

CENG 311
Computer Architecture Lecture Notes

Ege Özkan

November 12, 2020

Chapter 1

Introduction - October 16, 2020

1.1 Four Key Current Directions

1. Fundementally secure/reliable/safe architectures
2. Fundementally energy-efficent and memory centric architectures
3. Fundementally low latency and predictable architectures
4. Architectures for AI/ML, Genomics, Medicine, etc.

1.2 Transformation Hierarchy

The order travels through different hierarchical levels until it reaches the electrons. From problems to algorithms, through program/language to system software to SW/HW Interface to lower hardware components.

Computer architecture was traditionally limited to SW/HW Interface and to Micro-architecture, but in the present day, computer architecture expands from algorithms to devices. This is because, to achieve the highest energy efficiency and performance, one must take the expanded view therefore co-designing across the hierarchy.

This way, once can specialize most of the components for a specific domain.

1.3 Computer Architecture

Computer architecture is the science and art of designin computing platforms to achibe a set of design goals. Designing a supercomputer is different from designing a smartphone, but many fundamental principles are similiar.

The computer architecture allows better systems to be built by making computer faster, cheaper, smaller and more reliable, it enables new applications and enables better solutions to be found.

Studying computer architecutre allows one to understand why computers work the way they do.

1.3.1 Computer Architecture Today

The present day industry has entered a paradigm shift to novel architectures, as many difficult problems motivate and cause a demand for novel architectures.

Chapter 2

Computer Abstractions - October 23, 2020

Computer architecture is the science and art of selecting and interconnecting hardware components to create a computer that meets functional, performance and cost goals. Computer systems come in many forms, from general purpose personal computers to supercomputers.

A system must be *Functional* (correct), *Reliable* (continue to perform correctly.), High performance, low cost, low power/energy consumption and it must be secure. Keep in mind that the word *correct* may mean different things in different accuracy levels.

2.1 Abstractions

The computer system structure consists of the Application software at the top, the system software in the middle and the hardware in the bottom. **Hardware-software interface** handles the translation.

The compiler takes a high level compiled language (such as C), converts it into an Assembly language, and the assembler takes the Assembler code and converts it into the machine code.

A Microarchitecture is an implementation of an ISA, there may be multiple implementations from the same ISA.

Levels of transformation create abstractions. High-level language programmer does not really need to know what the ISA is and how a computer executes instructions. Abstractions also improve productivity, decisions made in the underlying level does not need to be considered.

Knowledge of the lower level may improve higher level design choices of a person. [For instance, consider the fact that modulo operator is *slower* in the hardware level, when writing C code, modulo operations can be transformed into bit-shift operations, which are *faster*, improving speed.]

2.2 the Von Neumann Architecture

the Von Neumann Architecture consists of a main memory, a CPU and the interconnection between them. Within the CPU, there is a control unit, and an arithmetic/logic unit. the Von Neumann Architecture is the *traditional*

structure of computers.

2.3 Main Memory

Collection of locations, each of which is capable of storing both instructions and data. Every location consists of an address, which is used to access the location, and the contents of the location. It is similar in structure to a programming array. [Interesting to note, many emulators actually use arrays to emulate the memory of simpler machines.]

2.4 Central Processing Unit

CPU consists of multiple parts, **control unit**, **arithmetic logic unit** and **registers**. Memory is *fetches/read* to the CPU, as CPU sizes tend to be significantly smaller than Memory sizes, that is why despite CPU being much faster, memory is used to read/fetch from and *write/store* to.

2.4.1 Program Compilation as Execution

C is a compiled language, which means C code such as `float value = x[i];` would get translated to assembly instructions such as `ld r0, addr[1]`. [Many C compilers also optimize your code, so there are changes made to it. (Some even causing bugs.)]

The resulting assembly code is stored in the memory as binary machine code. These instructions are fetched from the memory to the CPU (ALU). A specialised register called **program counter** (PC), *points* to the instruction to be executed. When ALU executes an instruction, program counter is incremented.

2.4.2 Metrics

Performance is generally valued in the CPU architecture, this is not as straightforward as one assumes, however. Consider an airplane, airplane performance can be defined as passenger capacity, cruising speed, cruising range or passengers per mph. As such, there are different *metrics* to optimize for.

For computers, **Latency** is the elapsed time to do a task, how long it takes to do a task while **Throughput** is total work done per unit time. Although they may *look* similar, they are not the same.

