#### Table of Content

Appendix A	2
Chapter 2	11
Chapter 3	16
Chapter 1	18
Chapter 4: Datapath	20
Chapter 4: Pipelining	25
Chapter 5	31
Week 12	42

## Appendix A

### Logic & Gates

An *asserted* signal is logically true, the *deasserted* is the opposite.

#### Two types of logic systems

**Combinational logic:** No memory in components, hence same output given same input.

**Sequential logic:** Memory in components, hence output depends on input and current memory state.



Figure 1: AND gate, OR gate, and inverter

The gates can be combined to form different forms of logic. An example of this is  $\overline{A} + \overline{B}$  which is equivalent to  $A \cdot \overline{B}$  by De Morgan's law, seen in Figure 2.

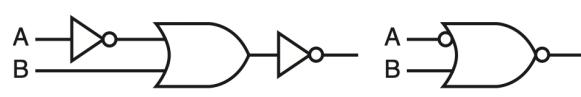


Figure 2: Logic gate implementation of example formula

### Decoders & Multiplexors

A **decoder** is a logic block that has an n-bit input and  $2^n$  outputs, where there is one unique true bit as output from a unique set of bytes of input.

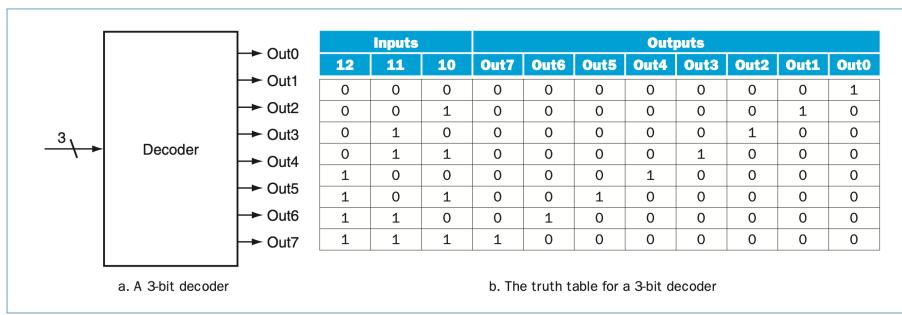


Figure 3: 3-bit input decoder that generates  $2^3 = 8$  different outputs (Out0 – Out7)

$$2^n \text{ outputs} \therefore \log_2(\text{output}) = \text{input bits}$$

Encoders are the other way around.

**Multiplexors** have a selector input (or control value), that will determine which inputs will become outputs.

In the case of the two-input MUX, its representation is the following,  $C = (A \cdot \bar{S}) + (B \cdot S)$ , using  $n$  (data inputs) AND gates, and one OR gate.

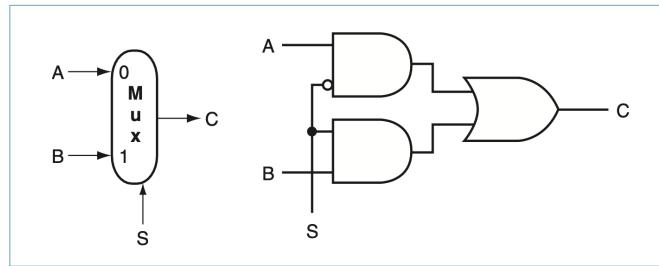


Figure 4: Two-input multiplexor that generates one output depending on the selector input  $S$

$$n \text{ (data inputs)} \therefore \log_2 n = S \text{ selector bits required to represent all inputs}$$

Often times a decoder generates  $n$  bits for a MUX, to be used as a selector signal.

## Buses

A collection of data lines that is treated as a single logical signal.

When showing a logic unit whose inputs and outputs are buses, the unit must be replicated a sufficient number of times to accommodate the width of the input.

You can use multiplexors to select between two buses, requiring  $n$  inputs to represent  $n$ -bit buses.

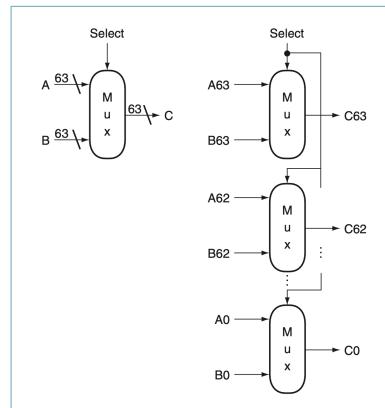


Figure 5: 1-bit multiplexors replicated 64 times to represent two 64-bit buses

## ALUs

### Operation done by the ALU

**Logic operations:** AND and OR gate operations, with NOR being available through an inversion of both input signals with **AInvert** and **BInvert** control signals.

**Arithmetic operations:** Addition and subtraction through the full adder, and **BInvert** control signal on one input for determining the type of operation.

The LEGv8 word is 64 bits wide, as such a 64 bit wide ALU is required (64 1-bit ALUs).

In its simplest form, a 1-bit logical unit for AND and OR operations simply requires a multiplexor and a one bit control signal to select between the two operations ( $2^1 = 2$ ).

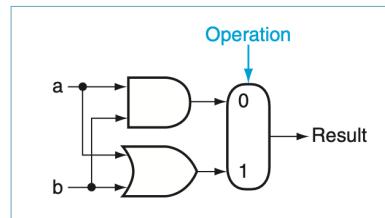


Figure 6: 1-bit logical unit for AND and OR operations

Implementing addition requires two input operands, one output, a **CarryIn** bit carried from the less-significant bits of the operation (i.e. another 1-bit logical unit), and a **CarryOut** bit to be carried forward to the next more significant bit.

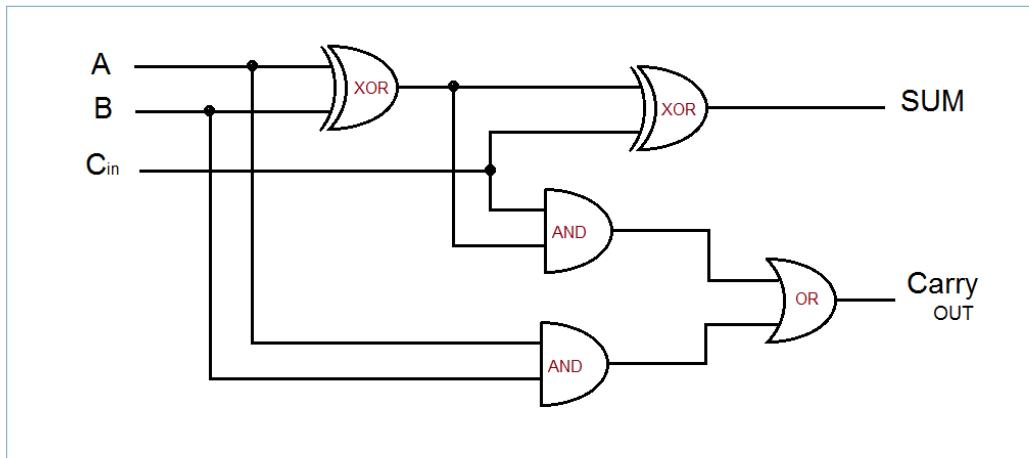


Figure 7: Full adder that performs mod 2 addition

The combination of the adder and the logic gates, coupled with a multiplexor with a control signal to determine the operation makes a complete 1-bit ALU, which can be seen in Figure 8.

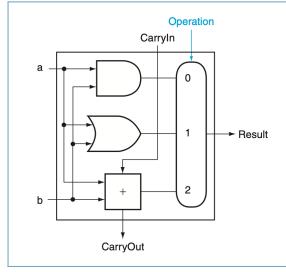


Figure 8: 1-bit alu with logical operations and addition

For expanding to a 64-bit ALU, the adders have to set up a ripple carry from the least to the most significant bit.

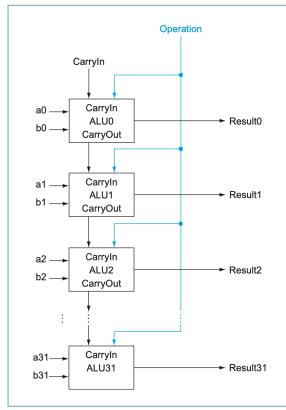


Figure 9: Ripple carry implemented for a 64-bit ALU

By inverting the second input (`BInvert` = 1, seen in Figure 10) and setting `CarryIn` to 1 in the least significant bit of the ALU, we get two's complement subtraction of  $b$  from  $a$ .

To implement a `NOR` function, existing components can be combined,  $(\overline{a + b}) = \overline{a} \cdot \overline{b}$  (DeMorgan's theorem), which means we need an `AND` and two inverters for both  $a$  and  $b$ , seen in Figure 10

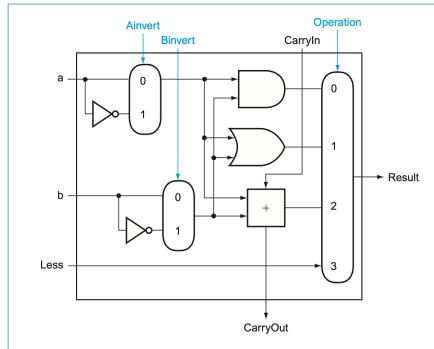


Figure 10: 1-bit ALU that performs subtraction, and NOR operations

On a 64-bit ALU we can use a zero flag to help with conditional branch instructions in LEGv8 (e.g.

CBZ), as they receive two inputs and require to test if the subtraction has a zero.

The following represents this with an inversion of an OR tree on all results from the subtraction considering a 64-bit subtraction, fully represented in Figure 11.

$$\text{Zero} = \overline{(R_0 + R_1 + R_2 + \dots + R_{63})}$$

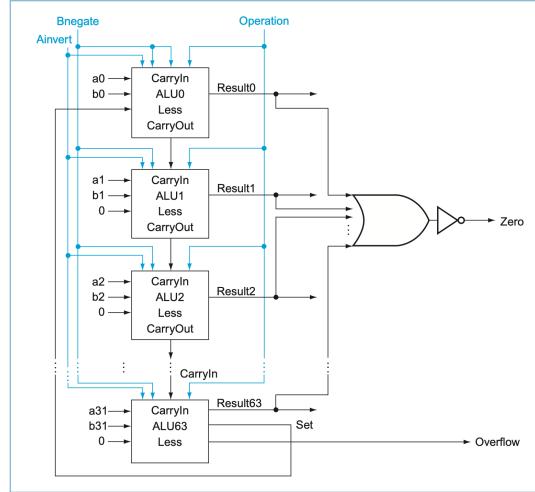


Figure 11: 64-bit ALU OR tree and an inverter for determining the Zero flag

For a generalized symbol of the ALU, Figure 12, where ALU operation is the control signal of the MUX that determines the type of operation.

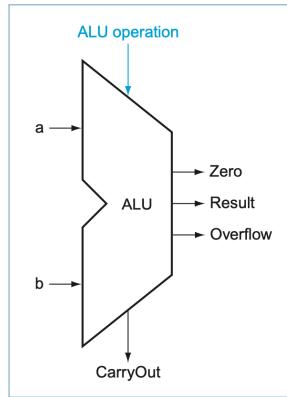


Figure 12: General symbol for an ALU or an adder

## Clocks

### Clocking methodology semantics

**Edge triggered clocking:** State changes occur on a clock edge.

**Synchronous system:** Type of memory system where data is read only when a clock signal indicates stability (i.e. non-changing value).

