

Computer Architecture

Lecture 2: Fundamentals, Memory Hierarchy, Caches

Prof. Onur Mutlu

ETH Zurich

Fall 2017

21 September 2017

Agenda for Today

- Finish up logistics from last lecture
- Why study computer architecture?
- Some fundamental concepts
- Memory hierarchy
- Caches

Takeaway From Lecture 1

Breaking the abstraction layers
(between components and
transformation hierarchy levels)

and knowing what is underneath

enables you to **understand** and
solve problems

Review: Major High-Level Goals of This Course

- Understand the principles
- Understand the precedents
- Based on such understanding:
 - Enable you to evaluate tradeoffs of different designs and ideas
 - Enable you to develop principled designs
 - Enable you to develop novel, out-of-the-box designs
- The focus is on:
 - Principles, precedents, and how to use them for new designs
- In Computer Architecture

A Note on Hardware vs. Software

- This course might seem like it is only “Computer Hardware”
- However, you will be much more capable if you master both hardware and software (and the interface between them)
 - Can develop better software if you understand the hardware
 - Can design better hardware if you understand the software
 - Can design a better computing system if you understand both
- This course covers the HW/SW interface and microarchitecture
 - We will focus on tradeoffs and how they affect software

What Do I Expect From You?

- **Required background:** Digital circuits course, programming, an open mind willing to take in many exciting concepts.
- **Learn the material thoroughly**
 - attend lectures, do the readings, do the exercises, do the labs
- **Work hard:** this will be a hard, but fun & informative course
- **Ask questions, take notes, participate**
- **Perform the assigned readings**
- **Come to class, participate**
- **Start early**
- If you want feedback, come to office hours
- Remember “**Chance favors the prepared mind.**” (Pasteur)



What Do I Expect From You?

- How you prepare and manage your time is very important
- There will be many lab and homework assignments
 - They will take time
 - Start early, work hard
- This will be a heavy course
 - However, you will learn a lot of fascinating topics and understand how a computing platform works
 - And, it will hopefully change how you look at and think about designs around you

How Will You Be Evaluated?

- Project assignments: 35%
 - Midterm exam: 25%
 - Final exam: 25%
 - Homeworks: 15%
-
- More on this later

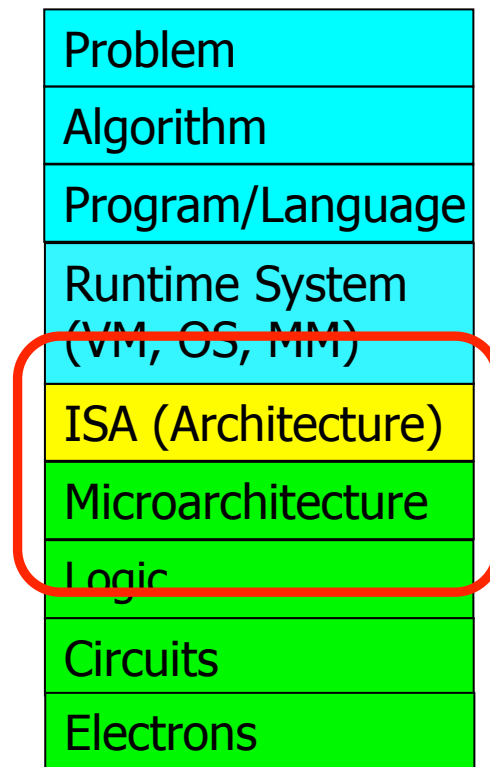
What Will You Learn

- **Computer Architecture:** The science and art of designing, selecting, and interconnecting hardware components and designing the hardware/software interface to create a computing system that meets functional, performance, energy consumption, cost, and other specific goals.
- **Traditional definition:** “The term *architecture* is used here to describe the attributes of a computer system as seen by a programmer, i.e., the conceptual structure and behavior as distinct from the organization of the hardware and controls, the logic design, and the physical implementation.” *Gene Amdahl, IBM Systems Journal*, 1964



Dr. Amdahl holding a 100gate LSI air-cooled chip. On his desk is a circuit board with the chips on it. This circuit board was for an Amdahl 470 V/6 (photograph dated March 1973).

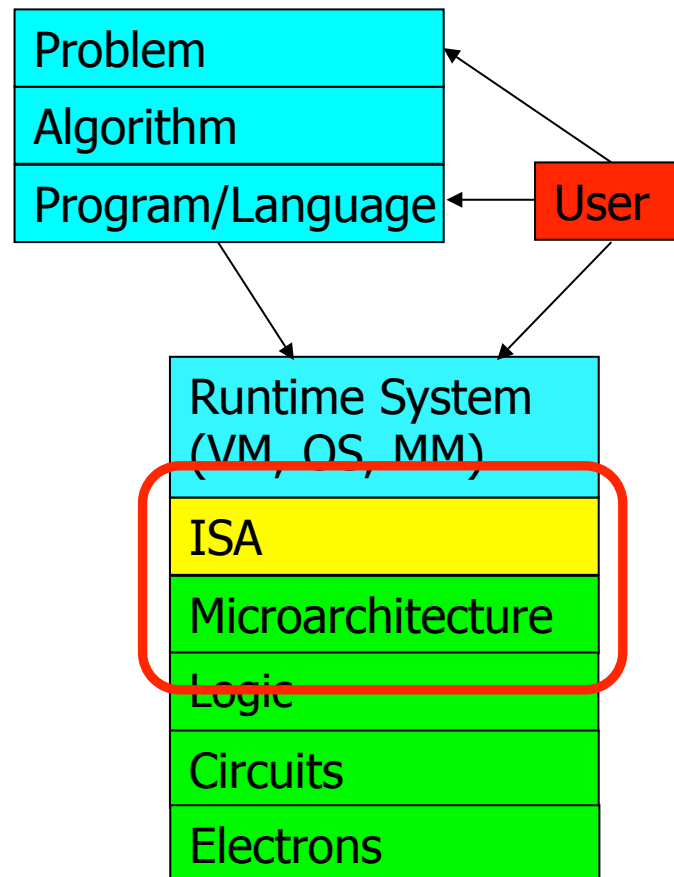
Computer Architecture in Levels of Transformation



- Read: Patt, "[Requirements, Bottlenecks, and Good Fortune: Agents for Microprocessor Evolution](#)," Proceedings of the IEEE 2001.

Levels of Transformation, Revisited

- A user-centric view: computer designed for users



- The entire stack should be optimized for user

What Will You Learn?

- Fundamental principles and tradeoffs in designing the hardware/software interface and major components of a modern programmable microprocessor
 - Focus on state-of-the-art (and some recent research and trends)
 - Trade-offs and how to make them
- How to design, implement, and evaluate a functional modern processor
 - Semester-long lab assignments
 - A combination of RTL implementation and higher-level simulation
 - Focus is functionality first (then, on “how to do even better”)
- How to think critically and broadly
- How to work efficiently

Course Goals

- Goal 1: To familiarize those interested in computer system design with both fundamental operation principles and design tradeoffs of processor, memory, and platform architectures in today's systems.
 - ❑ Strong emphasis on fundamentals, design tradeoffs, key current/future issues
 - ❑ Strong emphasis on looking backward, forward, up and down
- Goal 2: To provide the necessary background and experience to design, implement, and evaluate a modern processor by performing hands-on RTL and C-level implementation.
 - ❑ Strong emphasis on functionality, hands-on design & implementation, and efficiency.
 - ❑ Strong emphasis on making things work, realizing ideas

Course Website

- <http://safari.ethz.ch/architecture>
- All slides, lecture videos, readings, assignments to be posted
- Plus other useful information for the course
- Check frequently for announcements and due dates

Homework 0

- Due Sep 27
 - https://safari.ethz.ch/farm/architecture_fs17/doku.php?id=homeworks
- Information about yourself
- All future grading is predicated on homework 0

