

# Early Adopter - Fundamentals of Computer Systems

---

Martha A. Kim  
Computer Science Department  
Columbia University

EduPar 2009  
May 16, 2011

# Computer Science @ Columbia University

---

## Columbia University

Four year, research-intensive university

Columbia College (CC) (~4k total enrollment)

School of Engineering and Applied Sciences (SEAS) (~1.5K total enrollment)

School of General Studies (GS) (~1.5K total enrollment)

## Computer Science Curriculum

Undergraduate majors from CC, SEAS, and GS

Two years of computer science core + two years in selected track)

### CS Core

- Intro programming
- Data structures and algorithms
- Complexity theory
- Fundamentals of computer systems
- Computational linear algebra
- Probability and statistics

### CS Tracks

- Foundations of CS  
(algorithms, computational complexity, scientific computing, and security)
- Systems  
(networking, programming languages, operating systems, and software systems)
- Artificial Intelligence  
(machine learning, robots, and systems capable of exhibiting "human-like" intelligence)
- Applications  
(interactive multimedia applications for the Internet and wireless networks)
- Vision and Graphics  
(vision, graphics, and advanced forms of human-computer interaction)



# Fundamentals of Computer Systems Overview

---

- Full title: CSEE 3827: Fundamentals of Computer Systems
- Required course, part of CS core (as well as EE and CE core)
- Students typically take this course in their second or third year
- Taught every semester
  - Fall enrollment typically ~40 students
  - Spring enrollment typically ~80 students
- Textbooks
  - Until Spring 2011: Logic and Computer Design Fundamentals, Mano & Kime + Computer Organization & Design, Patterson & Hennessy
  - Currently: Digital Design and Computer Architecture, Harris & Harris
- Website: <http://www.cs.columbia.edu/~martha/courses/3827/sp11>



# Course Syllabus - Part I: Digital Logic

Approximate lecture hours spent  
Bloom level

<b>Information Representation</b> BCD 2's Complement ASCII general coding	<b>1</b>	<b>A</b>	<b>Sequential Circuit Design</b> Latches Flip-Flops Timing Constraints Registers	<b>3</b>	<b>A</b>
<b>Boolean Logic and Algebra</b> Algebraic rules Standard forms (SOP, POS) Karnaugh Maps Timing	<b>5</b>	<b>A</b>	<b>Finite State Machine Design</b> Moore Mealy	<b>3</b>	<b>A</b>
<b>Combinational Circuit Design</b> Multiplexers Encoders and Decoders Arithmetic Circuits	<b>3</b>	<b>A</b>	<b>Storage Architectures</b> Base cell abstraction Array organization Register v. SRAM v. DRAM	<b>2</b>	<b>C</b>



# Course Syllabus - Part II: Machine Organization

<b>Instruction Set Architectures</b> Benefit RISC v. CISC	1	C	<b>Pipeline MIPS Processor</b> Datapath & control Data hazards Control hazards Performance analysis	3	A
<b>MIPS ISA</b> Arithmetic instructions Load & stores Branches and control Calling conventions Stack management Instruction formats	5	A	<b>Caches</b> Memory hierarchy Performance analysis Direct mapped cache N-way associative cache Fully associative cache LRU replacement	2	C
<b>Performance Analysis</b> Throughput v. latency CPI, etc. Speedup & Amdahl's Law	1	A	<b>Advanced Microarchitectures</b> CPI, and other components Speedup and Amdahl's Law	1	K
<b>Single Cycle MIPS Processor</b> Datapath Control Performance analysis	2	A	<b>Modern Architectures</b> Multicore (motivation) GPUs	1	K



# General Strategy for Incorporating Parallelism

---

- The syllabus already contained a number of parallelism related concepts
- To emphasize them for the students, made special effort to highlight parallelism and other recurring themes in CS as they arose over the semester, and to draw analogies to where students may already have seen them.