A combinational logic block receives an input and then generates an output for a state element which is updated on a clock edge.

An edge-triggered methodology allows a state element to be read and written in the same clock cycle without creating a race condition.

For this to work, the clock cycle must be long enough for the state element to have received a stable input before the next active clock edge.

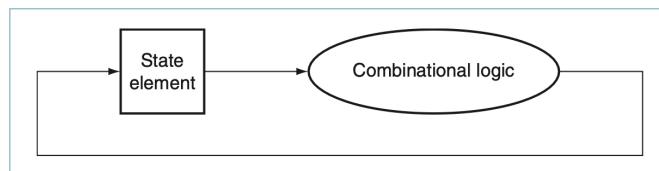


Figure 13: *Edge-triggered state element to be read and written to in one active clock edge*

One such state element is the register file.

## Flip-flops & Latches

### Types of clocked memory elements

**Flip-flops:** Edge-triggered element that changes the stored state only at a clock edge.

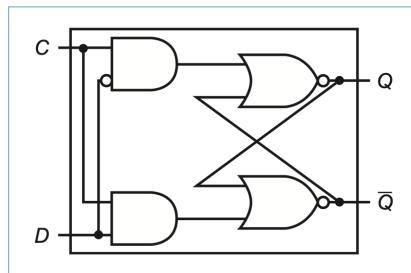
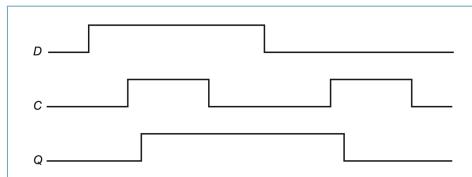
**Latches:** Level-sensitive element that changes the stored state at any time the clock is asserted.

Flip-flops are built upon latches and are going to be used in edge-triggered systems.

A **D flip-flop** or **D latch** is used for storing the value of one data input signal, in the internal memory, at the clock edge.

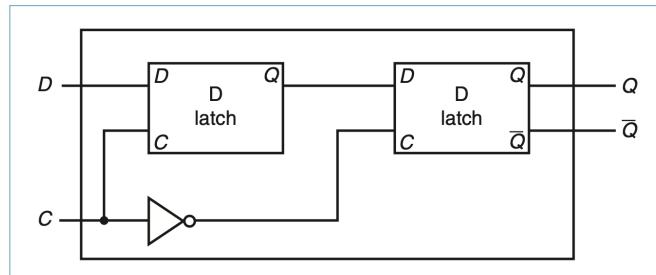
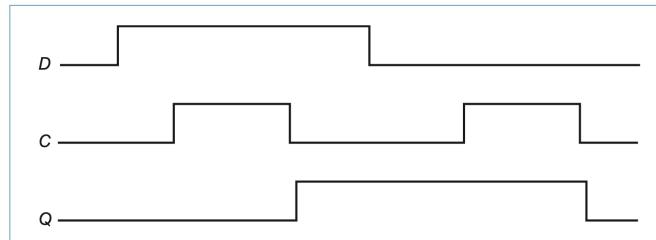
To implement a **D latch**, it requires two inputs, the data to be stored **D**, and the clock signal **C**, producing two outputs, the value of the internal state **Q**, and its complement **Q̄**.

The implementation has cross-coupled NOR gates that store the state value unless **C** is asserted, in which case **D** replaces the value of **Q** and is stored.

Figure 14: *D latch, composed of crossed NOR gates and a SR latch*Figure 15: *Progression of a D latch, assuming output is initially deasserted*

To implement a **D flip-flop**, with a falling-edge trigger, we can use two D latches, master and slave. Master sets input D when C is asserted. When C falls, master is closed, but slave is open and gets its input from master's Q.

In this sense, the rising-edge represents when the master takes in the D value, and the falling-edge represents when the slave takes in the master's D producing the final Q.

Figure 16: *D flip-flop with a falling-edge trigger made from two D latches, master and slave*Figure 17: *Progression of a D flip-flop with a falling-edge trigger, where output is initially deasserted*

The minimum time D must retain a valid input is the setup time plus the hold time (after edge).

## Register files

A **register file** consists of a bunch of **registers** that can be read and written to, and a **WriteReg** control signal (clock).

For writing it requires the control signal, the number of the register to write to (**Write register**), and the data to write (**Write data**).

For reading it requires the numbers of the registers to read from (**Read register number 1 & 2**), and it outputs the read contents from two registers (**Read data 1 & 2**).

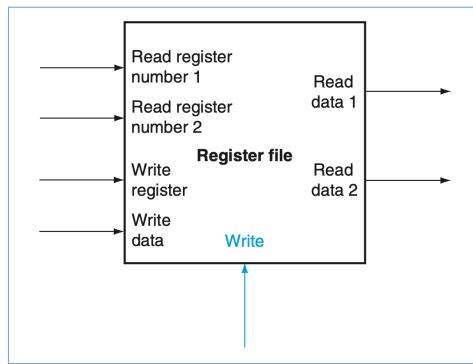


Figure 18: *Register file with two read ports and one write port*

The implementation for the write port consists of a decoder that will select one of the  $n - 1$  registers that will be ANDed with the **WriteReg** signal to act as the **C** input for the registers (**D** flip-flops).

The **D** input for every register is the **Write data** input from the reg. file.

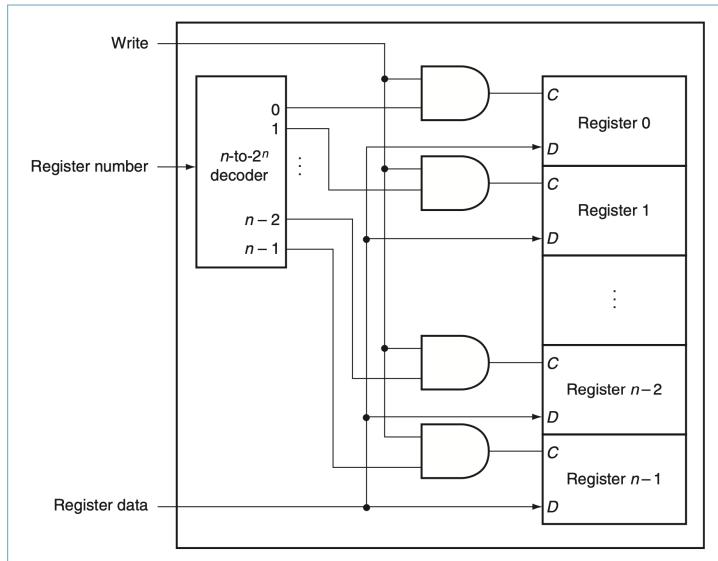


Figure 19: *Write implementation in the register file*

The implementation for the two read ports consists of using the stored state of the registers ( $Q$  output), as inputs for two different MUXes that use the **Read register number 1 & 2 reg. file** inputs as control signals to output the information of the two registers specified.

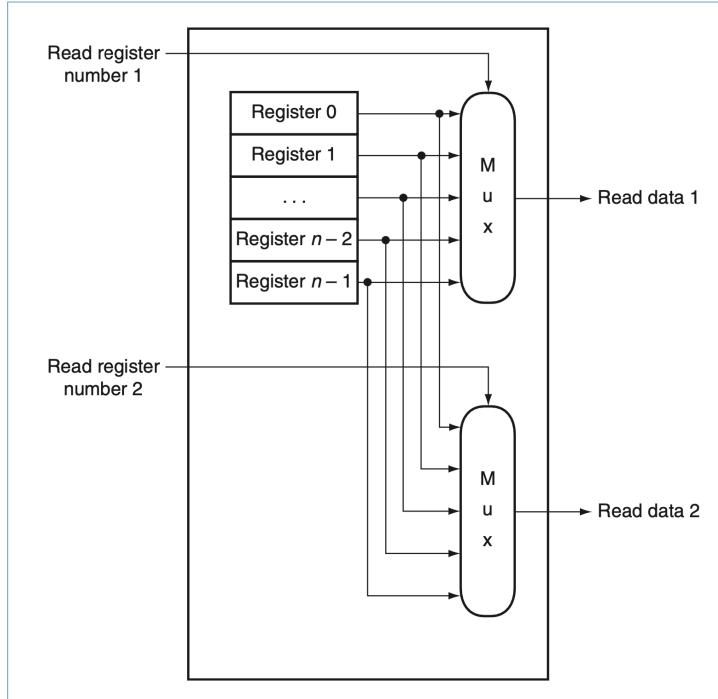


Figure 20: *Read implementation in the register file*

## Chapter 2

The size of a register in LEGv8 is 64 bits, which are denominated as doublewords (8 bytes).

A word, the natural unit of access in a computer, is 32 bits (4 bytes).

### LEGv8 Assembly

#### Relevant LEGv8 instructions

Addition:	ADD X1, X2, X3
Subtraction:	SUB X1, X2, X3
Add immediate:	ADDI X1, X2, #20
Subtract immediate:	SUBI X1, X2, #20
Load register:	LDUR X1, [X2, #20] // Load mem. at addrs. X2 + 20
Store register:	STUR X1, [X2, #20] // Store at addrs. X2 + 20
Logical AND:	AND X1, X2, X3
Logical Inclusive OR:	ORR X1, X2, X3
Logical Exclusive OR:	EOR X1, X2, X3
Bitwise Left:	LSL X1, X2, #20
Bitwise Right:	LSR X1, X2, #20
Conditional branch is 0:	CBZ X1, L0 // Immediate in operand is in word bytes
Conditional branch not 0:	CBNZ X1, L0 // Immediate in operand is in word bytes
Branch:	B L0 // Immediate in operand is in word bytes

The above list excludes instructions that set flags, e.g. ADDS, SUBS, SUBIS, ADDIS, and zero independent conditional branching, e.g. B.cond.

### LEGv8 architecture design

There are 32 64-bit registers, limited as larger number of registers implies increasing clock cycle time as electronic signals take longer as they must travel further. It also makes instruction formats more constrained.

Many architectures have alignment restriction that establish that words must start at addresses that are multiples of 4 and doublewords at multiples of 8.

## Signed and unsigned numbers

### Denomination of bits in a bit string

**Least significant bit:** The rightmost bit in a bit string. In a doubleword this is bit 63.

**Most significant bit:** The leftmost bit in a bit string. In a doubleword this is bit 0.

A **sign and magnitude** representation of a signed number has one bit set aside to represent the sign of the integer. Issues come with having a **-0** and **+0**, and more steps needed to determine the sign by the adder.

For making the hardware simple, **two's complement** representation is used, which used leading 0s to denote positives, and leading 1s to denote negatives. This means that the most significant bit can be used as a sign bit, making conversions the following way.

$$(x_{63} \cdot -2^{63}) + (x_{62} \cdot 2^{62}) + (x_{61} \cdot 2^{61}) + \dots + (x_1 \cdot 2^1) + (x_0 \cdot 2^0)$$

As an example,

$$11111111\ 11111111\ 11111111\ 11111111\ 11111111\ 11111111\ 11111111\ 11111100_{\text{two}} = -4_{\text{ten}}$$

Overflow occurs when the sign bit gets overridden, i.e. 0 when negative, 1 when positive.