Heads Up

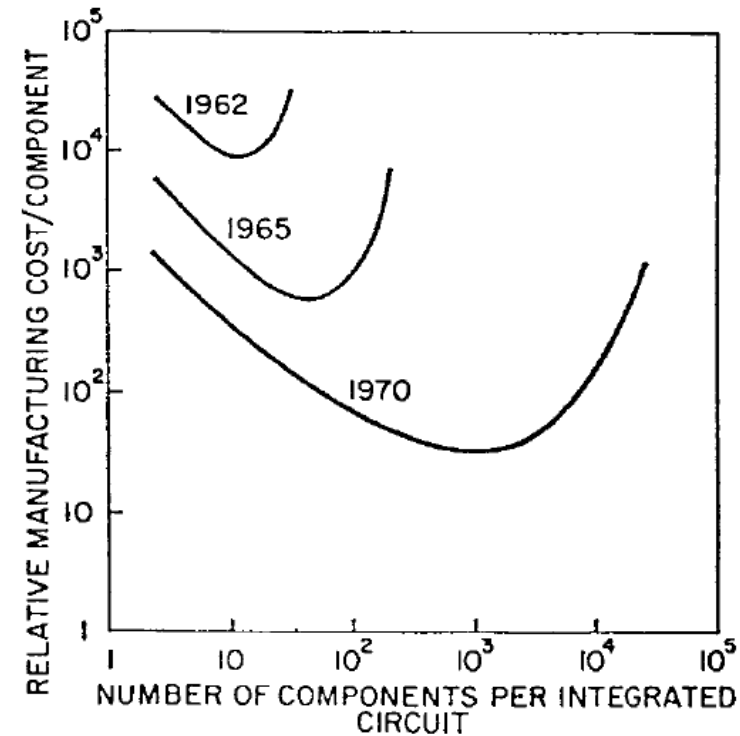
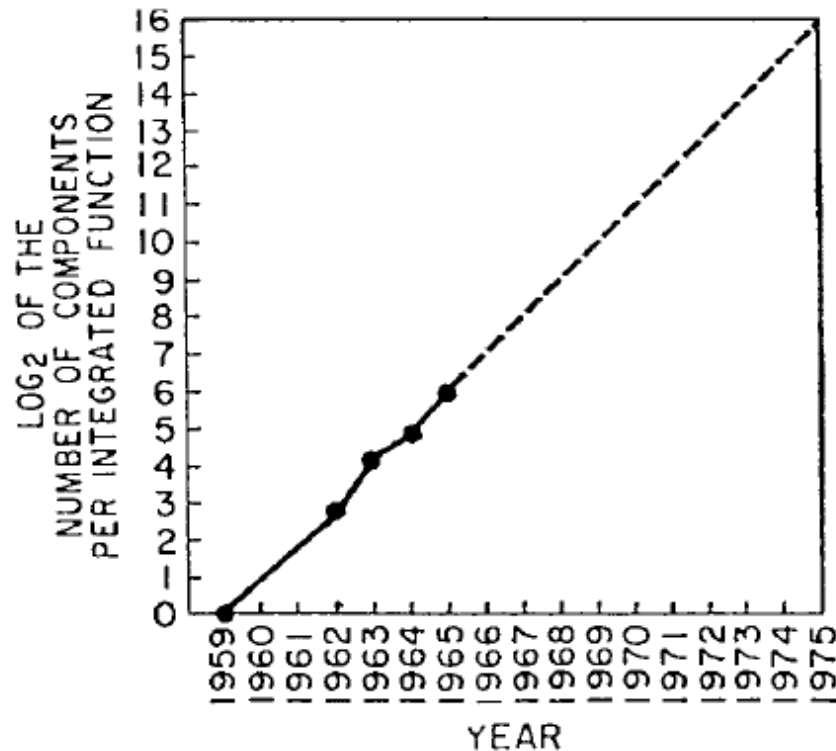
- We will have a few required review assignments
 - Due likely end of next week
- HW1 will be out early next week
 - Due in ~2 weeks
- Lab 1 will be out mid next week
 - Due in ~2 weeks
- Check the website. Will also be announced in lecture

Why Study Computer Architecture?

What is Computer Architecture?

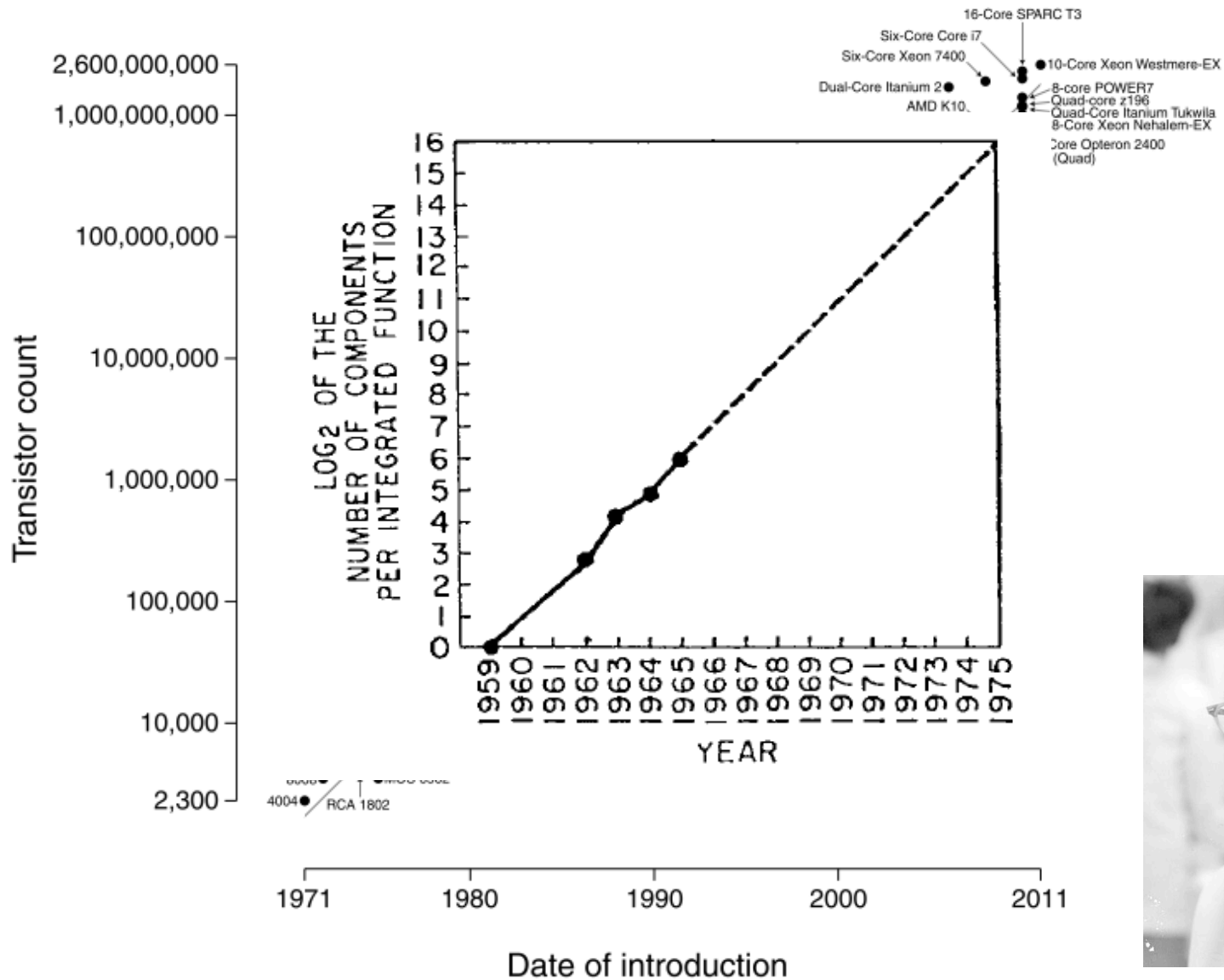
- The science and art of designing, selecting, and interconnecting hardware components and designing the hardware/software interface to create a computing system that meets functional, performance, energy consumption, cost, and other specific goals.
- We will soon distinguish between the terms *architecture*, and *microarchitecture*.
 - Actually, we have, in Digital Circuits course

An Enabler: Moore's Law



Moore, “Cramming more components onto integrated circuits,”
Electronics Magazine, 1965. Component counts double every other year

Microprocessor Transistor Counts 1971-2011 & Moore's Law



Number of transistors on an integrated circuit doubles ~ every two years

Recommended Reading

- Moore, “Cramming more components onto integrated circuits,” Electronics Magazine, 1965.
- Only 3 pages
- A quote:
“With unit cost falling as the number of components per circuit rises, by 1975 economics may dictate squeezing as many as 65 000 components on a single silicon chip.”
- Another quote:
“Will it be possible to remove the heat generated by tens of thousands of components in a single silicon chip?”

What Do We Use These Transistors for?

- Your readings for this week should give you an idea...
- Required
 - Patt, “Requirements, Bottlenecks, and Good Fortune: Agents for Microprocessor Evolution,” Proceedings of the IEEE 2001.
- Required for Review as part of HW 1
 - Moscibroda and Mutlu, “Memory Performance Attacks: Denial of Memory Service in Multi-Core Systems,” USENIX Security 2007.
 - Liu+, “RAIDR: Retention-Aware Intelligent DRAM Refresh,” ISCA 2012.
 - Kim+, “Flipping Bits in Memory Without Accessing Them: An Experimental Study of DRAM Disturbance Errors,” ISCA 2014.

Why Study Computer Architecture?

