Introduction

COE 501

Computer Architecture

Dr. Muhamed Mudawar

Computer Engineering Department King Fahd University of Petroleum and Minerals

Welcome to COE 501

- Professor: Muhamed F. Mudawar
- ✤ Office: Building 22, Room 410-2
- ✤ Office Phone: 4642
- Schedule and Office Hours:
 - ♦ <u>http://faculty.kfupm.edu.sa/coe/mudawar/schedule/</u>
- Course Web Page:
 - ♦ <u>http://faculty.kfupm.edu.sa/coe/mudawar/coe501/</u>
- ✤ Email:
 - ♦ mudawar@kfupm.edu.sa

Modules and Topics

- 1. Defining computer architecture, classes of computers, technology trends, performance, power, and cost.
- 2. Instruction set architectures, simple pipelines, hazards, data forwarding, control, branch prediction, exceptions.
- 3. Memory hierarchy, DRAM, cache organization, cache optimizations, virtual memory, cache performance.
- 4. Instruction-Level Parallelism, complex pipelining, multiple issue, VLIW approach, software approaches to ILP.
- 5. Data-level parallelism, SIMD, Loop-level parallelism, GPU.
- 6. Thread-level parallelism, multithreaded cores, multiprocessor architectures, shared and distributed memory models, cache coherence, and synchronization.

Grading Components



Assignments / Project 20%

Midterm Exam30%

Final Exam30%

Textbook

Computer Architecture

A Quantitative Approach

♦ Sixth Edition

♦ John Hennessy & David Patterson

♦ Morgan Kaufmann Publishers, 2019

John L. Hennessy | David A. Patterson

Sixth Edition

COMPUTER Architecture



What is Computer Architecture?

Original Definition:

The attributes of a [computing] system as seen by the programmer, i.e., the conceptual structure and functional behavior, as distinct from the organization of the data flows and controls, the logic design, and the physical implementation. (Amdahl, Blaaw, and Brooks, 1964)

Today, this is known as Instruction-Set Architecture

Abstraction Layers in Computing



Abstraction Layers in Computing

Application

Algorithms, Programming Languages

Compiler, Linker, Run-time Libraries

Operating System, Virtual Machines

Instruction Set Architecture (ISA)

Microarchitecture

Logic Gates, Circuits

Devices, Layout

Physics

Focus of this course + Impact on Software and Compiler

ISA vs. Microarchitecture

- Instruction Set Architecture (ISA)
 - ♦ Class of ISA: register-memory or register-register architectures
 - ♦ Programmer visible state (Register and Memory)
 - ♦ Addressing Modes: how memory addresses are computed
 - ♦ Data types and sizes for integer and floating-point operands
 - ♦ Instructions, encoding, and operation
 - ♦ Exception and Interrupt semantics
- Microarchitecture / Organization
 - ♦ Tradeoffs on how to implement the ISA for speed, energy, cost
 - Pipeline width and depth, cache organization, peak power, bus width, execution order, memory hierarchy, etc.

The Power of Abstraction

Abstraction:

- ♦ A higher level only needs to know about the interface to the lower level, not how the lower level is implemented
- Example: a high-level language programmer does not need to know what the ISA is and how a computer executes instructions

Abstraction improves productivity

- ♦ No need to worry about decisions made in underlying levels
- Example: programming in Java vs. C vs. assembly vs. binary vs. specifying control signals of each transistor every cycle
- Then, why would you want to know what goes on underneath or above?

Crossing the Abstraction Layers

- As long as everything goes well, not knowing what happens underneath or above is not a problem.
- What if
 - ♦ The program you wrote is running slow?
 - ♦ The program you wrote consumes too much energy?
- ✤ What if
 - ♦ The hardware you designed is too hard to program?
 - The hardware you designed is too slow because it does not provide the right primitives to the software?
- One goal of this course is to understand how a processor works underneath the software layer and how decisions made in hardware affect the programmer

Architecture Continually Changing



Compatibility: Ability of a new architecture to run older applications. Cost of software development makes compatibility a major force in market.

Changing Definition

- Computer Architecture's Changing Definition
- ✤ 1950s to 1960s:
 - ♦ Computer Architecture Course = Computer Arithmetic
- ✤ 1970s to mid 1980s:
 - ♦ Computer Architecture Course = Instruction Set Design,

Especially ISA appropriate for compilers

- ✤ 1990s until today:
 - ♦ Computer Architecture Course =

Design of CPU, memory system, I/O system, Multiprocessors

So, What is Computer Architecture?

Application

Algorithms, Programming Languages

Compiler, Linker, Run-time Libraries

Operating System, Virtual Machines

Instruction Set Architecture (ISA)

Microarchitecture

Logic Gates, Circuits

Devices, Layout

Physics

In its broadest definition, computer architecture is the design, organization, and implementation of a computing system that allows us to execute software applications efficiently using available manufacturing technologies, meeting price, power, and performance goals.

Computers Then ...



ENIAC : Electronic Numerical Integrator And Computer

Built by Eckert and Mauchly at the University of Pennsylvania (1943-45)

First general-purpose electronic computer.

Cost \$500,000 (\$6,300,000 today)

Around 18000 vacuum tubes, 7200 diodes, 1500 relays, 70000 resistors, 10000 capacitors, and around 5 million hand-soldered joints.