Using two's complement, sign extension, used for conditional branching, just requires to copy the sign bit  $m - n$  bits over, where  $m$  is the new size of the bit string.

### Power of 2 prefix definitions

**Kibibyte (Kib):**  $2^{10}$  bytes.

**Mebibyte (Mib):**  $2^{20}$  bytes.

**Gibibyte (Gib):**  $2^{30}$  bytes.

**Tebibyte (Tib):**  $2^{40}$  bytes.

**Pebibyte (Pib):**  $2^{50}$  bytes.

## Instruction format fields

Opcode 11 bits	Rm 5 bits	shamt 6 bits	Rn 5 bits	Rd 5 bits
-------------------	--------------	-----------------	--------------	--------------

Table 1: *R-type format*, ADD, SUB, LSL, LSR

Opcode 11 bits	Address 9 bits	Op2 2 bits	Rn 5 bits	Rt 5 bits
-------------------	-------------------	---------------	--------------	--------------

Table 2: *D-type format*, LDUR, STUR

Opcode 10 bits	Immediate 12 bits	Rn 5 bits	Rd 5 bits
-------------------	----------------------	--------------	--------------

Table 3: *I-type format*, ADDI, SUBI

Opcode 8 bits	Address 19 bits	Rd 5 bits
------------------	--------------------	--------------

Table 4: *Conditional branch format*, CB(N)Z

Opcode 6 bits	Address 26 bits
------------------	--------------------

Table 5: *Unconditional branch format*, B

All instruction field formats are 32 bits, thus the name *32-bit instructions*.

The **opcode** denotes the operation and format of an instruction.

The register operands, **Rn**, **Rm**, and **Rd** are 5 bits to represent the 32 registers ( $\log_2 32 = 5$ ).

The **address** field can represent a region of  $\pm 2^8$  bytes around the base register **Rn**.

The **shamt** field represents the shift amount used by bit-shift instructions.

## Conditional Statement & Loops

### Simple conditional statement implementation

```

1      CBNZ X3, Else      // if (X3 == 0) {
2      ADD  X0, X1, X2      //     X0 = X1 + X2;
3      B    Exit          // }
4 Else:                           // else {
5      SUB  X0, X1, X2      //     X0 = X1 - X2;
6 Exit:                         // }
```

### Simple while loop implementation

```

1      ADDI X0, XZR, #10   // X0 = 10;
2 Loop:
3      CBZ  X0, Exit      // while (X0 != 0) {
4      SUBI X0, X0, #1      //     X0 -= 1;
5      ADD  X1, X2, X0      //     X1 = X2 + X0;
6      B    Loop          // }
7 Exit:
```

For a more complete example, the following converts the C code into LEGv8 assembly instructions.

## C code loop into LEGv8 ASM

```

1 // C
2 int i = 0;
3 int k = 10;
4 while (save[i] == k) {
5     i += 1;
6 }
7
8 // X25: Address of save[]
9
10 // Textbook LEGv8 ASM
11     ADDI X22, XZR, #0    // EOR X22, X22, X22
12     ADDI X14, XZR, #10
13 Loop:
14     LSL  X10, X22, #2
15     ADD  X10, X10, X25
16     LDUR X11, [X10, #0]
17     SUB  X12, X10, X14
18     CBNZ X12, Exit
19     ADDI X22, X22, #1
20     B    Loop
21 Exit:
22
23 // Alternate LEGv8 ASM
24     ADD  X10, XZR, X25
25     ADDI X11, XZR, #10
26 Loop:
27     LDUR X12, [X10, #0]
28     SUB  X12, X12, X11
29     CBNZ X12, Exit
30     ADDI X10, X10, #8
31     B    Loop
32 Exit:

```

You can "throw away" the result by writing into XZR, the zero register.

**Set Flag Branching**

The conditional branch instruction, B.cond, allows for .cond to be used for signed comparisons, EQ (equals), NE (not equal), LT (less than), LE (less than or equal), GT, or GE.

Also it allows for unsigned comparisons, LO, LS (lower or same), HI, or HS (higher or same).

The flag for comparison is set by instructions like, ADDS, ADDIS, or SUBS, but they are limited in number as they create dependencies that obstruct pipelining execution.

A use case is bounds checking, that is,  $0 \leq x < y$ .

## Bounds checking shortcut in LEGv8 ASM

```
1 // Pseudo C
2 if (X20 >= X11 || X20 < 0) {
3     goto Error;
4 }
5
6 // LEGv8 ASM
7     SUBS XZR, X20, X11
8     B.HS Error
9     B    Exit
10 Error:
11     // Error handling
12 Exit:
```

Signed negative numbers look massive when looked at from an unsigned comparison in two's complement, and it also allows for checking if it is below a number (X11 in example).

## Chapter 3

### Floating point representation

Compromise must be found between the size of the *fraction* and the *exponent*.

#### Tradeoff of floating-point representation

**Precision:** Increased by an increase in the size for the *mantissa* or *fraction*.

**Range:** Increased by an increase in the size for the *exponent*.

A LEGv8 implementation of **floating-point** numbers has **s** as the sign bit, **exponent** an 8-bit with bias, and **fraction** a 23-bit number.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
<b>s</b>	<b>exponent</b>																														
1 bit	8 bits																														
	23 bits																														

Figure 21: *LEGv8 floating-point representation*

$$(-1)^S \times F \times 2^E$$

Range of representation is  $2_{\text{ten}} \times 10^{38}$  considering  $2^7$  (exponent 8-bit rep. divided by 2 for bias)  $\approx 38$

Overflow entails the exponent being too large for the **exponent** (E) field to represent it. Underflow is the same situation but the **exponent** field representing a negative value.

To reduce the chances of underflow or overflow we use **double precision**, which is a floating-point value represented in a 64-bit doubleword.

63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32
<b>s</b>	<b>exponent</b>																														
1 bit	11 bits																														
	20 bits																														
	<b>fraction</b>																														
	32 bits																														

Figure 22: *LEGv8 double precision floating-point representation*

In the case of an overflow/underflow an interrupt (unscheduled disruption call) saves the address of the instruction (to resume after correction) and jumps to a predefined address that deals with the exception.

## IEEE 754 Floating-Point Standard Specifications

The floating-point can always have a leading  $1.0$ , as such IEEE assumes it as a hidden bit increasing the representation for the **fraction** by 1 (23 to 24 in single p., or 52 to 53 in double p.).

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

$$(-1)^S \times (1 + (s_1 \times 2^{-1}) + (s_2 \times 2^{-2}) + \dots) \times 2^E$$

IEEE 754 has **NaN** for invalid operations (e.g.  $\frac{0}{0}$ ). Infinities can be represented through the largest exponent representations ( $+/-$ ) instead of interrupting.

The sign bit is the most significant bit to easily process integer comparisons.

For representing the exponent without needing a sign bit, IEEE 754 uses a bias of 127 for single, and 1023 for double, that subtracts from the total possible representation to denote negative and positive exponents.

So for single precision, an exponent of  $-1$  is  $-1 + 127_{\text{ten}} = 126_{\text{ten}}$ , and  $+1$  is  $1 + 127_{\text{ten}} = 128_{\text{ten}}$

## Floating-Point Addition

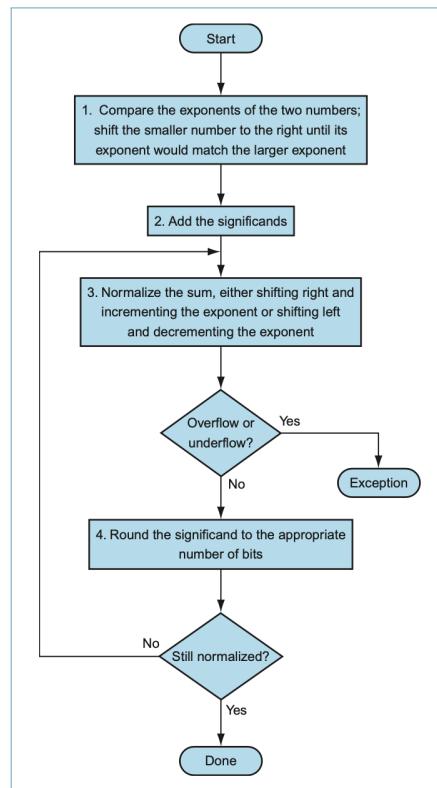


Figure 23: *Floating-point addition logic path*

# Chapter 1

## Design principles

### → **Design with Moore's Law**

Moore's Law states that an integrated circuit resources double every 18–24 month.

Thus, computer architects must anticipate where the technology will be when the design finishes rather than design for where it starts.

### → **Abstraction to Simplify Design**

Abstractions to hide lower-level details to simplify higher-level representations (more productive).

### → **Common case fast**

If you know what the common case is, enhancing it over the rarer cases is more optimal, as it involves more and is usually easier to enhance.

### → **Performance via Parallelism**

Computing operations in parallel (at the same time) increases performance over a time period.

### → **Performance via Pipelining**

Pipelining is a pattern of parallelism that involves increasing the throughput of a certain process (i.e. moving up the chain faster).

### → **Performance via Prediction**

If you can predict successfully to a degree you can anticipate instead of wait. In some cases it can be more beneficial to guess than to stall when misprediction costs are reasonable.

### → **Memory hierarchy**

The most expensive and fast memory per bit is at the top, and the opposite is at the bottom. Caches (top) can give the illusion of fast main memory (bottom).

### → **Dependability via Redundancy**

Systems become dependable when there is redundancy to account for solving failures and failure detections.

## Amdahl's Law

Related to the design principle of *making the common case fast*, you only benefit from improving the part you actually use, which is illustrated in **Amdahl's law**.

### Amdahl's Law

**Law:** Performance enhancement possible with a given improvement is limited by the amount that the improved feature is used.

$$S_{\text{total}} = \frac{1}{(1-f) + \frac{f}{S_E}}$$

Splitting up the original time into the unimprovable part  $1 - f$  and the improvable part  $f$ , speeding up the  $f$  portion by  $S_E$  yields a time ratio of

$$\frac{T_{\text{new}}}{T_{\text{old}}} \rightarrow (1-f) + \frac{f}{S_E}$$

As  $S_E \rightarrow \infty$ , the maximum overall speedup is  $\frac{1}{1-f}$ .

Thus, taken as the ratio between the old and new speeds, the formula for Amdahl's Law emerges.

## Chapter 4: Datapath

### Logic Design Elements

**Combinational elements:** Do not store any information internally, such as an ALU, and logic gates.

**State elements:** Have internal memory, a *state*, that remembers even in power loss, such as registers or memory.

### Branching in the Datapath

Conditional branching is done by incrementing the PC (program counter, keeps track of the current instruction) by the word (4 bytes) aligned offset (64-bits).

B offset      CBZ    Rd, offset

The offset varies in bits by instruction, requiring a sign-extension (from 32 to 64 bits), and 2 left shifts (word aligned multiple of 4 =  $2^2$ ).