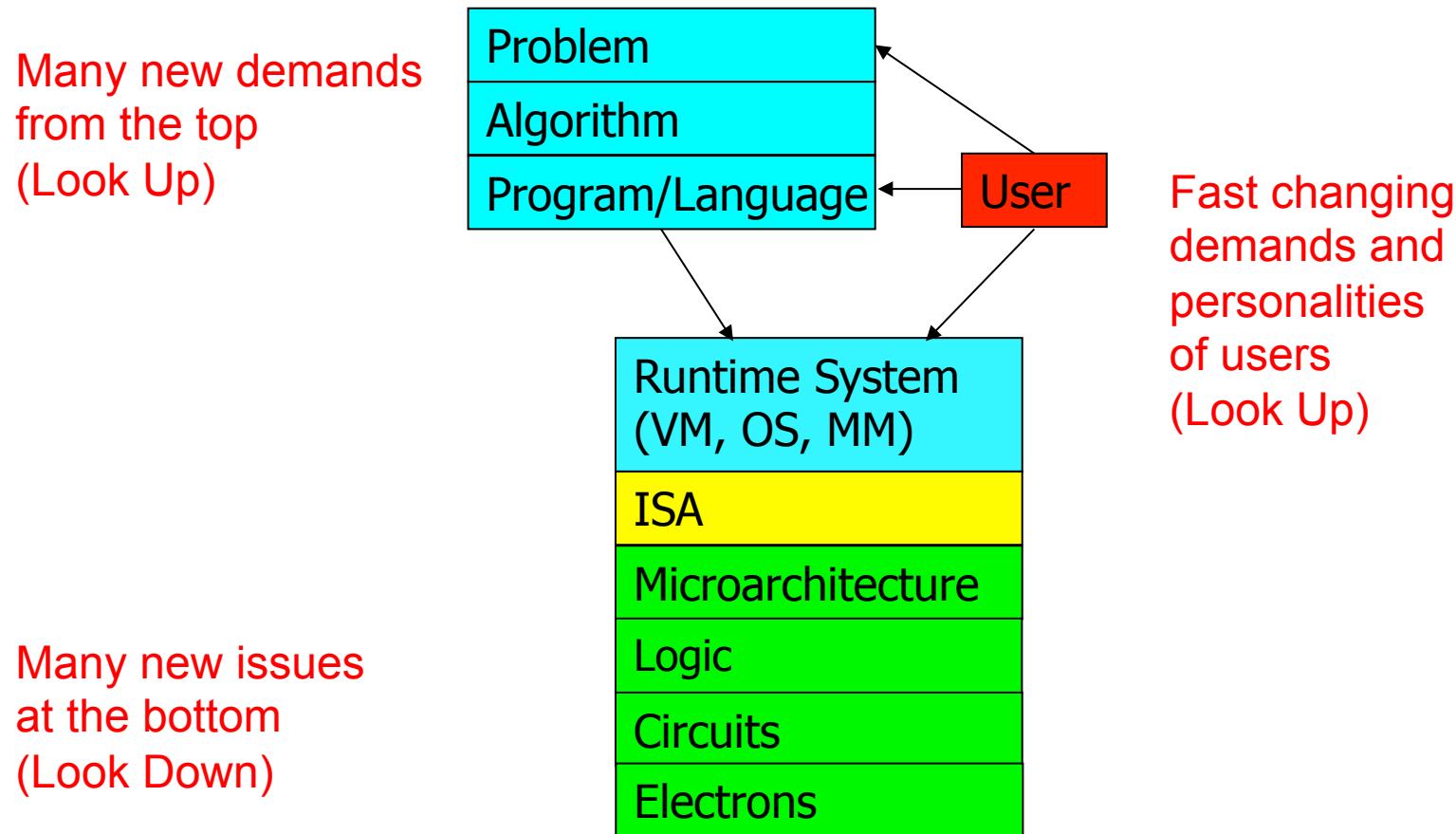
- **Enable better systems:** make computers faster, cheaper, smaller, more reliable, ...
 - By exploiting advances and changes in underlying technology/circuits
- **Enable new applications**
 - Life-like 3D visualization 20 years ago? Virtual reality?
 - Self-driving cars?
 - Personalized genomics? Personalized medicine?
- **Enable better solutions** to problems
 - Software innovation is built on trends and changes in computer architecture
 - > 50% performance improvement per year has enabled this innovation
- **Understand why computers work the way they do**

Computer Architecture Today (I)

- Today is a very exciting time to study computer architecture
 - Industry is in a large paradigm shift (to multi-core and beyond) – many different potential system designs possible
 - **Many difficult problems** *motivating and caused by* the shift
 - Power/energy constraints → heterogeneity?
 - Complexity of design → multi-core and heterogeneity?
 - Difficulties in technology scaling → new technologies?
 - Memory wall/gap → processing in memory?
 - Reliability wall/issues → new technologies?
 - Programmability wall/problem
 - Huge hunger for data and new data-intensive applications
 - No clear, definitive answers to these problems
-

Computer Architecture Today (II)

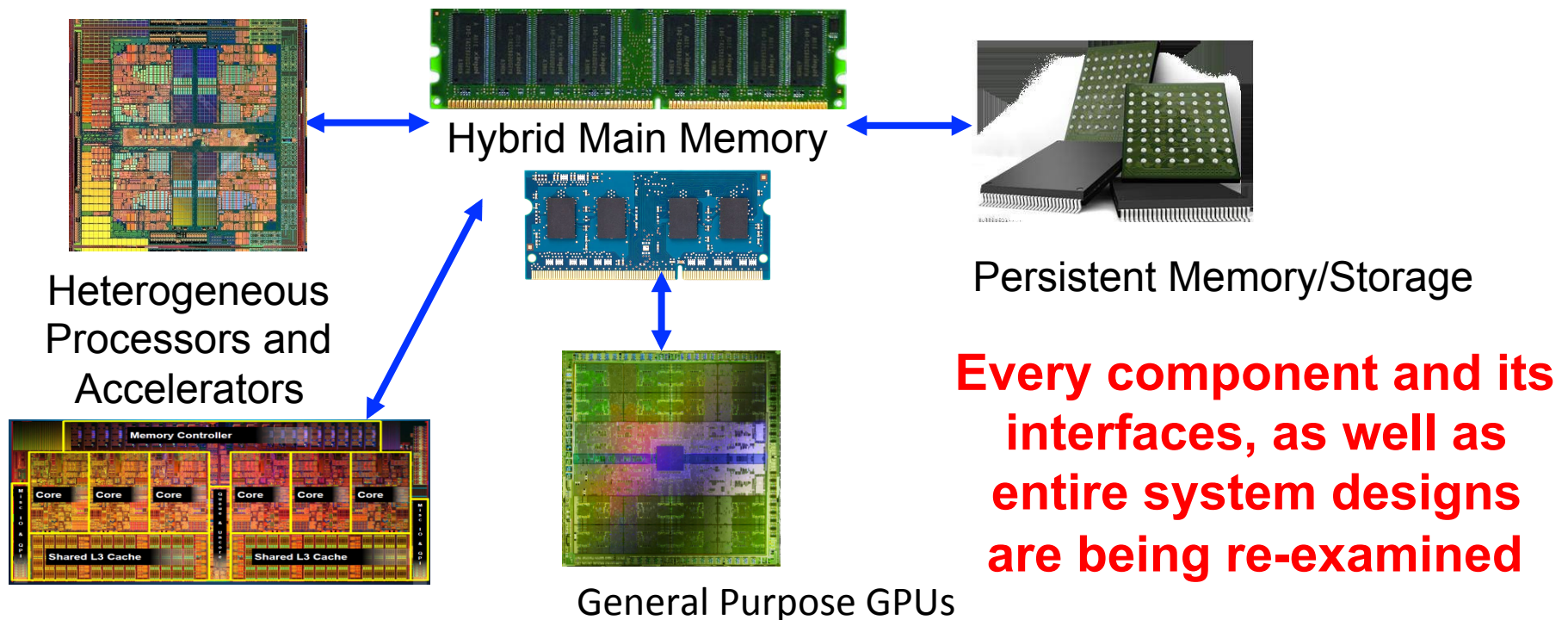
- These problems affect all parts of the computing stack – if we do not change the way we design systems



- No clear, definitive answers to these problems

Computer Architecture Today (III)

- Computing landscape is very different from 10-20 years ago
- Both UP (software and humanity trends) and DOWN (technologies and their issues), FORWARD and BACKWARD, and the resulting requirements and constraints



Computer Architecture Today (IV)

- You can revolutionize the way computers are built, if you understand both the hardware and the software (and change each accordingly)
- You can invent new paradigms for computation, communication, and storage
- Recommended book: Thomas Kuhn, “[The Structure of Scientific Revolutions](#)” (1962)
 - Pre-paradigm science: no clear consensus in the field
 - Normal science: dominant theory used to explain/improve things (business as usual); exceptions considered anomalies
 - Revolutionary science: underlying assumptions re-examined

Computer Architecture Today (IV)

- You can revolutionize the way computers are built, if you understand both hardware and software (and change each accordingly)

- You can improve communication

- Recommendation: Scientific Illustration

- Pre-prepare
- Normal scientific things (books)
- Revolutionary



ure of
eld
improve
anomalies
examined

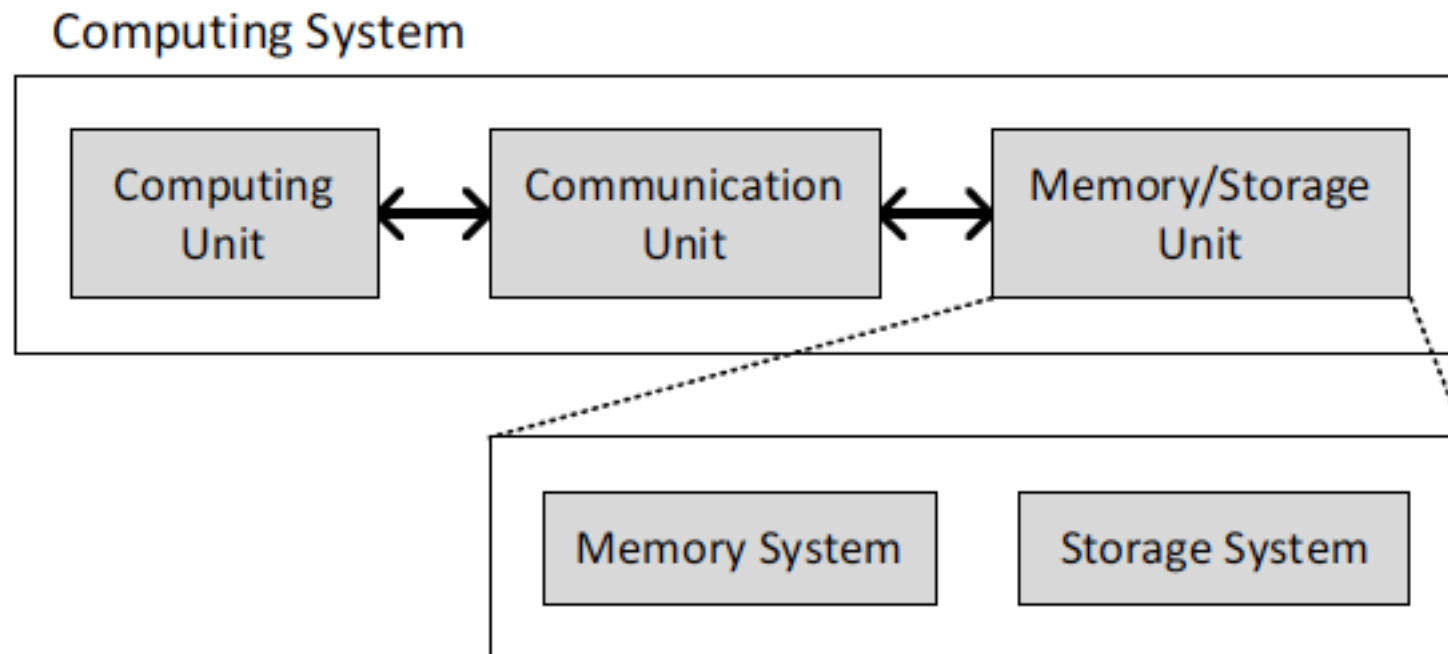
... but, first ...