30 tons, 167 square meters, consumed 150 kw of power

Addition = 200 µs, division = 6 ms, read-in 120 cards per minute

Not very reliable (one vacuum tube fails about every two days)

Major Technology Inventions



Classes of Computers

- Personal Mobile Devices
 - ♦ Cell phones and tablet computers, battery-operated
 - ♦ Emphasis on energy efficiency, real-time performance, and cost
 - ♦ Use of flash memory storage for energy and size requirement
- Desktop Computers
 - ♦ Battery-operated laptop computers to high-end workstations
 - ♦ Emphasis on price-performance and graphics performance
- Servers
 - ♦ Backbone of large-scale enterprise computing
 - ♦ Emphasis on availability, scalability, and performance
 - ♦ Failure of a server is more catastrophic than a desktop computer

Classes of Computers - cont'd

- Clusters / Warehouse Scale Computers
 - ♦ Clusters are collections of servers connected by a network
 - ♦ Example: search engine, social networking, online shopping
 - ♦ Used for "Software as a Service (SaaS)"
 - ♦ Emphasis on availability and price-performance
 - Supercomputers are related but the emphasis is on floating-point performance and fast internal networks
- Embedded Computers / Internet of Things (IOT)
 - Found everywhere: printers, networking devices, autonomous cars, flying UAVs, game players, digital cameras, TVs, etc.
 - ♦ Widest spread of processing power and cost

Moore's Law



Cramming More Components onto Integrated

Circuits, Gordon Moore, Electronics, 1965

Number of transistors per integrated circuit

doubles every N months ($12 \le N \le 24$)



Moore's Law: Transistor Count on IC chips



Data source: Wikipedia (https://en.wikipedia.org/wiki/Transistor_count)

The data visualization is available at OurWorldinData.org. There you find more visualizations and research on this topic.

Licensed under CC-BY-SA by the author Max Roser.

Feature Size Scaling

Feature Size: minimum gate length of a transistor (source to drain)



Source: Weste and Harris, CMOS VLSI Design Lecture Slides

Technology Trends

- Integrated circuit technology
 - ♦ Transistor density: increases 35% per year (Moore's Law)
 - ♦ Die size: increases 10% to 20% per year (less predictable)
 - ♦ Integration overall: increases 40% to 55% per year
- DRAM capacity: 40% down to 25% per year (2000 to 2014)
 - Recently slowing down: 8 Gbits (2014), 16 Gbits (2019), 32 Gbits ???
- ✤ Flash storage capacity: increases 50% to 60% per year
 - \diamond 8X to 10X cheaper per bit than DRAM
- Magnetic disk capacity: 40% per year (2004 to 2011)
 - ♦ Recently slowing down to only 5% per year
 - ♦ 10X cheaper per bit than Flash
 - \diamond 100X cheaper per bit than DRAM

Growth of Capacity per DRAM Chip



Microprocessor Trend



Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten New plot and data collected for 2010-2017 by K. Rupp

Example of a Multicore Processor

IBM Power8 Processor

22 nm, 650 mm2 die 12 cores

16 exec pipe per core

Cache Hierarchy 32 KB Instruction Cache 64 KB Data Cache 512 KB L2 per core 96 MB eDRAM shared L3

Memory Interface Dual Memory controllers 230 GB/sec bandwidth



Processor Performance over 40 Years



The End of the Uniprocessor Era

- ✤ Late 1970's saw the emergence of the microprocessor
- New RISC architectures appeared in early 1980's
 - ♦ RISC = Reduced Instruction Set Computer
- RISC architectures raised performance using
 - ♦ Instruction Level Parallelism through pipelining & multiple-issue
 - ♦ Cache memory reduced memory latency
 - ♦ Instruction execution pipelines and caches improved clock rates
 - \diamond 17 years of improved performance at a rate of ~50% per year
- Single processor performance dropped in 2003
 - ♦ Lack of more Instruction Level Parallelism (ILP) to exploit efficiently
 - ♦ Power dissipation put a limit on higher clock frequencies

The End of Moore's Law

- If Moore's law was to continue until 2050, engineers have to build transistors that are smaller than a hydrogen atom!
- Moore's law might end in 2021 to 2030. However, transistors will keep getting better and more energy efficient.



Parallelism in Computer Architecture

Computer Architecture can exploit parallelism in four ways:

- 1. Instruction-Level Parallelism (ILP)
 - Through pipelining and multiple instruction issue each cycle
 Limited ILP in real programs, difficult to exploit efficiently
- 2. Data-Level Parallelism (DLP)
 - ♦ Single instruction operates on multiple data (SIMD)
 - ♦ Vector architectures and graphics processing units (GPUs)
- 3. Thread-Level Parallelism (TLP)
 - ♦ Can execute parallel threads on multiple cores (or same core)
 - ♦ Parallel threads communicate through shared memory
- 4. Request-Level Parallelism
 - \diamond Parallelism among independent processes specified by the OS

Conventional Wisdom in Computer Architecture

- Old Conventional Wisdom: Power is free, Transistors are expensive.
- New Conventional Wisdom: "Power wall" Power is expensive, Transistors are free (Can put more on chip than can afford to turn on)
- ✤ Old CW: Multiplication is slow, Memory access is fast
- New CW: "Memory wall" Memory is slow, multiplication is fast (200 clock cycles to DRAM memory, 4 clock cycles for multiplication)
- Old CW: Sufficiently increasing Instruction Level Parallelism via compilers, pipelining, out-of-order execution, speculation, …
- ✤ New CW: "ILP wall" law of diminishing returns on more HW for ILP
- ✤ Old CW: Uniprocessor performance 2X in 1.5 years
- New CW: Uniprocessor performance now 2X in ~7 years Change in chip design: switching to multiple "cores"
- Simpler processors are more energy efficient than complex ones

And in Conclusion ...

- Computer Architecture is more than ISAs and RTL
- It is about the interaction of hardware and software
- Design of appropriate abstraction layers
- Computer architecture is shaped by technology and applications
 - ♦ History provides lessons for the future
- Going from sequential to parallel computing
 - ♦ Requires innovation in many fields, including computer architecture