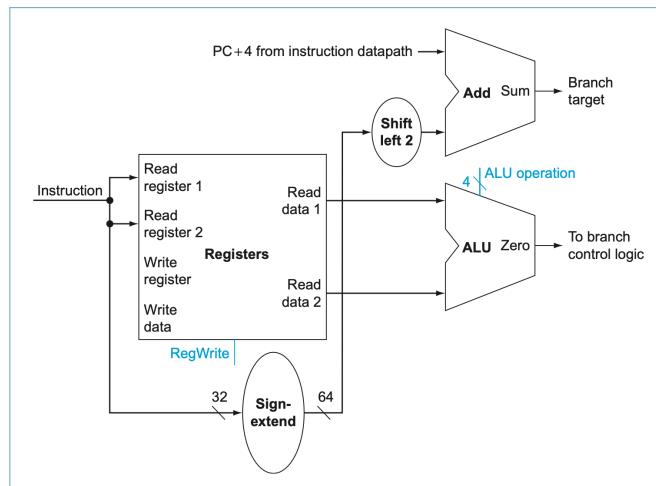


Figure 24: *Subset datapath for unconditional and conditional branching*

Control logic is used to decide whether the incremented PC or branch target should replace the PC, based on the zero output of the ALU, seen in Figure 24.

### Data memory unit and sign extensions

The sign extend takes the whole 32-bit instruction and selects the specific offset field size (e.g. 9 bits for loads and stores (small offset needed), 19 bits for CBZ), and fills upper bits with the most

significant bit of the extracted address.

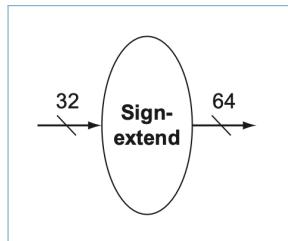


Figure 25: *Sign extension from 32 bits to 64 bits*

For memory instructions (e.g. LDUR, STUR), we need to write and read from memory using the data memory unit.

The structure takes in the address of memory to access, write input (when storing), read output (when loading), and read and write control signals.

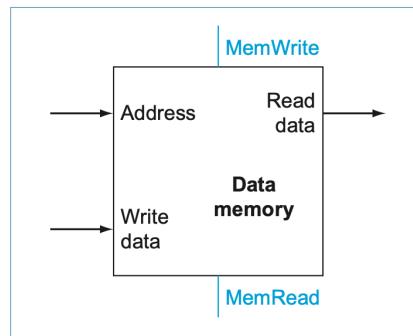


Figure 26: *Memory data unit diagram*

## R-type, Branch, and Memory Instruction Datapath

Building the datapath for R-type (ALU operations), conditional branching, and memory instructions in a single clock cycle.

### MUX Controlled by the ALUSrc Signal

**R-type:** Uses the value to do an ALU operation with the first register read output.

**Load/Store (D-type):** Uses the sign extended offset with the first output, that holds the memory address to load/store from. In this case, the second read register output is still relevant in storing, as it holds the data to be written.

## MUX Controlled by the MemtoReg Signal

**R-type:** Skips the reading stage from the data memory unit and sets the output of the ALU to be the `Write Data` input of the register file.

**Store (D-type):** Reads at address + offset from the data memory unit, and redirects the output to the Write Data input of the register file.

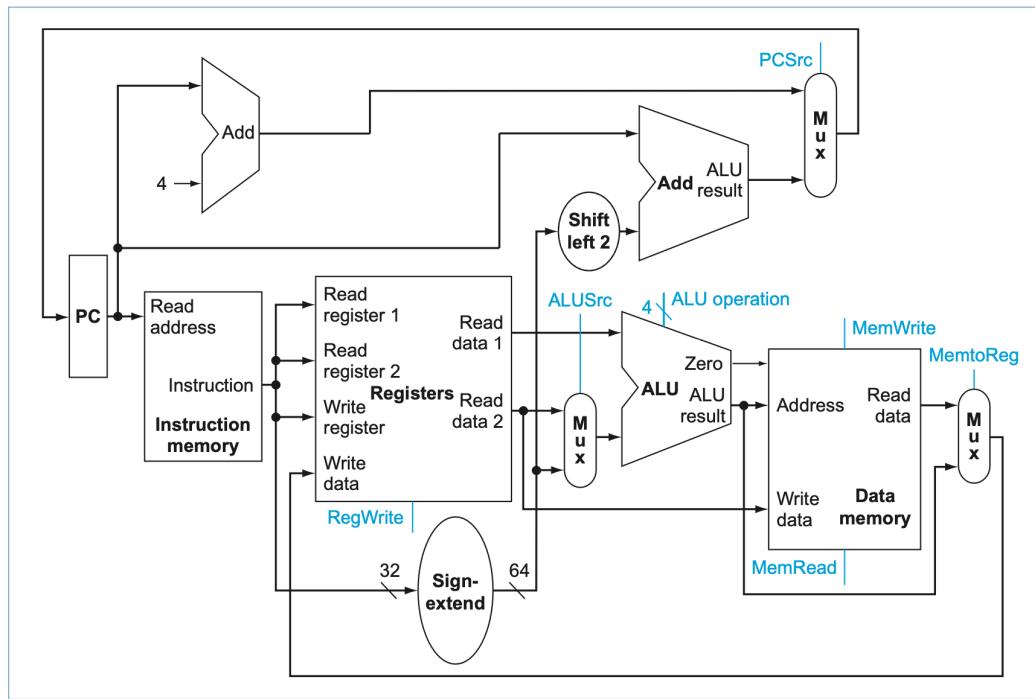


Figure 27: Datapath for R-type, load/store, and branching instructions

## Instruction Partitions by Instruction Memory

Considering the R-type, load/store, and conditional branching formats, they are broken down as seen in Figure 28.

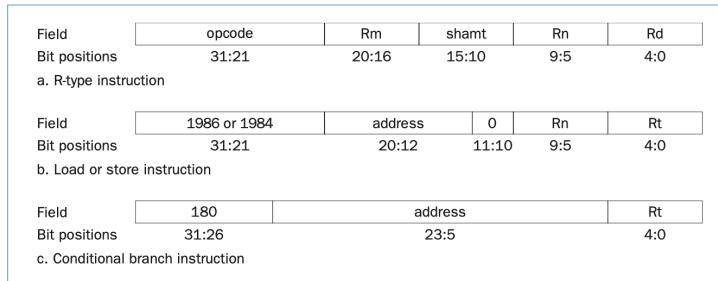


Figure 28: *Memory data unit diagram*

*Simplicity favors regularity:* The more uniform the design, the simpler the control logic for decoding the instruction (Instruction Memory).

## Generating the ALU Control Input

The ALU control is made up from the opcode (11-bit) and a 2-bit control signal ALUOp.

The opcode field only has three bits that matter when determining function (30, 29, 24).

ALUOp encodes for addition ( $0\ 0$ ), passing input B considering CBZ ( $X\ 1$ ), and other functions ( $1\ X$ ).

$$\text{ALU operation (4)} = \text{ALUOp (2)} + \text{opcode (11 but uses 3)}$$

ALUOp		Opcode field											Operation
ALUOp1	ALUOp0	I[31]	I[30]	I[29]	I[28]	I[27]	I[26]	I[25]	I[24]	I[23]	I[22]	I[21]	
0	0	X	X	X	X	X	X	X	X	X	X	X	0010
X	1	X	X	X	X	X	X	X	X	X	X	X	0111
1	X	1	0	0	0	1	0	1	1	0	0	0	0010
1	X	1	1	0	0	1	0	1	1	0	0	0	0110
1	X	1	0	0	0	1	0	1	0	0	0	0	0000
1	X	1	0	1	0	1	0	1	0	0	0	0	0001

Figure 29: Truth table for the 4 ALU control bits, with additional "don't care entries"

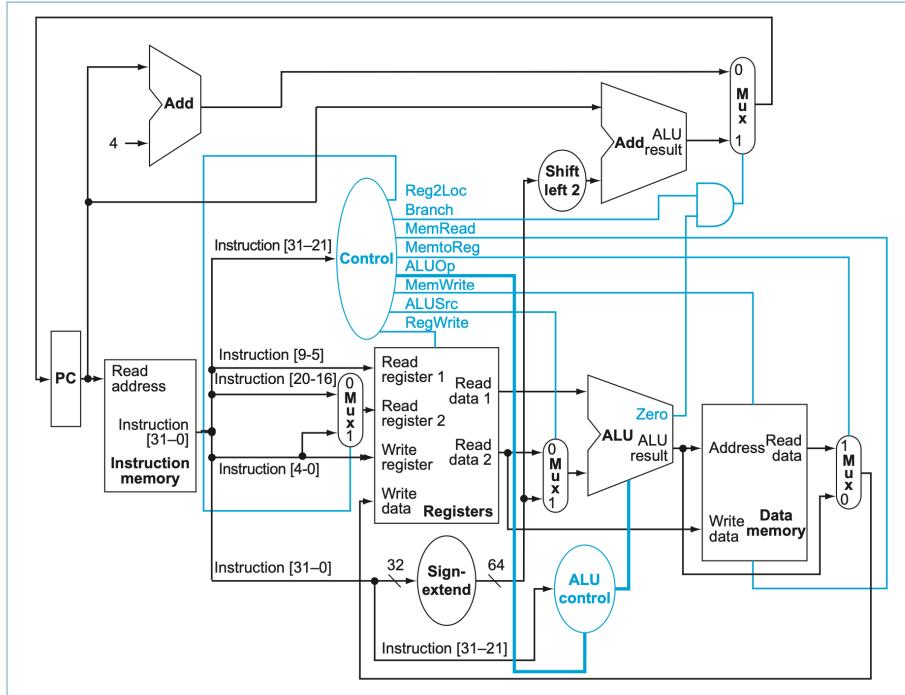


Figure 30: Datapath for R-type, D-type, and branching instructions with all 9 control lines

Consider in Figure 30, the PCSrc signal, that determines whether to increment PC by 4 or by the

offset specified, does not originate from the **control unit**, but rather the **branch signal** from the control unit being asserted and the zero output of the ALU is asserted.

Instruction	Reg2Loc	ALUSrc	MemtoReg	RegWrite	MemRead	MemWrite	Branch	ALUOp1	ALUOp0
R-format	0	0	0	1	0	0	0	1	0
LDUR	X	1	1	1	1	0	0	0	0
STUR	1	1	X	0	0	1	0	0	0
CBZ	1	0	X	0	0	0	1	0	1

Figure 31: Setting control lines according to desired instruction (determined by opcode)

## Implementing Unconditional Branching into Datapath

Adding a new signal **UncondBranch** to the control unit and ORing it with the resulting **PCSrc** control signal that emerges from the branch control signal.