Theme	In this course	Seen before
Abstraction	Hierarchical circuit design	Functional decomposition of software
Interface v. implementation	ISA v. microarchitecture	Public v. private methods and fields in Java
Serial v. parallel	Combinational v. sequential circuits	?

- Finally, made the occasional passing reference to relevant advanced feature of current topic (e.g., the value of non-blocking cache design when teaching caches)



# Example 1: Circuits Compute in Parallel

- When teaching **combinational circuits**, heavy emphasis on parallel nature of their operation.
  - Logic gates are always computing, all at once, but not necessarily synchronized
  - Signals propagate constantly, but not in synchronized fashion
- **Observation:** students assume synchronous circuitry (likely from sequential programming training) before they've ever seen sequential synchronous circuits.

## Combinational Circuit Timing Analysis

A circuit has many paths from input to output

Signals are constantly propagating along these paths

Correct signals start propagating from correct input values

Some signals are faster than others (they race)

Computation done when last correct signal arrives at output

This is the critical path

### Ripple-Carry adder circuit depth

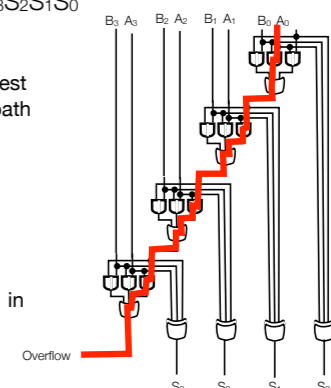
$$A_3A_2A_1A_0 + B_3B_2B_1B_0 = S_3S_2S_1S_0$$

- Depth of a circuit is the longest (most gates to go through) path

- Overflow has depth 8

- $S_3$  has depth 7

- In general,  $S_i$  has depth  $2i+1$  in Ripple-Carry Adder



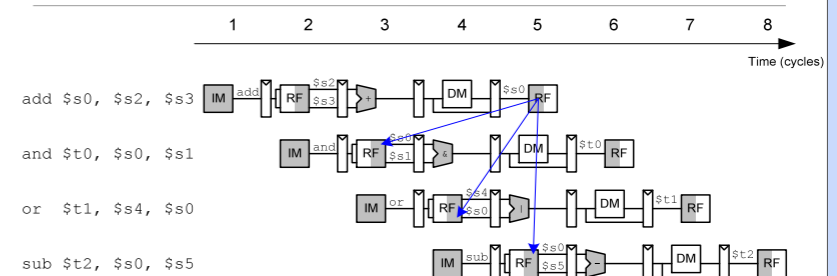
# Example 2: Pipelines = Parallel Execution

- Emphasized that the benefit of pipelining was **improved throughput** thanks to parallel **instruction execution**
- Extended discussion of throughput v. latency, and how parallel resources can boost the former. Employed analogy to web host adding servers. Can serve more clients, but latency is unchanged for each client.
- **Observation:** as with first serial then parallel programming instructions, students struggle to shift their mental model from a serial (single cycle processor) to a parallel processor with multiple instructions in-flight at once.

## Pipeline Hazards

- Hazards arise because a pipeline executes multiple instructions at once
- At runtime processors detect conflicts between concurrent instructions
- When possible resolve with bypassing (allowing all to proceed in parallel)
- When not possible stall processor (making execution more serial)

### Data Hazard



#### • Handling them:

- Insert nops in code at compile time
- Rearrange code at compile time
- Forward data at run time
- Stall the processor at run time