- Let's understand the fundamentals...
- You can change the world only if you understand it well enough...
 - Especially the past and present dominant paradigms
 - And, their advantages and shortcomings – tradeoffs
 - And, what remains fundamental across generations
 - And, what techniques you can use and develop to solve problems

Fundamental Concepts

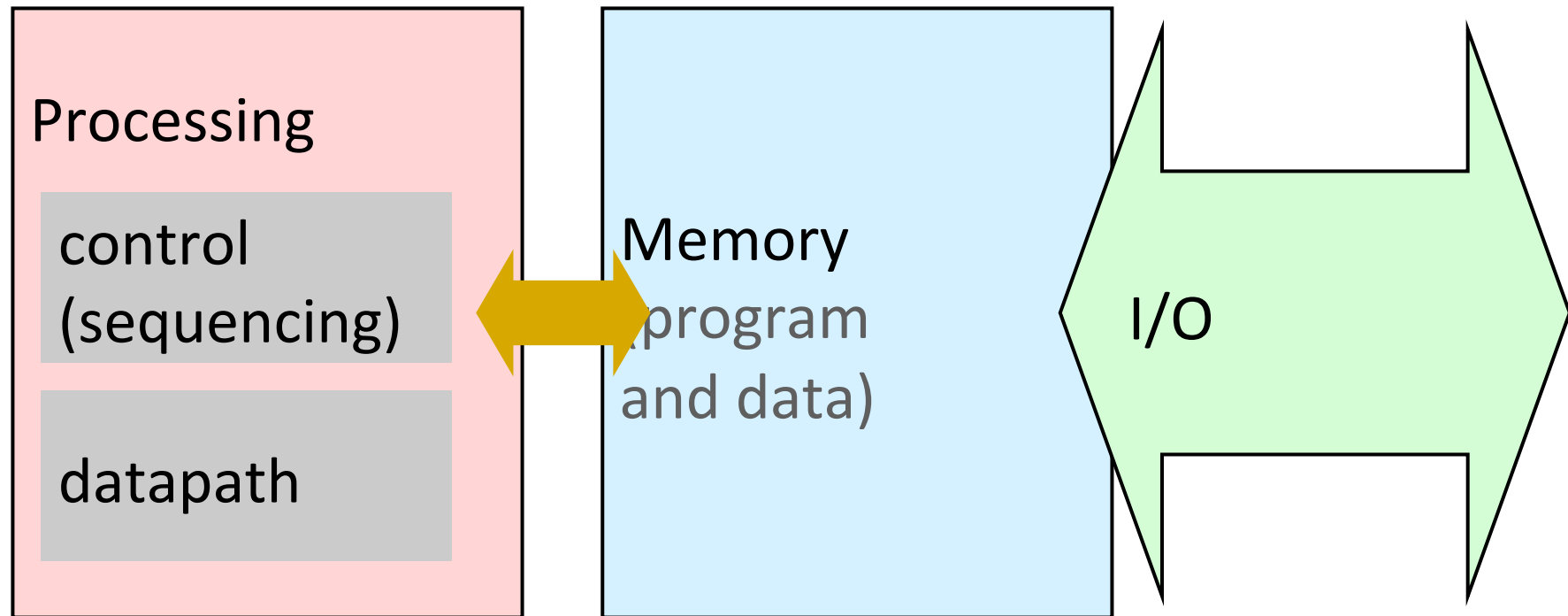
What is A Computer?

- Three key components
- Computation
- Communication
- Storage (memory)



What is A Computer?

- We will cover all three components



The Von Neumann Model/Architecture

- Also called *stored program computer* (instructions in memory). Two key properties:
- Stored program
 - Instructions stored in a linear memory array
 - Memory is unified between instructions and data
 - The interpretation of a stored value depends on the control signals
- Sequential instruction processing
 - One instruction processed (fetched, executed, and completed) at a time
 - Program counter (instruction pointer) identifies the current instr.
 - Program counter is advanced sequentially except for control transfer instructions

The Von Neumann Model/Architecture

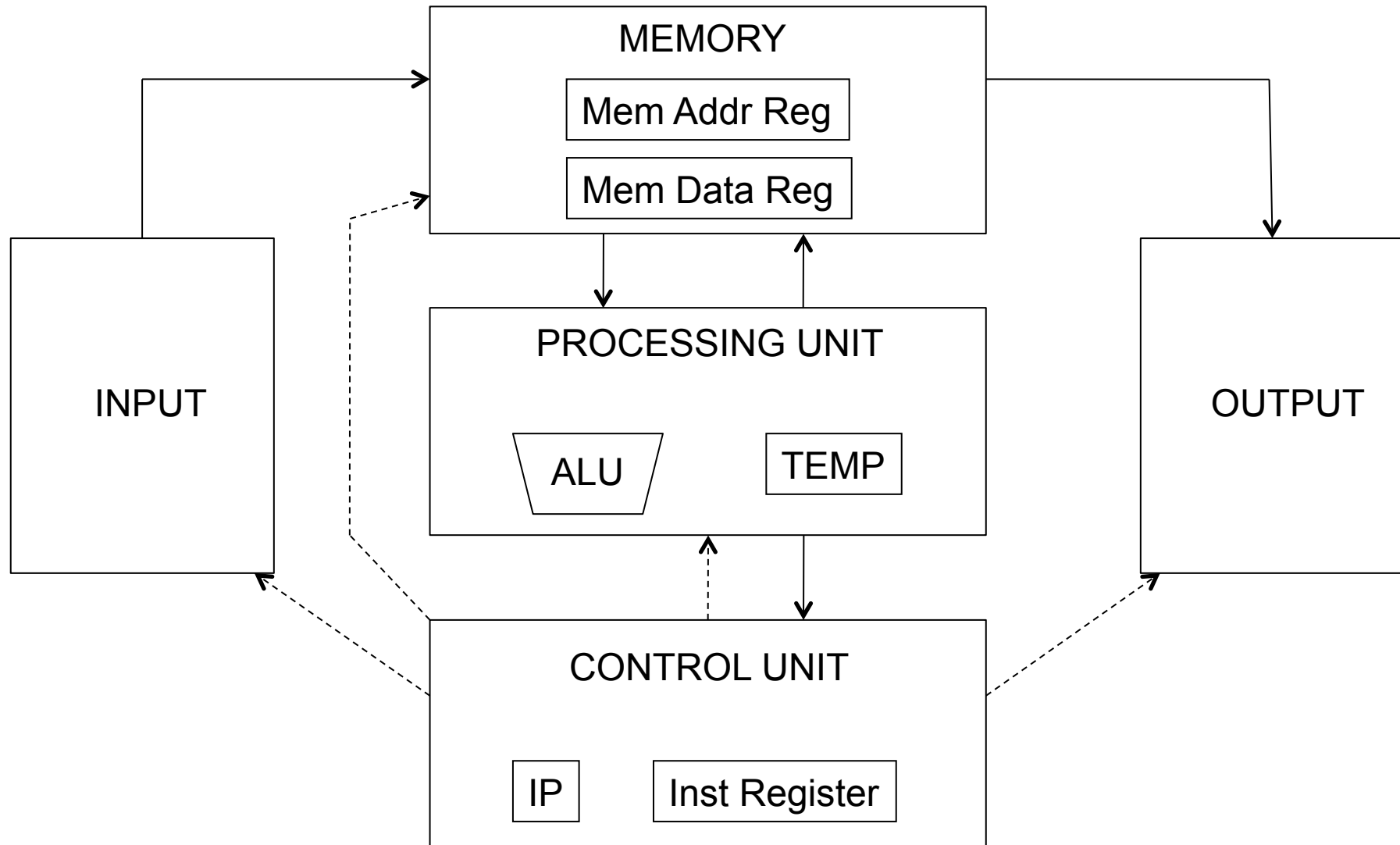
- Recommended readings

- Burks, Goldstein, von Neumann, "Preliminary discussion of the logical design of an electronic computing instrument," 1946.
- Patt and Patel book, Chapter 4, "The von Neumann Model"

- Stored program

- Sequential instruction processing

The Von Neumann Model (of a Computer)



The Von Neumann Model (of a Computer)