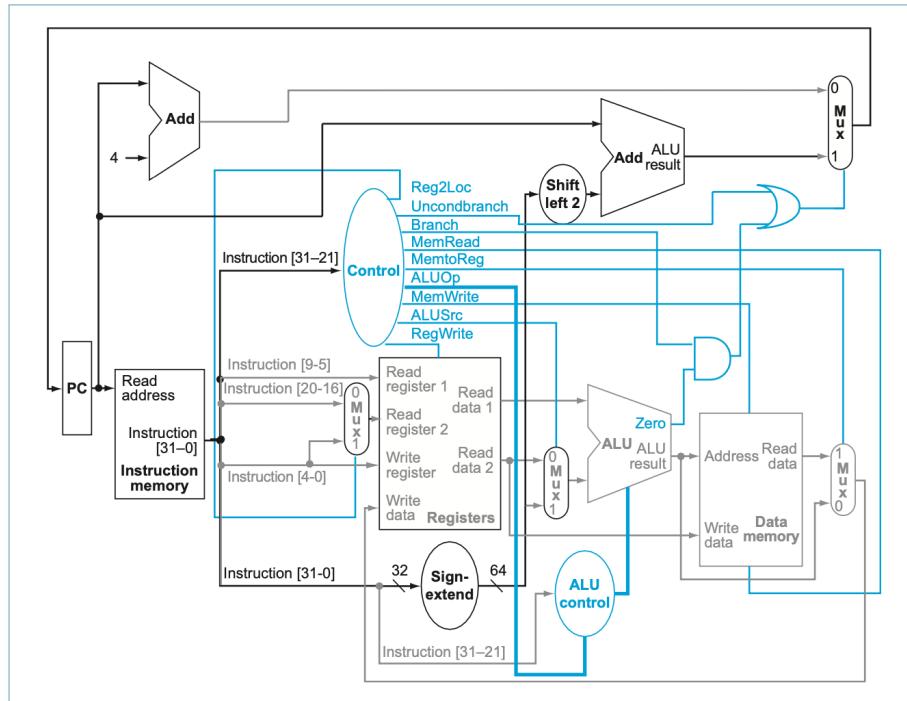


Figure 32: Datapath with UncondBranch control signal

Single cycle implementations have the limitation that the clock cycle needs to be greater than the time it takes to complete the longest possible path (load instruction: uses the IM, RegFile, ALU, DataMemory, and the RegFile again).

## Chapter 4: Pipelining

Pipelining overlaps instructions in execution.

### LEGv8 Pipelining Steps

1. **IF**: Fetch instructions from memory.
2. **ID**: Decode instructions and read registers.
3. **EX**: Executes the op. or calculates the address.
4. **MEM**: Access data memory (if needed).
5. **WB**: Writing back into register (if needed).

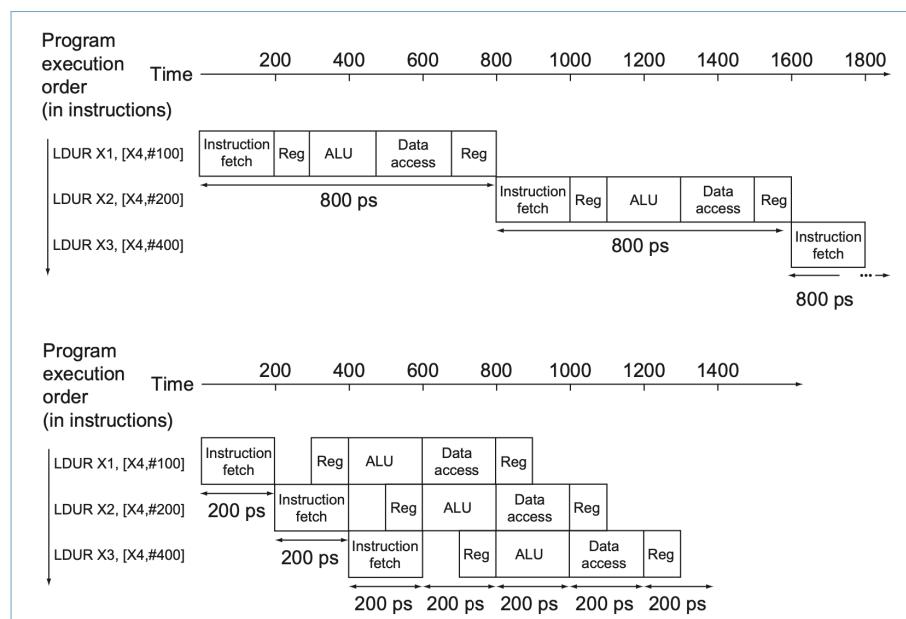


Figure 33: *Single-cycle execution vs Pipelined execution, showing an improvement from 800ps to 200ps*

Pipelining improves performance by increasing instruction *throughput* instead of decreasing execution times by instruction.

### Pipelining the Datapath

The separation of the five stages can be seen in the datapath below.

The left to right flow has some exceptions. Write backs involve putting a result in the middle of the path. The selection of the branch that involves +4 or jumping.

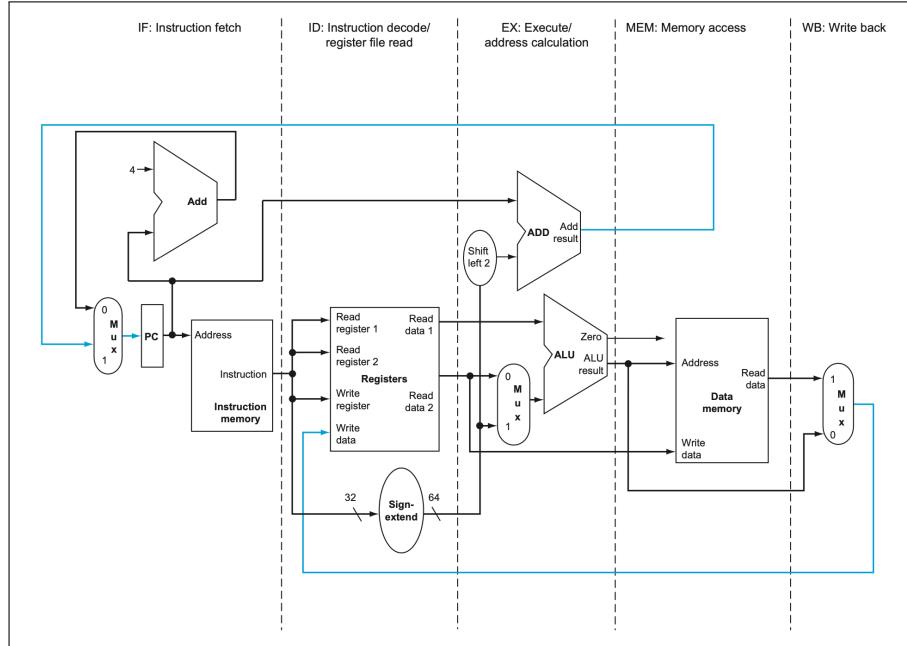


Figure 34: Five stages separation in Datapath

Pipeline registers separate the five stages and prevent data loss, keeping stability for newer instructions to use.

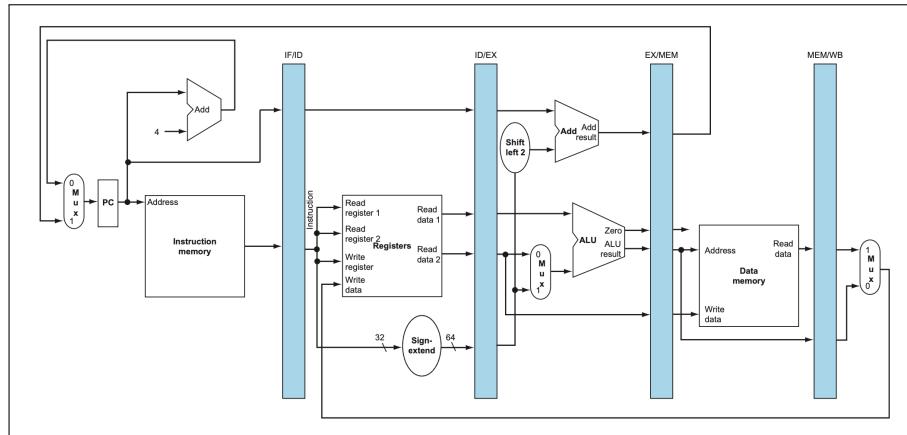


Figure 35: Pipeline registers in Datapath: IF/ID, ID/EX, EX/MEM, MEM/WB

Each logical component (Inst. Mem, RegFile, ALU, etc) must be used only once in a single pipeline stage to avoid *structural hazards*.

As instructions, like loads, need to keep track of a write back register,  $Rt$ , which means that this information needs to be propagated all the way upto the MEM/WB stage.

Propagation also happens with the PC, for calculating branching.

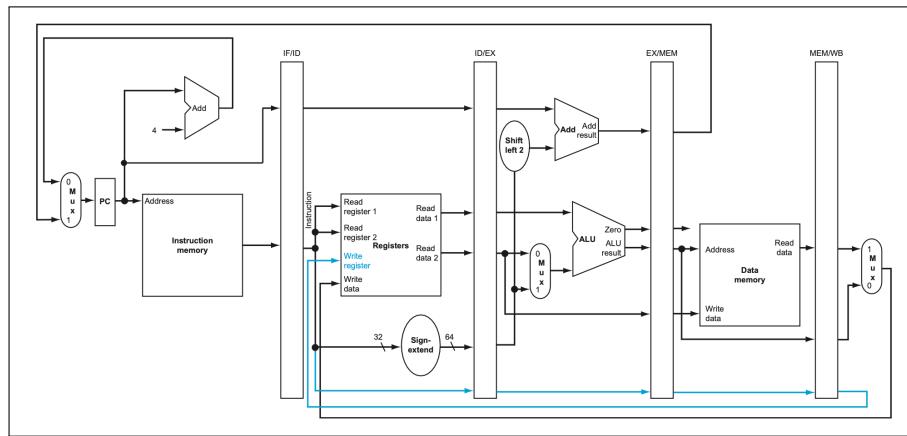


Figure 36: Pipelined Datapath with propagates write back register

## Pipelining the Control

As each control line is associated with one component active in one pipeline stage, the five stages can be used to classify the existing control signals. This means that pipelined registers can be expanded to fit the control signals.

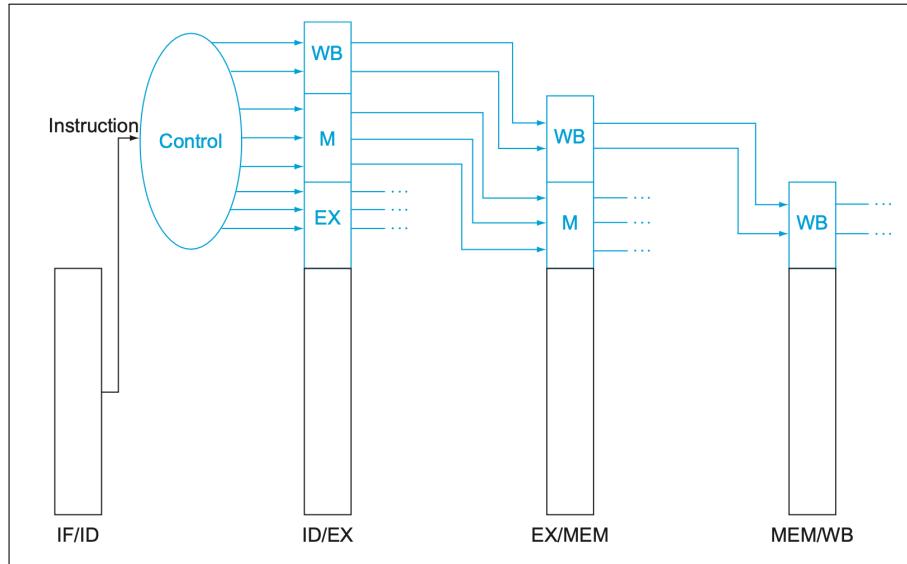


Figure 37: Control signals with applied Pipelining

Considering that the control signals start at EX, the control information can be created at ID.