In this course, while discussing processor performance, the focus will primarily be on the execution time for a single job (latency). It can be measured in different ways, **Execution time** includes all aspects, the total response time while **CPU time** is the direct time spent on processing a given job, this discounts IO time and other jobs' shares.

$$\text{Performance} = \frac{1}{\text{Execution Time}} \quad (2.1)$$

Relative performance is the performance comparison between two computers. Saying X is n times faster than Y is the same as dividing their performances or the inverse of their execution times.

If Computer A runs a program in 10 seconds, and B in 15 seconds, A is said to be 1.5 times faster than B.

2.4.3 CPU Clock

CPU clock determines when events take place in the hardware. [a clock cycle is similar to a Minecraft Tick.] a clock period is the duration of a clock cycle. Clock period has the unit of time and Clock frequency, which is the cycles

per second, has the unit of Hz.

This has the effect that the CPU Time is the:

$$\text{CPU Time} = \text{Cycle Count} \times \text{Cycle Time} \quad (2.2)$$

Where Cycle Count count of cycles needed for program to run, CPU Time is the CPU time of the program, Cycle Time is the time it takes for a single CPU cycle.

$$\text{CPU Time} = \frac{\text{Cycle Count}}{\text{Clock Rate}} \quad (2.3)$$

Where Clock Rate is the clock rate. Therefore, increasing the clock rate, or decreasing the cycle time will improve performance.

For instance, if a program runs in 10 seconds on a computer A, which has a 2GHz clock rate, which means the clock cycles needed for this program in computer A is $2 \times 10^9 \text{Hz} \times 10\text{s}$ Clock cycles.

If computer B runs this program in 6 seconds, but 1.2x more clock cycles, to calculate its clock rate, one just has to multiply the Clock rate above with $\frac{1.2}{6\text{s}}$, which results in the calculation Clock Rate = 4GHz

Different instructions take different amounts of time on different machines, division generally takes more time than addition, floating point operations take longer than integer ones, accessing memory takes more time than accessing registers. So one can assume that the number of cycles will equal the number of instructions for simplicity, but it would be incorrect.

Instruction count for a program is determined by a program, ISA and compiler, but the average cycles per instruction is

determined by CPU hardware, if different instructions have different CPI, average CPI is affected by instruction mix.

Clock Cycles equal Instruction count \times Cycles per instruction, whereas the CPU time depends on the product of Instruction count, CPI and Clock Cycle time, keep in mind Clock Cycle Time can be changed with $\frac{1}{\text{Clock Rate}}$.

Imagine an example where:

Computer	Cycle Time	CPI
A	250ps	2.0
B	500ps	1.2

Here, $t_a = \text{CC} \times \text{CCT} = I \times 2 \times 250$ and $t_b = I \times 1.2 \times 500$ where I is the instruction count, comparing them, $\frac{t_a}{t_b} = \frac{500I}{600I}$, shows that computer a is $\frac{6}{5}$ times faster for a program.

Another example, this time with compilers:

	A	B	C
CPI for class	1	2	3
IC in sequence 1	2	1	2
IC in sequence 2	4	1	1

Here the compiler writer can choose between different code generation sequences for three classes A, B and C. But which code sequence is faster, and what is the CPI for each sequence?

For the first sequence, there is a total of $1 \times 2 + 2 \times 1 + 3 \times 2 = 10$ Clock Cycles spent, for the second sequence, there is a total of $1 \times 4 + 2 \times 1 + 3 \times 1 = 9$. Since the compiler will execute them in the same computer, the clock rate is equal, and hence, the time being spent is directly related to the number of clock cycles, therefore, the second sequence is much more beneficial.

One could also calculate the average CPI, this is calculated using:

$$CPI_A = \frac{\text{Clock Cycles}}{\text{Instruction Count}} \quad (2.4)$$

$$= \sum_{i=1}^n \left(CPI_i \times \frac{\text{Instruction Count}_i}{\text{Instruction Count}} \right) \quad (2.5)$$

$$\text{Clock Cycles} = \sum_{i=1}^n (CPI_i \times \text{Instruction Count}_i) \quad (2.6)$$

Where CPI_A is the weighted average CPI.
Overall:

$$\text{CPU Time} = \frac{\text{Instructions}}{\text{Program}} \times \frac{\text{Clock cycles}}{\text{Instruction}} \times \frac{\text{Seconds}}{\text{Clock Cycle}} \quad (2.7)$$

Determining the Values

CPU Execution Time can be determined by running the program.

Clock cycle time is usually published with documentation.

Instruction Count via the software tools that profile execution or by simulators.

CPI however, depends on a wide variety of design details including the memory system and the processor structure, therefore it is much harder to determine.