- Q: Is this the only way that a computer can operate?
- A: No.
- Qualified Answer: No, but it has been the dominant way
 - i.e., the dominant paradigm for computing
 - for N decades

The Dataflow Model (of a Computer)

- Von Neumann model: An instruction is fetched and executed in **control flow order**
 - As specified by the **instruction pointer**
 - Sequential unless explicit control flow instruction

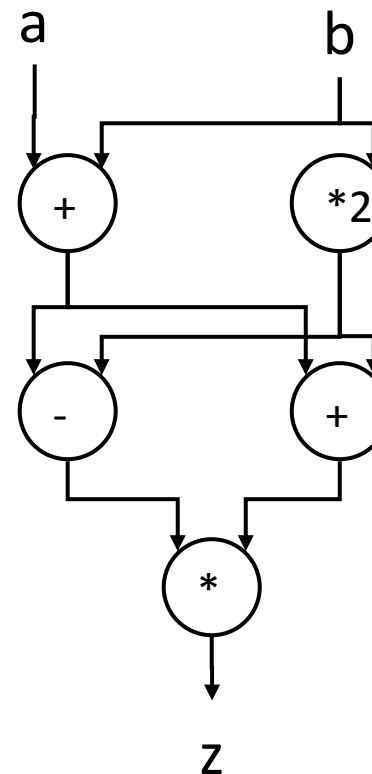
- Dataflow model: An instruction is fetched and executed in **data flow order**
 - i.e., when its operands are ready
 - i.e., there is **no instruction pointer**
 - Instruction ordering specified by data flow dependence
 - Each instruction specifies “who” should receive the result
 - An instruction can “fire” whenever all operands are received
 - Potentially many instructions can execute at the same time
 - Inherently more parallel

Von Neumann vs Dataflow

- Consider a Von Neumann program
 - What is the significance of the program order?
 - What is the significance of the storage locations?

```
v <= a + b;  
w <= b * 2;  
x <= v - w  
y <= v + w  
z <= x * y
```

Sequential

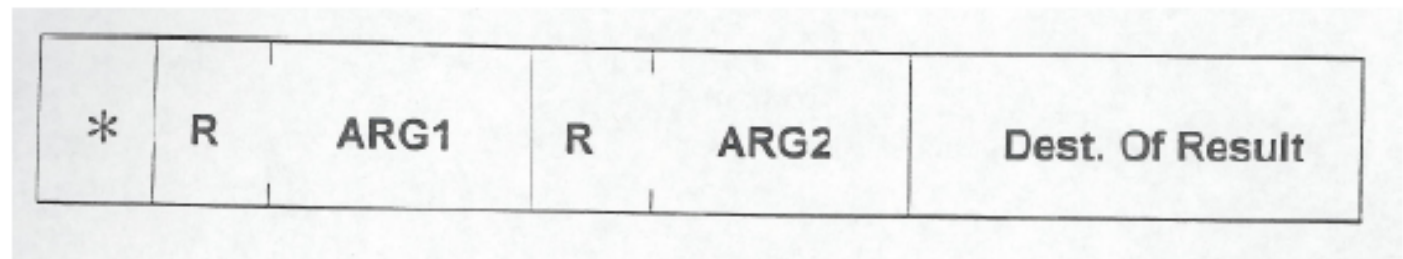
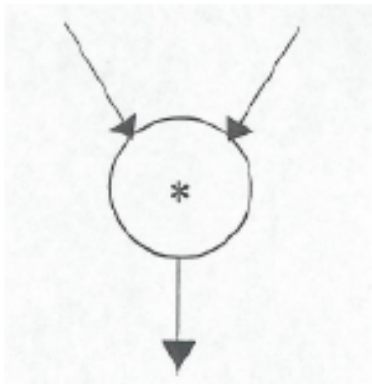


Dataflow

- Which model is more natural to you as a programmer?

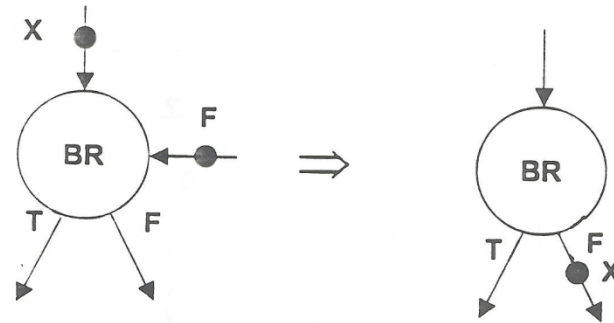
More on Data Flow

- In a data flow machine, a program consists of data flow nodes
 - A data flow node fires (fetched and executed) when all its inputs are ready
 - i.e. when all inputs have tokens
- Data flow node and its ISA representation

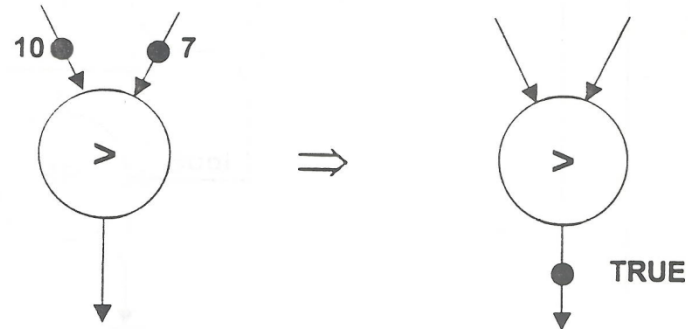


Data Flow Nodes

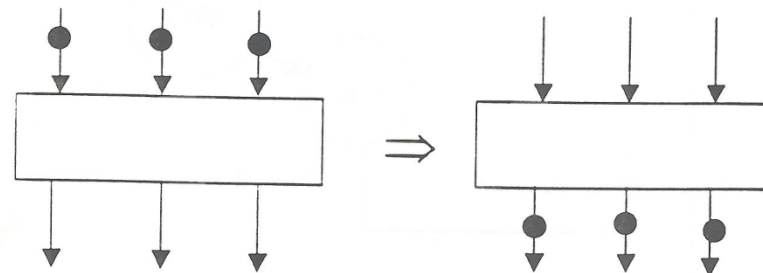
**Conditional*



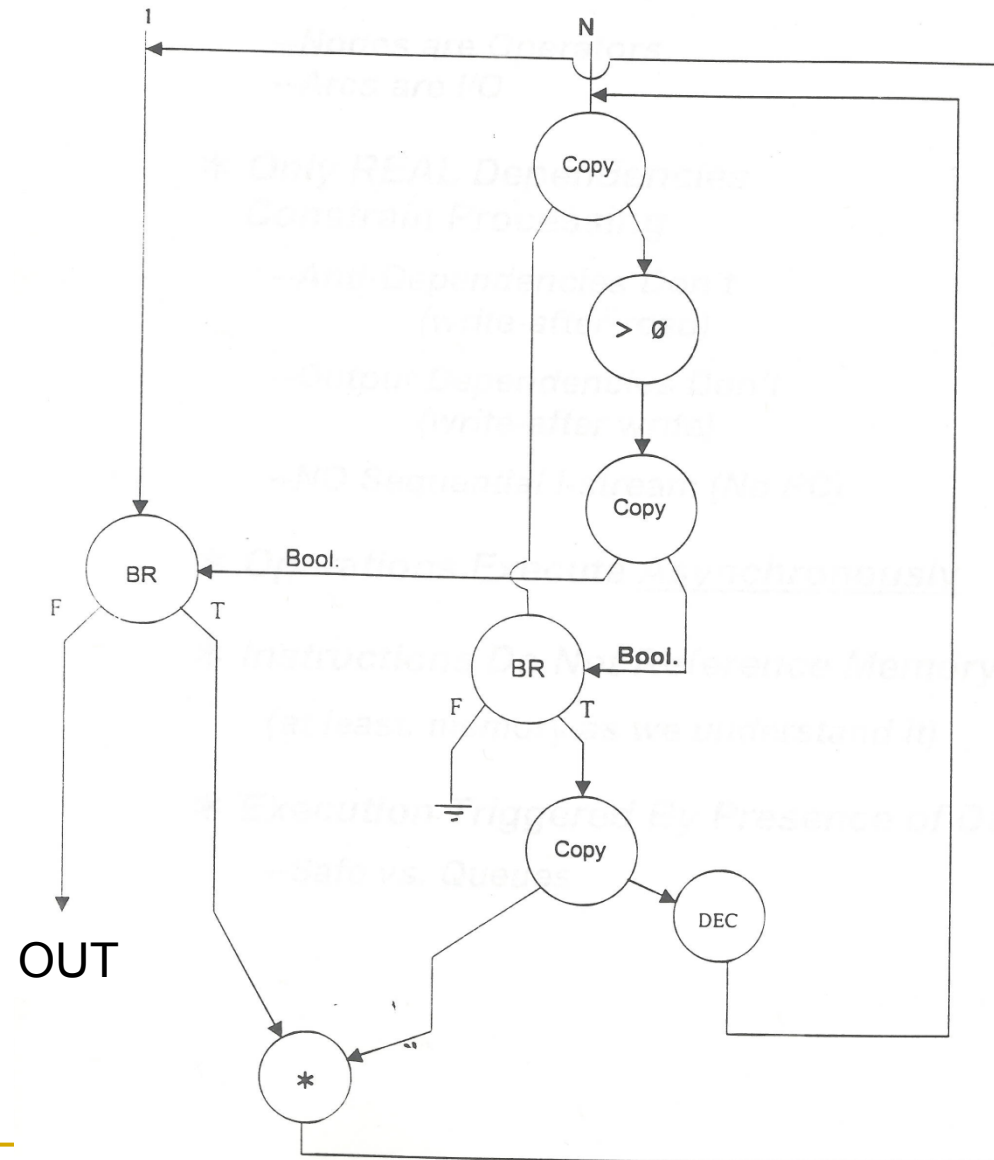
**Relational*



**Barrier Synchron*



An Example Data Flow Program



ISA-level Tradeoff: Instruction Pointer

- Do we need an instruction pointer in the ISA?
 - Yes: Control-driven, sequential execution
 - An instruction is executed when the IP points to it
 - IP automatically changes sequentially (except for control flow instructions)
 - No: Data-driven, parallel execution
 - An instruction is executed when all its operand values are available ([data flow](#))
- Tradeoffs: MANY high-level ones
 - Ease of programming (for average programmers)?
 - Ease of compilation?
 - Performance: Extraction of parallelism?
 - Hardware complexity?