The Reg2Loc signal is directly encoded into the instruction field as a bit that determined whether to use the R<sub>m</sub> for R-type or R<sub>t</sub> for cond. branching and load/store instructions.

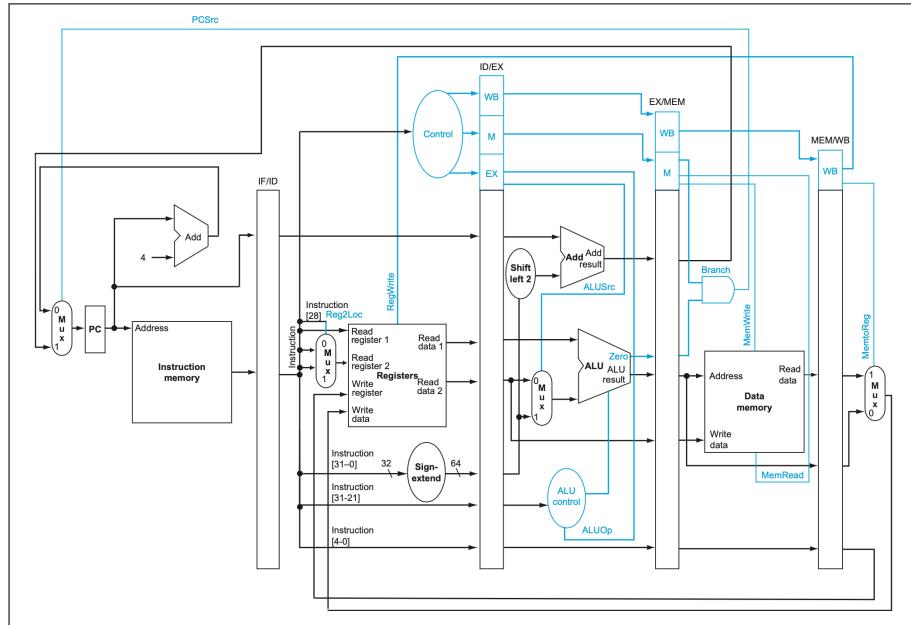


Figure 38: *Pipelined Datapath with Pipelined Control Signals*

# Pipeline Hazards

→ Structural Hazard

When one element cannot support two instructions at the same time.

Imagined if there was no distinction between IM and Data Memory, which would cause a first instruction to access data, while a fourth one is trying to fetch it at the same time.

## → Data Hazards

When a new instruction needs data from another instruction that has not yet finished executing.

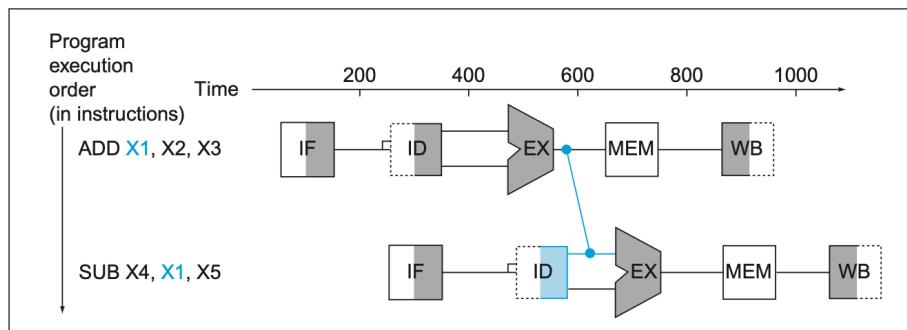
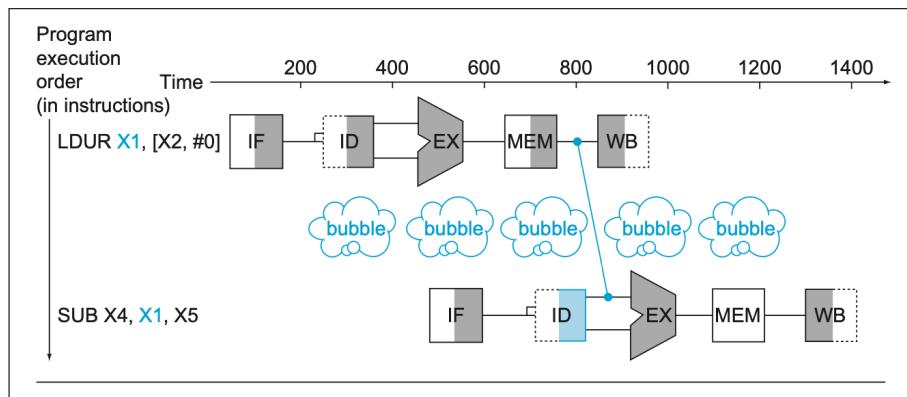
Data dependence from one instruction to an earlier one.

ADD X19, X0, X1  
SUB X2, X19, X3 // X19 depends on ADD instruction's execution

Stalling would cause a waste of 3 clock cycles, and compilers' optimizations are not enough.

A solution is to **forward** or bypass, which entails, in the example, to supply the result of the ALU directly into the input of the subtraction (or RegFile write).

There are cases where pipeline stalls or bubbles are present in forwarding, such as the **load-use data hazard**, in which the data has not yet been loaded when they are needed by another instruction.

Figure 39: *forwarding implementation in ADD and SUB example*Figure 40: *Load-use hazard with one cycle stalling to forward data*

**True data dependences** are when results are used later as operands.

**False data dependences** are when results are not passed through yet are clobbered, such as reading a register and then writing to the same register (write clobbers the read, but can be solved by using another register name).

### Types of False Data Dependences

**Anti-dependence:** Reading a register and then writing (clobbering) to the same register (WAR).

**Output dependence:** Writing to a register and then writing to the same register (WAW).

ADD	R1, R2, R4	; $R1 \leftarrow R2 + R4$
SUB	R4, R1, R5	; $R4 \leftarrow R1 - R5$
AND	R6, R1, R6	; $R6 \leftarrow R1 \text{ and } R6$
OR	R8, R1, R9	; $R8 \leftarrow R1 \text{ or } R9$
XOR	R6, R1, R11	; $R6 \leftarrow R1 \text{ xor } R11$

Figure 41: Green are True dependences, Red are Anti-dependences, and Blue are Output-dependences

For stalling you can use NOP instructions that have no effect other than acting as bubbles in the pipeline, in order to conserve the forwarding.

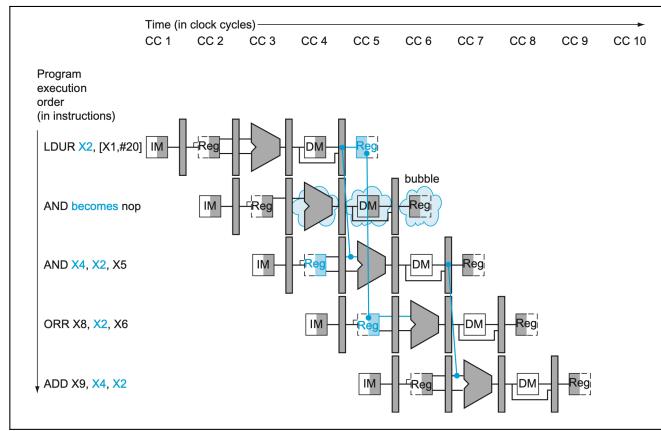


Figure 42: Example of NOP use for stalling and later forwarding

### → Control Hazards

When the fetched instruction is not the one needed (e.g. encountered branching).

Prediction is used to assume rather than wait for the branch to ascertain the outcome, such as always not taking the branch.

Dynamic prediction differs in result depending of what happened before (keeping a history).

The penalty for an incorrect prediction is that the pipeline must restart from the proper branch address, stalling the process.

**Simple Branch Prediction Methodologies**

**Branch not taken:** Never takes a branch, and proceeds full speed when branch is not resolved (*branch taken* is the opposite).

**1-bit Dynamic branch:** Keeps a 1-bit history of whether the branch was taken or not, and replicates latest behavior.

Another two ways to predict are using 2-bit prediction, and tournament schemes (multiple predictions for a branch and they "compete").

With 2-bit predictions, there are 4 states of decisions: strongly taken, weakly taken, strongly not taken, weakly not taken. If "strongly", then if it is the opposite it will become "weakly", and if it is "weakly" then if it is the opposite it will become the other "weakly" (like two-step verification).

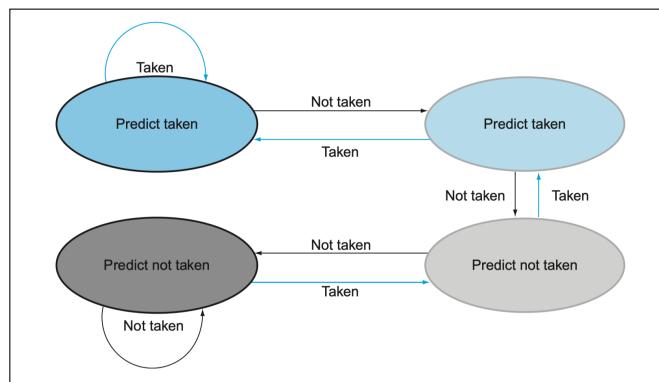


Figure 43: 2-bit Branch Prediction Logic

## Chapter 5

The memory hierarchy has elements increasing in sizes and access times the further away from the processor they are. Derived from the locality principle.

**Principle of Locality**

**Temporal:** If an element is referenced, then it will tend to be referenced again very soon.

**Spacial:** If an element is referenced, then elements around it will tend to be referenced again very soon.

*Elements in this case refers to data locations.*

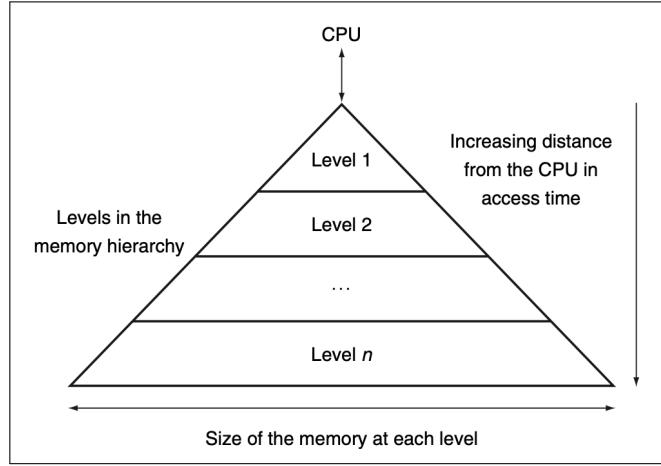


Figure 44: 2-bit Branch Prediction Logic

## SRAM

Volatile memory arrays with a read/write port that have a fixed access time to any data.

They don't need refreshing (access time is close to cycle time) and they require minimal power to retain charge in standby.