Computer Benchmarks are programs or set of programs used to evaluate computer performance, benchmarks allow us to make performance comparisons based on execution times, they can vary greatly in terms of their complexity and their usefulness, benchmarks should:

- Be representative of the type of applications that run on the computer.

- Not be overly dependent on one or two features of a computer.

Amdahl's Law

Improving an aspect of a computer and expecting a proportional improvement in overall performance. Where n is the improvement factor:

$$T_{\text{improved}} = \frac{T_{\text{affected}}}{n} + T_{\text{unaffected}} \quad (2.8)$$

For instance if multiply accounts for 80 seconds of 100, how much improvement in multiply performance to get 2 times overall improvement. [ie to make 100 seconds to fall to 50 seconds]

$$50 = \frac{80}{n} + 20 \quad (2.9)$$

Must improve by 2.6 times.

However if multiply accounts for 80 seconds of 100 seconds and we want five times overall improvement. This is impossible, since the unaffected part also takes 20 seconds.

Summary of Performance Evaluation

Good benchmarks, such as SPEC can provide an accurate method for evaluating and comparing computer performance. Amdahl's law provides an efficient method for determining speedup due to an enhancement, and one must make the common case fast.

Chapter 3

Instructions: Language of the Computer - October 30, 2020

It was already established that the Compilers and Assemblers bring high-level languages to binary machine language program. The Instruction Set Architecture (ISA) consists of Instructions and Memory.

ISA's instructions are controlled by Op-codes, Addressing Modes, Data Types, Instruction Types and Formats, Registers and Condition codes. An ISA interacts with the memory's address space, it has to deal with its addressability, alignment and virtual memory management.

A microarchitecture is a specific implementation of an ISA under specific design constraints and goals. It is the things done in hardware without exposure to the software. These range from pipelining to voltage/frequency scaling.

3.1 MIPS Instruction Set

MIPS is a simple ISA.

3.1.1 Assembler Instructions

Add and Subtract

Takes two sources and has one destination.

```
add a, b, c
```

Where **a** gets **b + c**, arithmetic operations occur **only on registers** in MIPS. **sub**, **mult** and **div** take the same form.

This is an important lesson on design, as a design principle, *simplicity favours regularity*.

Consider the code

$$f = (g + h) - (i + j)$$

Compiles to the following MIPS Assembly

```
add t0, g, h
add t1, i, j
sub f, t0, t1
```

Register Operands

Observe that the arithmetic instructions use register operands. MIPS has a 32x32-bit register file, that is, 32 register with 32-bit register size. Used for frequently accessed data, numbered 0 to 31. In Assembly, they are named **\$t0** to **\$t9** for temporary values and **\$s0** to **\$s7** for saved variables.

Smaller is faster.

In the same C code, consider that all variables are stored in **\$s0** to **\$s4**.

(keep in mind that above MISP codes were not assemble-able).

```
add $t0, $s1, $s2
add $t1, $s3, $s4
sub $s0, $t0, $t1
```

Memory Operands

Main memory is used for composite data (arrays, structures, dynamic data) to apply arithmetic operations (loading and storing values), it is byte-addressed. Words are aligned in memory (its address must be a multiple of 4.) [This design decision is shared with x86, this is why sometimes stuff is padded in C in weird ways.]

A **word** is a 4-byte construct, (32-bits) that correspond to an integer. Do consider that big-endianness or little-endianness of an architecture will impact in which order a word, or an integer, will be stored in the memory in which order.

MIPS is a big-endian architecture.[Meanwhile x86 is a little-endian architecture]

```
lw d, off(b)
sw s, off(b)
```

d is the destination register, **b** is the base register **off** is offset value, **b + off** forms the memory address to be accessed. **lw** is the **load word** instruction, while **sw** is the **store word** instruction.

Consider the C code `g = h + A[8];`, **g** is in `$s1`, **h** in `$s2` and the base address of the **A** array is in `$s3`.

```
lw $t0, 32($s3)
add $s1, $s2, $t0
```

Observe that the $8 \times 4 = 32$ since $4i$ is the offset value for a specific index value in C. This

is due to the fact words occupy 4 bytes, and they align by four too.

Consider the similar C code `A[12] = h + A[8];` which also stores the result, this assembles into:

```
lw $t0, 32($s3)
add $s1, $s2, $t0
sw $t0, 48($s3)
```

!/!*/!*

Since registers are significantly faster to access than memory, compiler should try to use registers for variables as much as possible. This would be correct even if register access time and memory access time was the same, [Which they aren't], since memory access requires more instructions.