ISA vs. Microarchitecture Level Tradeoff

- A similar tradeoff (control vs. data-driven execution) can be made at the microarchitecture level
- ISA: Specifies how the programmer sees instructions to be executed
 - Programmer sees a sequential, control-flow execution order vs.
 - Programmer sees a data-flow execution order
- Microarchitecture: How the underlying implementation actually executes instructions
 - Microarchitecture can execute instructions in any order as long as it obeys the semantics specified by the ISA when making the instruction results visible to software
 - Programmer should see the order specified by the ISA

Let's Get Back to the Von Neumann Model

- But, if you want to learn more about dataflow...
- Dennis and Misunas, "A preliminary architecture for a basic data-flow processor," ISCA 1974.
- Gurd et al., "The Manchester prototype dataflow computer," CACM 1985.
- A later lecture or course
- If you are really impatient:
 - <http://www.youtube.com/watch?v=D2uue7izU2c>
 - <http://www.ece.cmu.edu/~ece740/f13/lib/exe/fetch.php?media=onur-740-fall13-module5.2.1-dataflow-part1.ppt>

The Von-Neumann Model

- All major *instruction set architectures* today use this model
 - x86, ARM, MIPS, SPARC, Alpha, POWER
- Underneath (at the microarchitecture level), the execution model of almost all *implementations (or, microarchitectures)* is very different
 - Pipelined instruction execution: *Intel 80486 uarch*
 - Multiple instructions at a time: *Intel Pentium uarch*
 - Out-of-order execution: *Intel Pentium Pro uarch*
 - Separate instruction and data caches
- But, what happens underneath that is *not* consistent with the von Neumann model is *not* exposed to software
 - Difference between ISA and microarchitecture

What is Computer Architecture?

- **ISA+implementation definition:** The science and art of designing, selecting, and interconnecting hardware components and designing the hardware/software interface to create a computing system that meets functional, performance, energy consumption, cost, and other specific goals.
- **Traditional (ISA-only) definition:** “The term *architecture* is used here to describe the attributes of a system as seen by the programmer, i.e., the conceptual structure and functional behavior as distinct from the organization of the dataflow and controls, the logic design, and the physical implementation.” *Gene Amdahl*, IBM Journal of R&D, April 1964

ISA vs. Microarchitecture

■ ISA

- ❑ Agreed upon interface between software and hardware
 - SW/compiler assumes, HW promises
- ❑ What the software writer needs to know to write and debug system/user programs

■ Microarchitecture

- ❑ Specific implementation of an ISA
- ❑ Not visible to the software

■ Microprocessor

- ❑ **ISA, uarch**, circuits
- ❑ “Architecture” = ISA + microarchitecture

Problem
Algorithm
Program
ISA
Microarchitecture
Circuits
Electrons

ISA vs. Microarchitecture

- What is part of ISA vs. Uarch?
 - Gas pedal: interface for “acceleration”
 - Internals of the engine: implement “acceleration”

- Implementation (uarch) can be various as long as it satisfies the specification (ISA)
 - Add instruction vs. Adder implementation
 - Bit serial, ripple carry, carry lookahead adders are all part of microarchitecture
 - x86 ISA has many implementations: 286, 386, 486, Pentium, Pentium Pro, Pentium 4, Core, ...

- Microarchitecture usually changes faster than ISA
 - Few ISAs (x86, ARM, SPARC, MIPS, Alpha) but many uarchs
 - *Why?*

ISA

- Instructions

- Opcodes, Addressing Modes, Data Types
- Instruction Types and Formats
- Registers, Condition Codes

- Memory

- Address space, Addressability, Alignment
- Virtual memory management

- Call, Interrupt/Exception Handling

- Access Control, Priority/Privilege

- I/O: memory-mapped vs. instr.

- Task/thread Management

- Power and Thermal Management

- Multi-threading support, Multiprocessor support



Intel® 64 and IA-32 Architectures
Software Developer's Manual

Volume 1:
Basic Architecture

Microarchitecture

- Implementation of the ISA under specific **design constraints and goals**
- Anything done in hardware without exposure to software
 - Pipelining
 - In-order versus out-of-order instruction execution
 - Memory access scheduling policy
 - Speculative execution
 - Superscalar processing (multiple instruction issue?)
 - Clock gating
 - Caching? Levels, size, associativity, replacement policy
 - Prefetching?
 - Voltage/frequency scaling?
 - Error correction?

Property of ISA vs. Uarch?

- ADD instruction's opcode
 - Number of general purpose registers
 - Number of ports to the register file
 - Number of cycles to execute the MUL instruction
 - Whether or not the machine employs pipelined instruction execution
-
- Remember
 - Microarchitecture: Implementation of the ISA under specific design constraints and goals

Design Point

- A set of design considerations and their importance
 - ❑ **leads to tradeoffs** in both ISA and uarch
- Considerations
 - ❑ Cost
 - ❑ Performance
 - ❑ Maximum power consumption
 - ❑ Energy consumption (battery life)
 - ❑ Availability
 - ❑ Reliability and Correctness
 - ❑ Time to Market
- Design point determined by the “Problem” space (application space), the intended users/*market*

Problem
Algorithm
Program
ISA
Microarchitecture
Circuits
Electrons

Application Space

- Dream, and they will appear...

Other examples of the application space that continue to drive the need for unique design points are the following:

- 1) scientific applications such as those whose computations control nuclear power plants, determine where to drill for oil, and predict the weather;
- 2) transaction-based applications such as those that handle ATM transfers and e-commerce business;
- 3) business data processing applications, such as those that handle inventory control, payrolls, IRS activity, and various personnel record keeping, whether the personnel are employees, students, or voters;
- 4) network applications, such as high-speed routing of Internet packets, that enable the connection of your home system to take advantage of the Internet;
- 5) guaranteed delivery (a.k.a. real time) applications that require the result of a computation by a certain critical deadline;
- 6) embedded applications, where the processor is a component of a larger system that is used to solve the (usually) dedicated application;
- 7) media applications such as those that decode video and audio files;
- 8) random software packages that desktop users would like to run on their PCs.

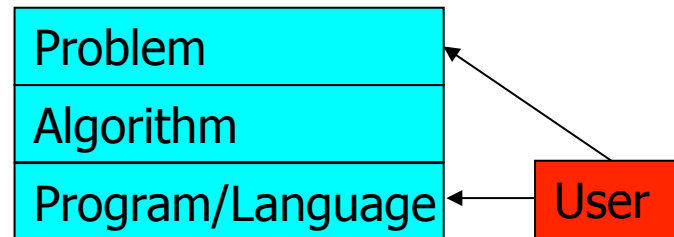
Each of these application areas has a very different set of characteristics. Each application area demands a different set of tradeoffs to be made in specifying the microprocessor to do the job.

Tradeoffs: Soul of Computer Architecture

- ISA-level tradeoffs
- Microarchitecture-level tradeoffs
- System and Task-level tradeoffs
 - How to divide the labor between hardware and software
- *Computer architecture is the science and art of making the appropriate trade-offs to meet a design point*
 - *Why art?*

Why Is It (Somewhat) Art?

New demands
from the top
(Look Up)