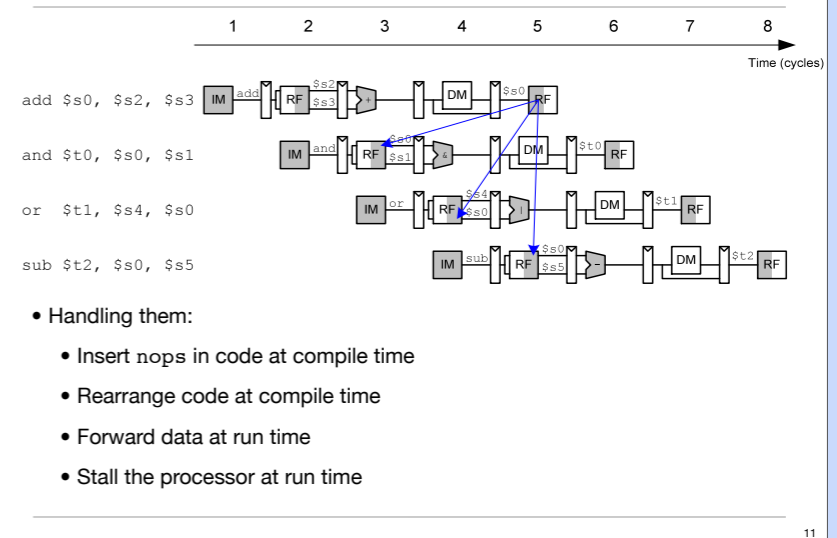
# Example 3: Advanced Microarchitectures

- At the end of the semester, build on course foundations to expose student to advanced microarchitectures. Emphasize that vast majority of **optimizations are designed with the goal of doing more work in parallel.**
- **Observation:** students are excited by topics they have heard the name of (GPUs, multicore). Despite these topics not being on the exam, students still asked more questions in these lectures than the others.

## Pipeline Hazards

- Hazards arise because a pipeline executes multiple instructions at once
- At runtime processors detect conflicts between concurrent instructions
- When possible resolve with bypassing (allowing all to proceed in parallel)
- When not possible stall processor (making execution more serial)

### Data Hazard

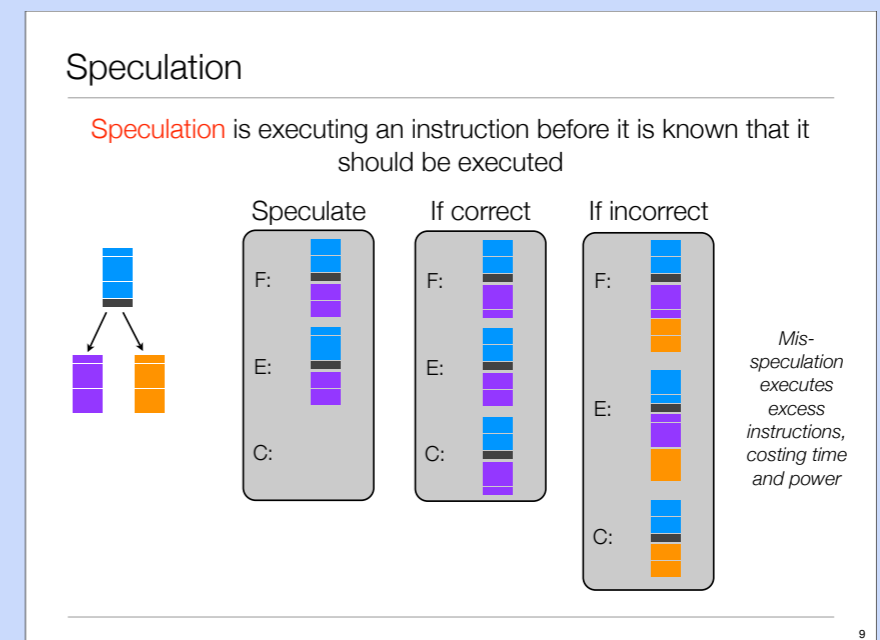


# Example 3: Advanced Microarchitectures

- At the end of the semester, build on course foundations to expose student to advanced microarchitectures. Emphasize that vast majority of **optimizations are designed with the goal of doing more work in parallel.**
- **Observation:** Given basic foundation and time constraints, must introduce these topics only at the conceptual level.
- **Observation:** students are excited by topics they have heard the name of (GPUs, multicore).