Immediate Operands

Immediate instructions negate the need for using load instructions by using scalars. Consider, **addi** which is the immediate addition instruction, compared to **add**, it does not need the usage of a load. For instance, this snippet:

```
li $t0, 32
add $s1, $s1, $t0
```

Can easily be rewritten as

```
add $s1, $s1, 32
```

Without using **li**. Which itself is an immediate operand that immediately loads a value to a register. This corresponds to the third design principle of MIPS, making the common case fast. Here, small constants are common, immediate operand avoids a load instruction.

Related to this, `$zero` read-only register (it cannot be overwritten) holds the value zero, adding with the zero register can be used to move values between registers for instance.

Logical Operations

Logical operations are bitwise operations used to perform boolean logic on values. These are **and**, **or**, **nor**, and their immediate equivalents **andi** and **ori**.

And is useful for masking bits. For instance, **anding** bits in a word with a word whose certain parts are zero would mask out these parts. Likewise, **or** is useful for swapping certain bits to 1 in a word.

MIPS does not have a not instruction, instead, one can **nor** (not or) with the special **\$zero** register. Since $\neg p = \neg(p \vee 0)$

Shift Operations

Another bitwise manipulation operations. **sll** and **srl** is used for left logical and right logical shifts respectively. Both of them fill the shifted parts with zero.

Shifting a number left i times is same as multiplying it by 2^i , and shifting a number i times right would mean dividing it by 2^i .

3.1.2 Instruction Representation

Instructions are kept as a series of high and low electronic signals. MIPS instructions are 32 bits long.

MIPS Instructions are classified into different types of instructions.

R-Type Instructions

Most well-known of these R-Types are arithmetic instructions. Shift amount refers to the amount of shift that will be applied to the number (For **sll** and **srl**.) Function code is used to further distinguish between

Part	Length	Explanation
op	6 Bits	OP-Code
rs	5 Bits	Source register
rt	5 Bits	Source register
rd	5 Bits	Destination register
shamt	5 Bits	Shift Amount
funct	6 Bits	Function Code

Table 3.1: The structure of R-Type Instructions

instructions as an extension of the OP Code.

Register numbers in the machine codes differ from their assembler counterparts a bit. Registers between **\$t0** and **\$t7** are translated to numbers 8 to 15, Registers **\$t8** and **\$t9** gets translated to 24 and 25 and registers between **\$s0** and **\$s7** are translated to 16 to 23.

[I should probably point out the instructor at this point said we should now the general instruction structure but not the register translations.]

Function field may look redundant, it exists due to the fact that all arithmetic operations are encoded with OP-Code 0, hence they are differentiated via their Function codes. [This whole deal occurs to keep the OP-Code length the same between R, I and J type operations. There are many more R type operations, instead of giving OP-Code more space OP-Code field is used to group them to certain categories, and Function field specifies the exact instruction. This also greatly simplifies the processor design, all instructions with OP-Code 0 gets sent to ALU, for instance.]

[You may point out, but Amber, doesn't all of the R type instructions go to ALU anyway? What is the point? Well, padawan, MIPS is designed to be extended, Sony PSP for instance has an VFPU in it as well, and R Type instruc-

tions on VFPU use a different OP-Code]

Assembler instructions like `add`, `sub`, `div`, `srl` and `sll` are translated to R-Type instructions.

I-Type Instructions

Part	Length	Explanation
<code>op</code>	6 Bits	OP-Code
<code>rs</code>	5 Bits	Register 1
<code>rt</code>	5 Bits	Register 2
constant/address	16 Bits	16 Bit Number.

Table 3.2: The structure of I-Type Instructions

Instructions in the I-Type are immediate arithmetic instructions, load/store instructions and branch instructions.

Meaning of registers change from instruction to instruction. The address is used to offset one of the registers.

Assembler Instructions like `addi`, `lw`, `sw` are translated to I-Type instructions.

J-Type Instructions

Part	Length	Explanation
<code>op</code>	6 Bits	OP-Code
<code>pseudo-address</code>	26 Bits	Shortened address

Table 3.3: The structure of J-Type Instructions

[These were not mentioned, but I am still putting this here for the sake of completeness, although I believe we will return to these in the future.]

J-Type instructions are reserved for Jump instructions, hence they have a large amount of space left for an address. This value, stored

at the `pseudo-address` field is the shortened address of the destination, its two least significant and four most significant bits are removed and are assumed to be the same as the current instruction's address. (Wikibooks - MIPS Assembly)