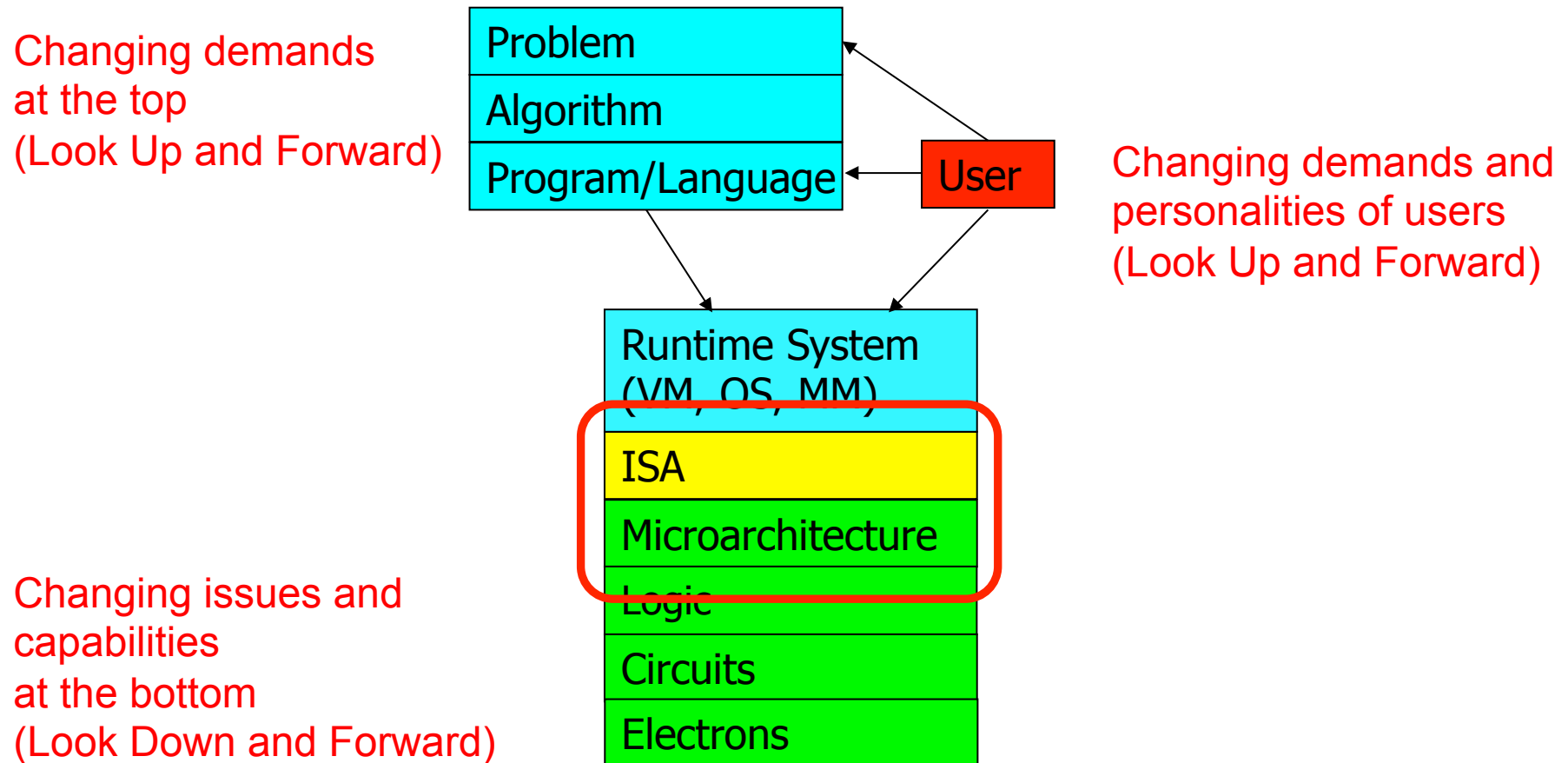


New demands and
personalities of users
(Look Up)

New issues and
capabilities
at the bottom
(Look Down)

- We do not (fully) know the future (applications, users, market)

Why Is It (Somewhat) Art?



- And, the future is not constant (it changes)!

How Can We Adapt to the Future

- This is part of the task of a good computer architect
- Many options (bag of tricks)
 - Keen insight and good design
 - Good use of fundamentals and principles
 - Efficient design
 - Heterogeneity
 - Reconfigurability
 - ...
 - Good use of the underlying technology
 - ...

We Covered a Lot of This in
Digital Circuits & Computer Architecture

One Slide Overview of Digital Circuits SS17

- Logic Design, Verilog, FPGAs
- ISA (MIPS)
- Single-cycle Microarchitectures
- Multi-cycle and Microprogrammed Microarchitectures
- Pipelining
- Issues in Pipelining: Control & Data Dependence Handling, State Maintenance and Recovery, ...
- Out-of-Order Execution
- Other Execution Paradigms
- Memory and Caches (very brief)

Covered Concurrent Execution Paradigms

- Pipelining
- Out-of-order execution
- Dataflow (at the ISA level)
- Superscalar Execution
- VLIW
- SIMD Processing (Vector and array processors, GPUs)
- Decoupled Access Execute
- Systolic Arrays

Digital Circuits Materials for Review (I)

- All Digital Circuits Lecture Videos Are Online:
 - <https://www.youtube.com/playlist?list=PL5Q2soXY2Zi-IXWTT7xoNYpst5-zdZQ6y>
- All Slides and Assignments Are Online:
 - http://www.syssec.ethz.ch/education/Digitaltechnik_17.html

Digital Circuits Materials for Review (II)

- Particularly useful and relevant lectures for this course
- Pipelining (Lecture 15)
 - <https://youtu.be/vBGVRURaxl8>
- Dependence handling (Lecture 16)
 - <https://youtu.be/B7bTbBRVdxA>
- Pipelining Issues (Lecture 17)
 - <https://youtu.be/C5ViR0dGILI>
- Out-of-order execution (Lecture 18)
 - <https://youtu.be/R5G05HstI3A>

This Course

- We will have more emphasis on
 - The memory system
 - Multiprocessing & multithreading
 - Parallel processing paradigms

- We will likely dig deeper on some Digital Circuits concepts (as time permits)
 - ISA
 - Branch handling
 - ...

Tentative Agenda (Upcoming Lectures)

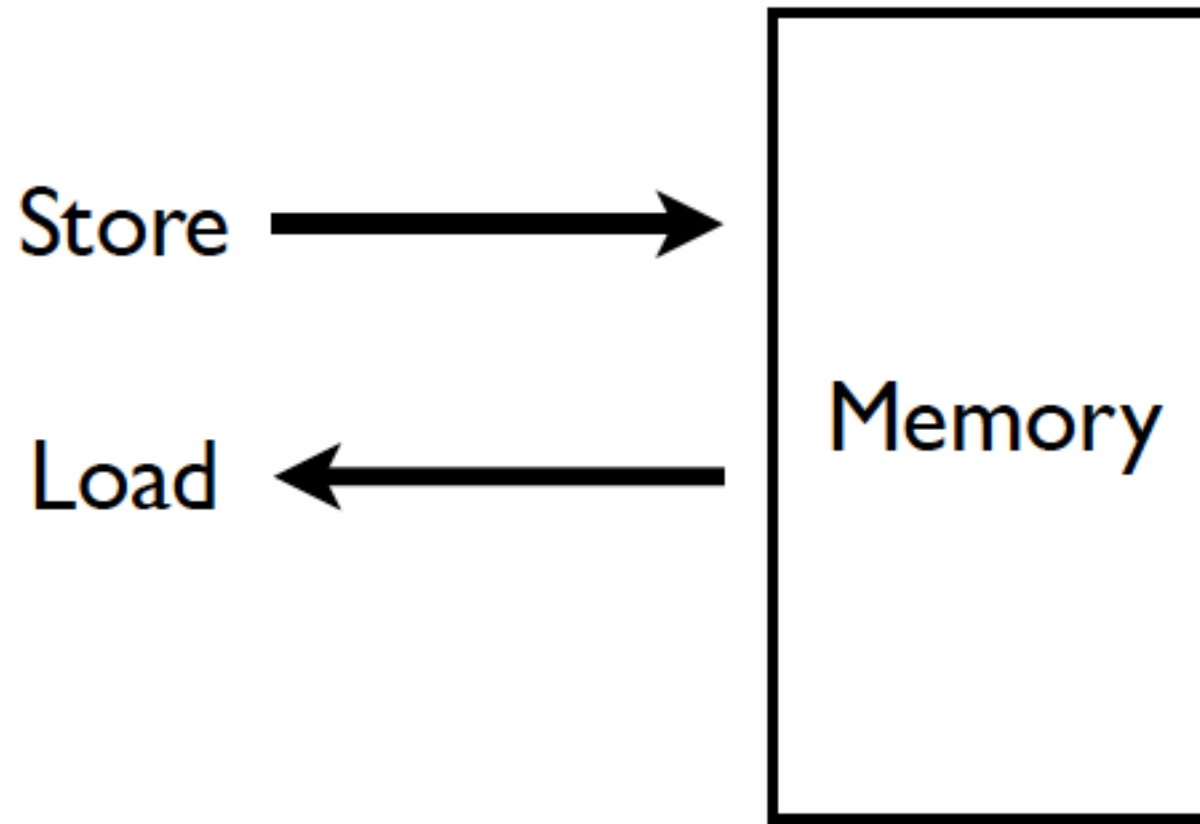
- The memory hierarchy
- Caches, caches, more caches (high locality, high bandwidth)
- Virtualizing the memory hierarchy
- Main memory: DRAM
- Main memory control, scheduling, interference, management
- Memory latency tolerance and prefetching techniques
- Non-volatile memory & emerging technologies

- Multiprocessors
- Coherence and consistency
- Interconnection networks
- Multi-core issues
- Multithreading

Optional Readings for Today & Next Week

- Memory Hierarchy and Caches
- Cache chapters from P&H: 5.1-5.3
- Memory/cache chapters from Hamacher+: 8.1-8.7
- An early cache paper by Maurice Wilkes
 - Wilkes, "Slave Memories and Dynamic Storage Allocation," IEEE Trans. On Electronic Computers, 1965.