*Used to be used as caches, but nowadays caches are integrated onto the processor chip.*

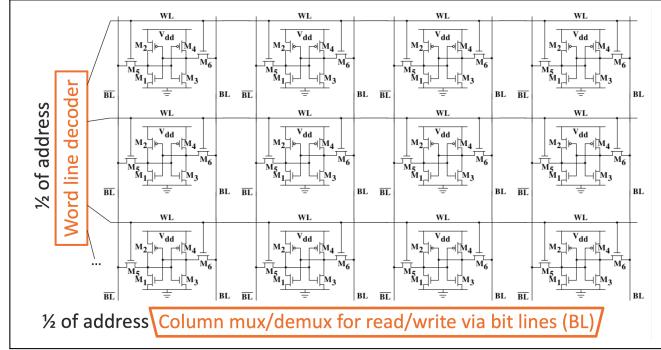


Figure 45: SRAM array bit addressing

For selection, there is a row (word-line MUX) and column decoder (column MUX/DEMUX).

→ Row: Takes the higher order  $k$  bits of the address and allows selection of  $2^k$  rows.

→ Column: Takes the remaining lower order  $m$  bits and allows selection of  $2^m$  columns.

## DRAM

Called *dynamic* because it requires periodical refreshing (read out and write in), in contrast to SRAM (static).

DRAM is volatile.

To save time, DRAM (made up of banks that can be accessed in parallel) refreshes entire rows at a time (two-level decoding) through fast buffers (SRAM-like).

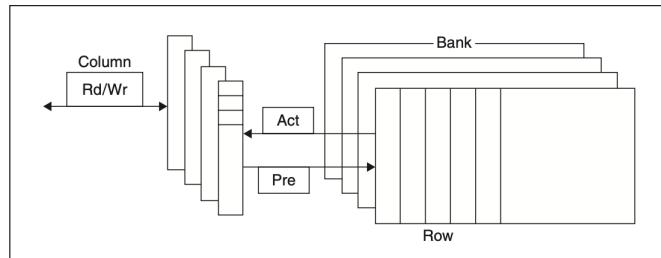


Figure 46: *DRAM Banks and Row Buffering*

Once in a row buffer, you can access all data of the row by changing the column address at no major cost.

**SDRAMs** share clock signal (synchronous), which has faster and more predictable access times.  
*No need for time negotiation overhead.*

## Flash Memory

Flash is non-volatile, RAM, and tough, particularly good for mobile devices.

Writes can wear out flash memory bits, thus it employs wear-leveling: Remapping used blocks to less used blocks.

## Disk Memory

Non-volatile memory which surface is made of tracks (circles), read/written to through a small coil called a *read-write head*.

The smallest unit of information read/written is a sector.

*Principle of locality helps reduce avg. seek time (how much it takes in finding the correct track)*

## Basics of Caches

**Cache** refers to storage that takes advantage of locality of access. It can be understood as being guided by the following questions.

- Is the referenced block already in the cache?
- If so, where in the cache can I find it?

## Direct Mapping Cache

Each memory location is mapped to one location in the cache.

$$(\text{Block Address}) \bmod (\text{Number of Blocks in the Cache})$$

Cache entries can be computed with the lower order k-bits, called the **index** (e.g. if number of blocks is 8, then you need  $\log_2 8 = 3$  bits).

**Tags** are additional identification to make sure that the word requested matches with the one stored in the cache. They are the upper bits remaining after indexing, stored in the cached entry, that gets compared with the request upon access (hit or miss).

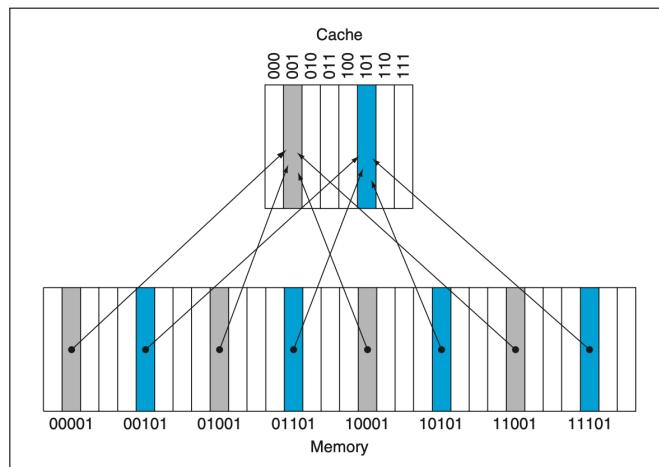


Figure 47: Direct Mapping Cache Block Addressing illustration

There is also a **valid bit** to distinguish whether the cache entry has valid data (e.g. upon processor start up), and a **dirty bit** that indicates writes to a cache line after it was brought in (for things like updating the reference in main memory).

Caches benefit from *temporal locality* as referenced words replace less than recently referenced ones. Note that the cache address starts with a byte offset (needed later to index into sets and pick which block is needed).

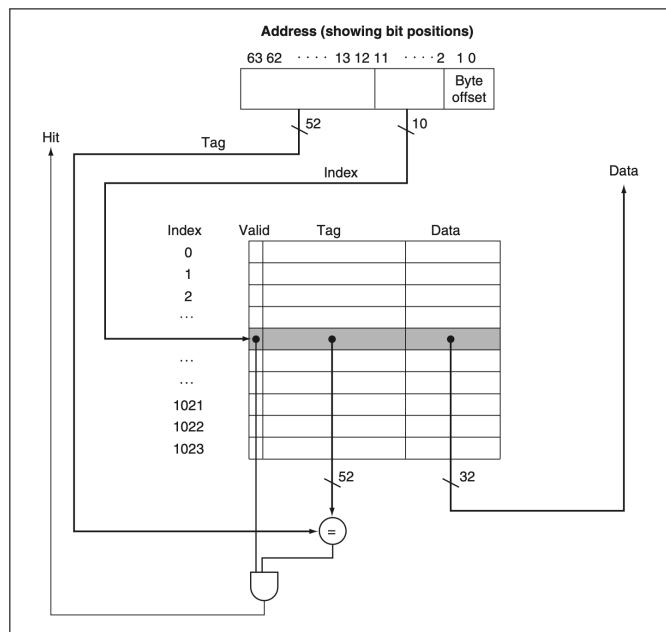


Figure 48: Diagram of Direct Mapping Cache and the relationship with the Address

The tag field will have the following size for 64-bit addresses, where there are  $n$  index bits ( $2^n$  blocks),  $m$  bits to address the  $2^m$  words, and  $m + 2$  for addressing one of the 4 bytes per word.

$$64 - (n + m + 2)$$

Thus the total number of bits in a directly-mapped cache is

$$2^n \times (\text{BlockSize} + \text{TagSize} + \text{ValidFieldSize})$$

$$2^n \times ((2^m \times 2^2 \text{ bytes} \times 2^3 \text{ bits}) + (64 - n - m - 2) + 1)$$

$$2^n \times ((2^m \times 32) + (64 - n - m - 2) + 1)$$

Larger blocks can exploit spatial locality to lower miss rates. When the improvement slows down, the miss penalty overwhelms the lower rates.

### Cache Write Handling

**Write-through:** Writes update both the cache and the next lower level in the memory hierarchy using a **write buffer** to avoid stalls due to overhead in writing to lower memory.

**Write-back:** Replaces the written value in the cache to lower memories when the block is replaced.

## Cache Misses

### The 3 C's

**Compulsory miss:** Unavoidable miss that occurs when first loading a block into the cache (e.g. on processor startup).

**Capacity miss:** Cache is too small to hold the required block.

**Conflict miss:** Competing for the same memory address due to cache design (even if the cache is large enough).

Compulsory misses can be improved by increasing the size per block in the cache (less blocks used by the program) or load through prediction.

Solutions to conflict misses is to add more slots/sets to a particular index (k-associative design).

To illustrate how performance is affected through hits and misses, AMAT is used.

$$AMAT = TimeForHit + MissRate \times MissPenalty$$

## Set & Full Associative Cache

For **set-associative**, there are n-sets (where  $n \geq 1$ ) that are mapped by the index, that contain k number of blocks that will perform tag comparisons against the requested block's tag.

**Full associative** does not perform a direct mapping by the index, but rather requires a search, as a block can be placed in any location of the cache (equivalent of using 1 set that holds all block locations).

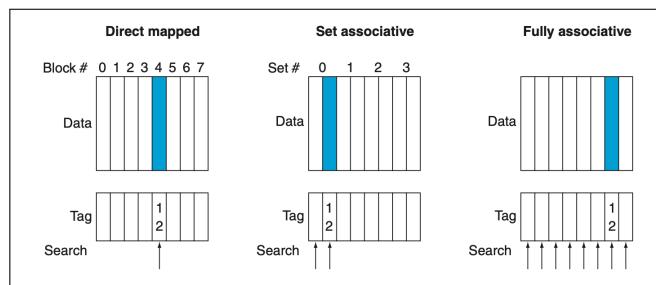


Figure 49: Example of Memory access with the three different Cache designs

A more detailed look into the memory hierarchy.

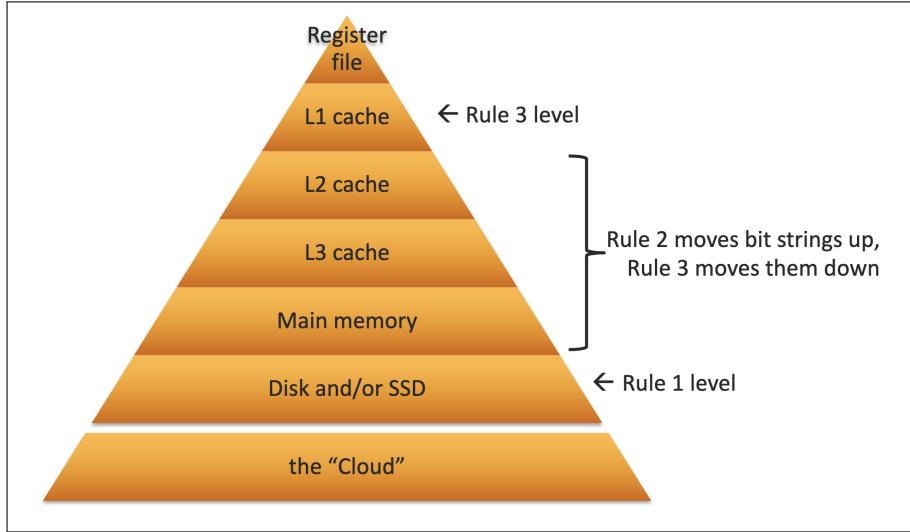


Figure 50: Example of Memory access with the three different Cache designs

## Virtual Memory

Main memory has limited space, so **virtual memory** treats it like a "cache relative to secondary storage" (handles managing the two memories).

To enforce protection between different programs, virtual memory allows all programs to have their own **address space**, which addresses will get translated into *physical addresses* (relocation).

If it is not mapped to main memory, then it gets mapped to secondary regions (e.g. disk).

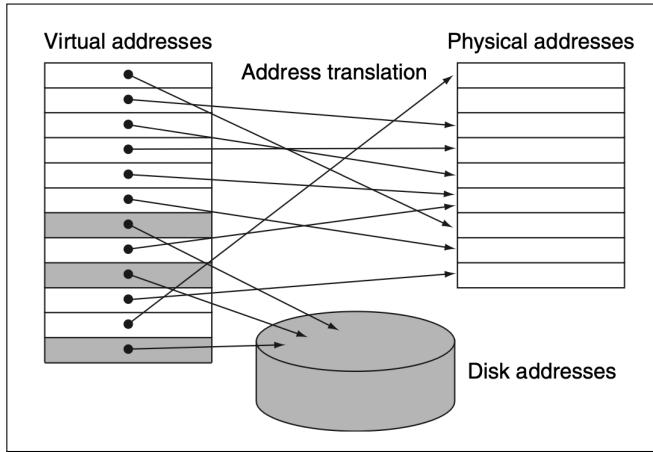


Figure 51: Pages mapped from Virtual addresses to Physical addresses

When a virtual memory block requested is not in main memory there is a **page fault**.