## Speculative Execution

- ILP is limited by branches (students learned this when we taught control hazards in pipelines)
- To find more instructions to keep a pipeline full, some processors guess the branch outcome
- Brief overview branch prediction strategies
- If the prediction is correct - hooray! the pipeline stayed full of useful work
- If the prediction was wrong - we lost time doing unnecessary work



# Example 4: Multicore

- The students have all heard of multicore, and often come in thinking more cores is better. They are very surprised to hear about why multicore can make things more difficult (on the software side) and then curious to understand why we went multicore anyway.

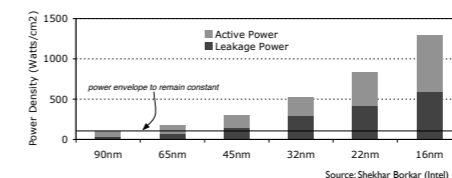
## Models of uniprocessor performance and power

- Must be extremely simple, but also scientifically sound
- Performance proportional to  $\text{SQRT}(\text{Area})$
- Power proportional to Area

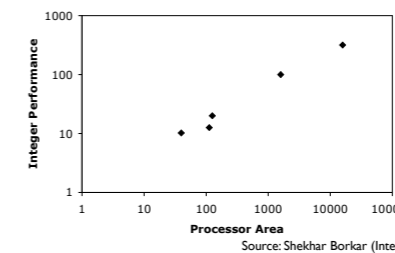
## Comparison of abstract processor scenarios

- Assume a base processor with unit performance and area
- Given 4x more transistors you have two options
  - Massive uncore with  $\text{Power}=4$ ,  $\text{Performance}=2$
  - 4way multicore with  $\text{Power}=4$ ,  $\text{Performance}=4$
- From efficiency standpoint multicores make more sense than massive cores
- However (bit caveat) performance relies on there being software to exploit it

Chip Area and Power Consumption



With leakage power dominating, power consumption roughly proportional to transistor count



Pollack's Law: Processor performance grows with sqrt of area

The Resulting Shift to Multicore



Perf = 1  
Power = 1

Perf = 2  
Power = 4

Perf = 4  
Power = 4

# Evaluation

---

- Evaluation primarily conducted via comparison to prior years offerings.
- Important concepts in parallelism that were emphasized in the first half of the semester deliberately tested on the midterm exam (example below).
- Midterm student survey indicated that despite a very full syllabus, addition of new material was acceptable thanks to integration with existing syllabus. I.e., One cannot bolt new topics on to crowded courses, but shifted emphasis on topics already integrated seems to be working OK.

## Dynamic circuit and timing behavior

Students were asked to (and successfully did) analyze the timing characteristics of an ALU containing four parallel computational paths.

2. (10 points) Design a 32-bit ALU (arithmetic logic unit) which computes either  $+$ ,  $-$ ,  $\times$ ,  $\div$  of two 32-bit values  $A$  and  $B$  putting the result on 32-bit output  $C$ . You may assume the presence of existing a 32-bit adder, subtractor, multiplier, and divider, and you may incorporate them in your design as black boxes. The particular operation is controlled by an additional two bit input  $S$ , where

*NB: In addition to the arithmetic modules, you may incorporate muxes, decoders, and other standard components as you find helpful.*

- if  $S = 00$ ,  $C = A + B$
- if  $S = 01$ ,  $C = A - B$
- if  $S = 10$ ,  $C = A \times B$
- if  $S = 11$ ,  $C = A \div B$

(a) Draw a schematic for your ALU design.

(b) Give an expression for the propagation delay ( $TP_{ALU}$ ) and contamination delay ( $TC_{ALU}$ ) of your ALU, in terms of the propagation and contamination delays of the sub-modules. For example,  $TP_+$  is the propagation delay of the adder, and  $TC_{MUX}$  is the contamination delay of a mux.