Virtual memory addresses are made of the *virtual page number* (high-order bits) and a *page offset* (low-order bits).

- The page offset remains unchanged when translating from virtual to physical addresses and is used to obtain information directly from page.
- The virtual page number (VPN) has more bits, and thus maps to more locations than the physical page number (main memory vs secondary memory). Illusion of unbounded memory.

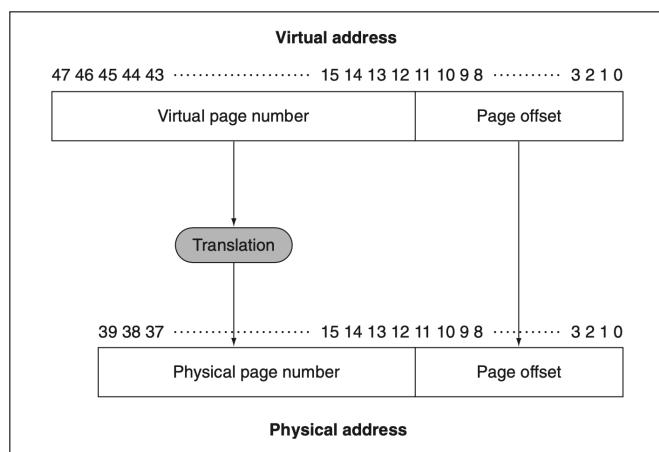


Figure 52: *Virtual address translation to Physical addresses*

Physical memory is stored in three parts: OS, page table, page frames (largest).

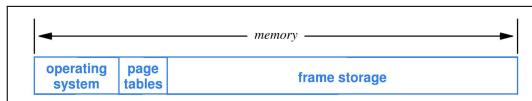


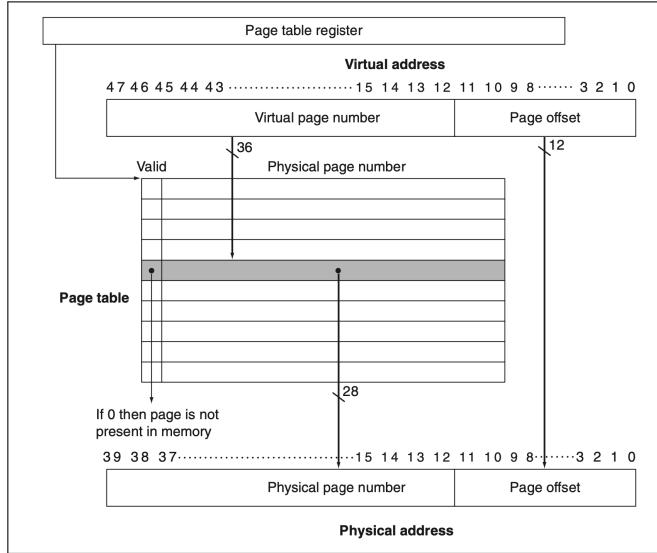
Figure 53: *Virtual address translation to Physical addresses*

## Page Tables

To find the corresponding physical page number (PPN) from a VPN, we use **page tables**, seen in the Figure 54.

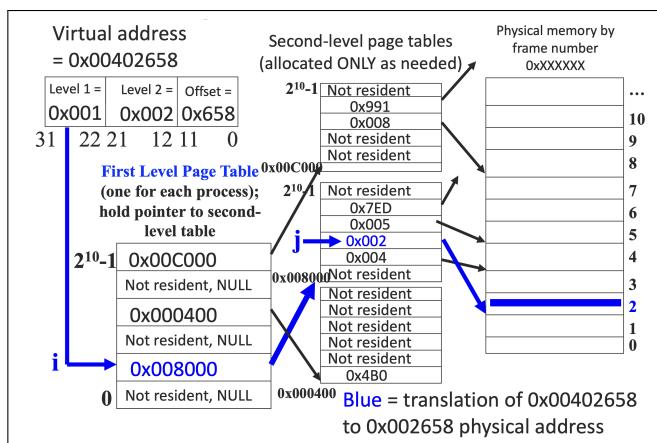
### Step by Step on Page Table Indexing

- 1. VPN gets aligned with Page Table:** VPN gets multiplied by the size of a PTE (page table entry), and then gets added to the base of the page table stored in the Page Table register (like an array entry).
- 2. Mapping and access of PPN:** There can be either a NOT PRESENT (causing a page fault) or a PPN (valid bit) indexed at the VPN offset calculated.

Figure 54: *Page Table illustration with a valid bit to determine entry*

The problem with design is that it will generate a massive page table (too big to fit in RAM), so the solution is to add multi-level page tables.

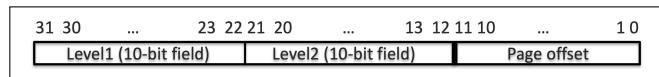
- First level page table (page-directory) is allocated always, and will point to second level page tables.
- Second level page tables are allocated as they are used, and point to another level or the physical memory address needed.

Figure 55: *Multi-level Page Table illustration*

A way to speed up this process is using the TLB (translation lookaside buffer) that acts like a full-associative cache into the page-directory.

Fully-associative as conflict misses will turn into capacity misses (capacity required is low).

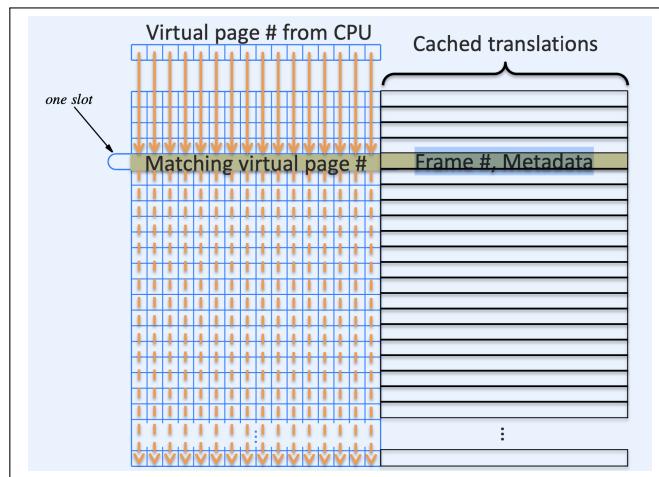
For improving TLB lookups we can use the CAM (content addressable memory) that, in constant

Figure 56: *Multi-level Page Table address field*

time, gets a key and does parallel comparisons will all CAM entries.

→ Upon hit it gives the PPN and metadata (e.g. SIGSEGV detection), and upon miss it uses the page-directory.

CAM is expensive and is organized as N m-bit slots, where N and m are chosen by need.

Figure 57: *CAM structure for constant time parallel access*

## Page Faults

The additional mappings by the VPN, that are not to main memory, map to the **swap space** (on disk or flash).

When a page fault occurs the OS will add the *least recently used* page (prediction).

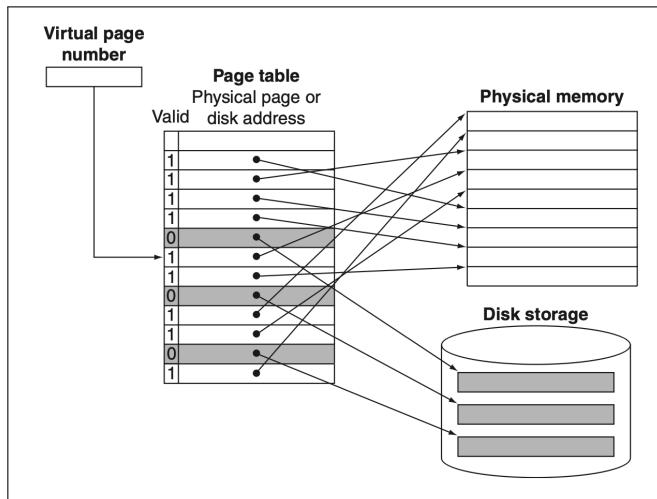


Figure 58: *Page Table that maps to main memory and swap space*

The MMU (memory management unit) makes up the page table register and the TLB, and will send the page fault signal (upon miss) to the OS, who then calls the page fault handler.

Handling page faults takes a lot of time.

## Bus Errors

Following of what was said about buses before, there are potential errors that can arise.

### Bus errors

**Address conflicts:** Two different elements/devices respond to the same address bus.

**Unassigned address:** No devices respond to a given address.

## Week 12 Slides

Interrupts are built in hardware, that saves the current pointer, jumps to code for the interrupt, and resumes execution from when *return from interrupt* is reached.

### Device Drivers

The **device driver** is software that communicates with a device: initializes the device, CSRs (control & status registers) are set to start I/O ops, and handles device interrupts.

#### Purpose of Device Drivers

**Encapsulation:** Creates an interface for wrapping up the details of how to talk to a device.

**Hiding:** Hides what happens under the hood.

**Independence:** No dependence on the design of any device (generic).

The structure of a device driver: lower half, upper half, and shared variables.

- Lower half: Handler code invoked during interrupt (direct communication to device).
- Upper half: Functions invoked by applications (I/O requests).
- Shared variables: Used by both halves with input and output buffers.

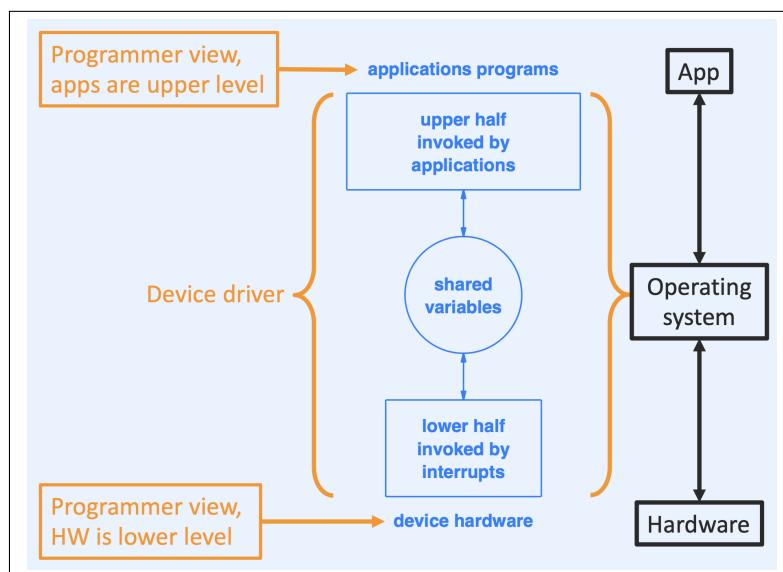


Figure 59: *Device driver structure*

### Types of Devices

**Character oriented:** Transfer one byte at a time (e.g. keyboards).

**Block oriented:** Transfer blocks of data at a time (e.g. disks, network interfaces).

## Output Request Queue

There is a queue shared by both halves (shared variables), where the top enqueues (when the application issues I/O) and the bottom dequeues the next request (starting it).