Memory (Programmer's View)



Abstraction: Virtual vs. Physical Memory

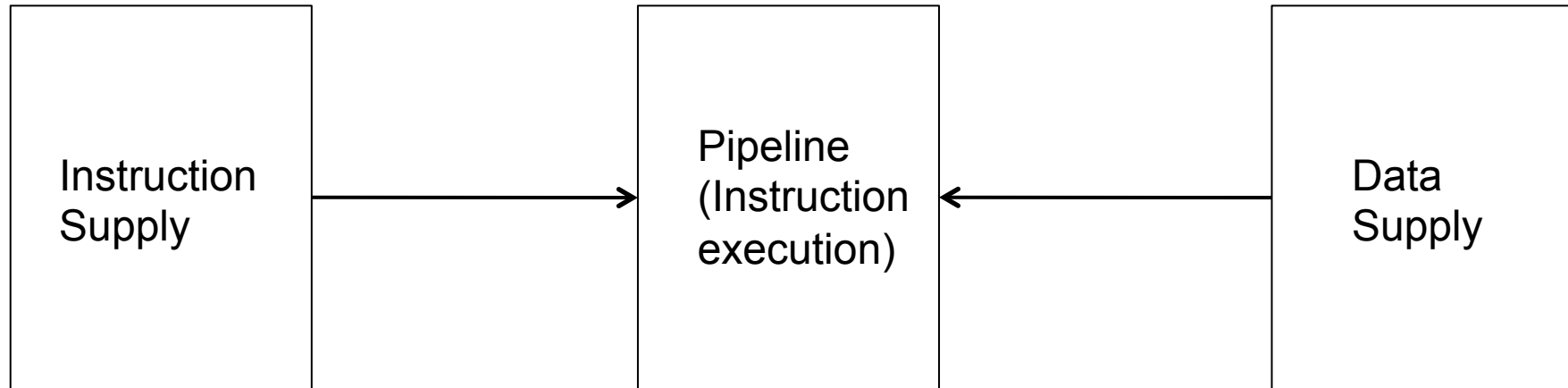
- **Programmer** sees **virtual memory**
 - Can assume the memory is “infinite”
 - Reality: **Physical memory** size is much smaller than what the programmer assumes
 - **The system** (system software + hardware, cooperatively) maps **virtual memory addresses** to **physical memory**
 - The system automatically manages the physical memory space **transparently to the programmer**
- + Programmer does not need to know the physical size of memory nor manage it → A small physical memory can appear as a huge one to the programmer → Life is easier for the programmer
- More complex system software and architecture

A classic example of the programmer/(micro)architect tradeoff

(Physical) Memory System

- You need a larger level of storage to manage a small amount of physical memory automatically
→ Physical memory has a backing store: disk
- We will first start with the physical memory system
- For now, ignore the virtual→physical indirection
- We will get back to it when the needs of virtual memory start complicating the design of physical memory...

Idealism



- Zero latency access

- Infinite capacity

- Zero cost

- Perfect control flow

- No pipeline stalls

- Perfect data flow
(reg/memory dependencies)

- Zero-cycle interconnect
(operand communication)

- Enough functional units

- Zero latency compute

- Zero latency access

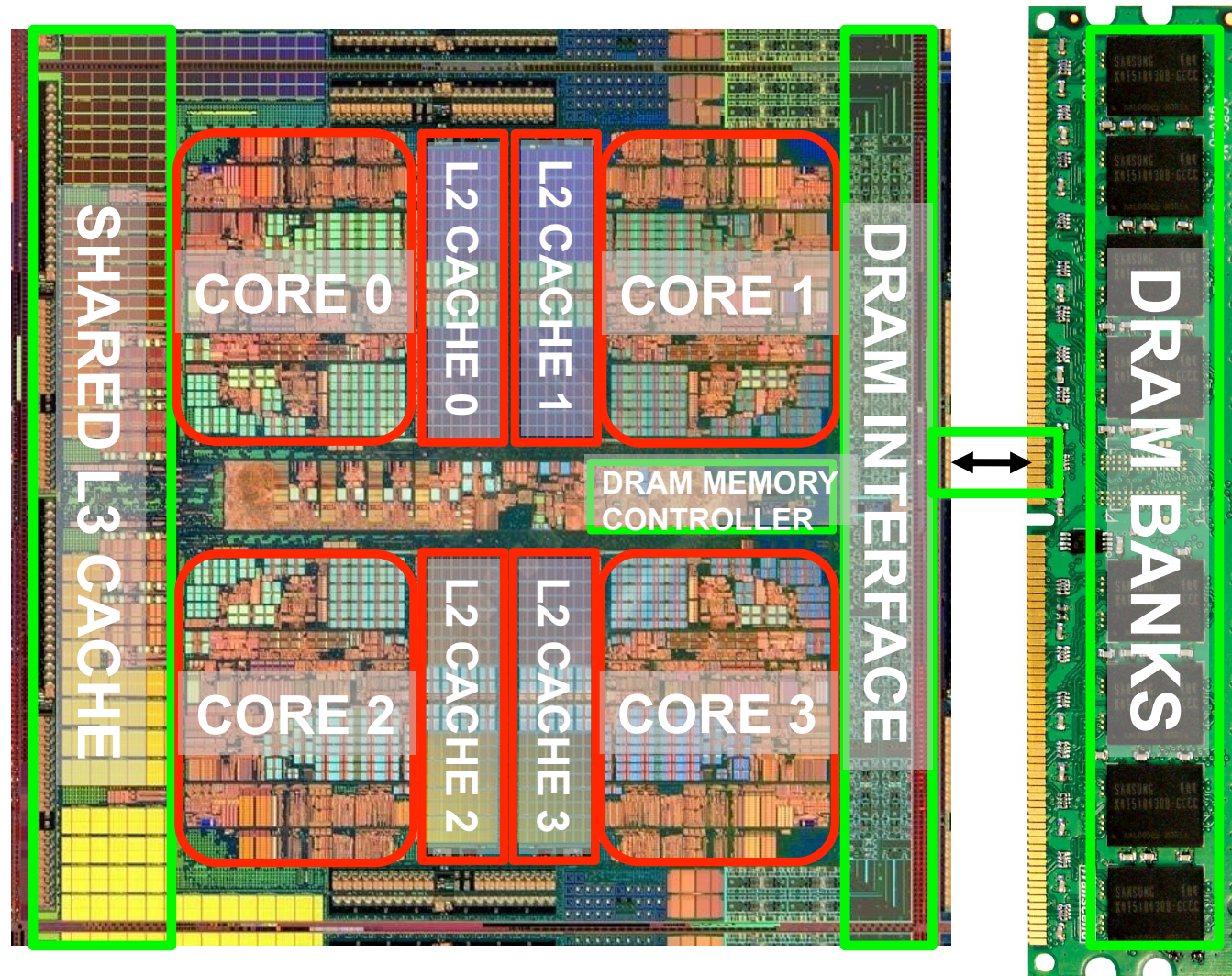
- Infinite capacity

- Infinite bandwidth

- Zero cost

The Memory Hierarchy

Memory in a Modern System



Ideal Memory

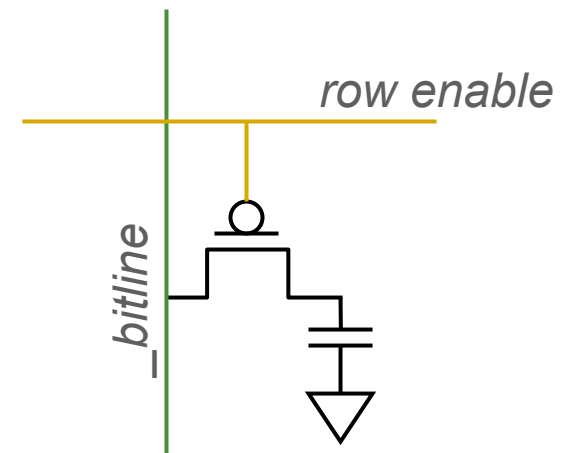
- Zero access time (latency)
- Infinite capacity
- Zero cost
- Infinite bandwidth (to support multiple accesses in parallel)

The Problem

- Ideal memory's requirements oppose each other
- Bigger is slower
 - Bigger → Takes longer to determine the location
- Faster is more expensive
 - Memory technology: SRAM vs. DRAM vs. Disk vs. Tape
- Higher bandwidth is more expensive
 - Need more banks, more ports, higher frequency, or faster technology

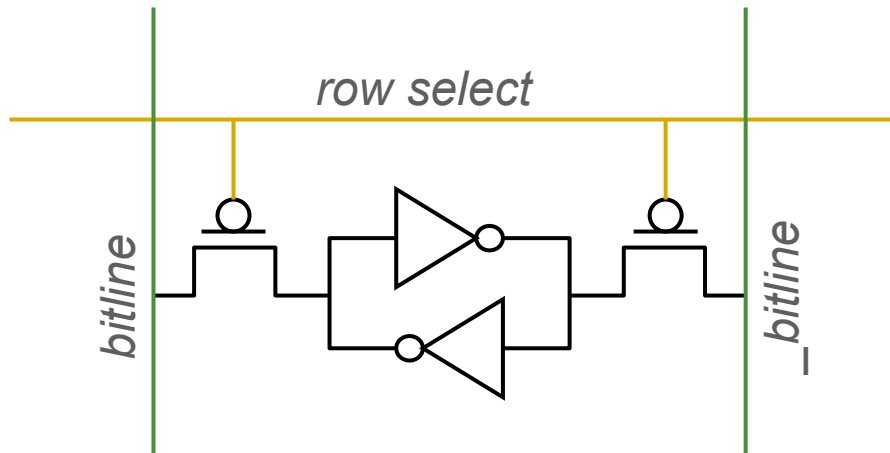
Memory Technology: DRAM

- Dynamic random access memory
- Capacitor charge state indicates stored value
 - Whether the capacitor is charged or discharged indicates storage of 1 or 0
 - 1 capacitor
 - 1 access transistor
- Capacitor leaks through the RC path
 - DRAM cell loses charge over time
 - DRAM cell needs to be refreshed

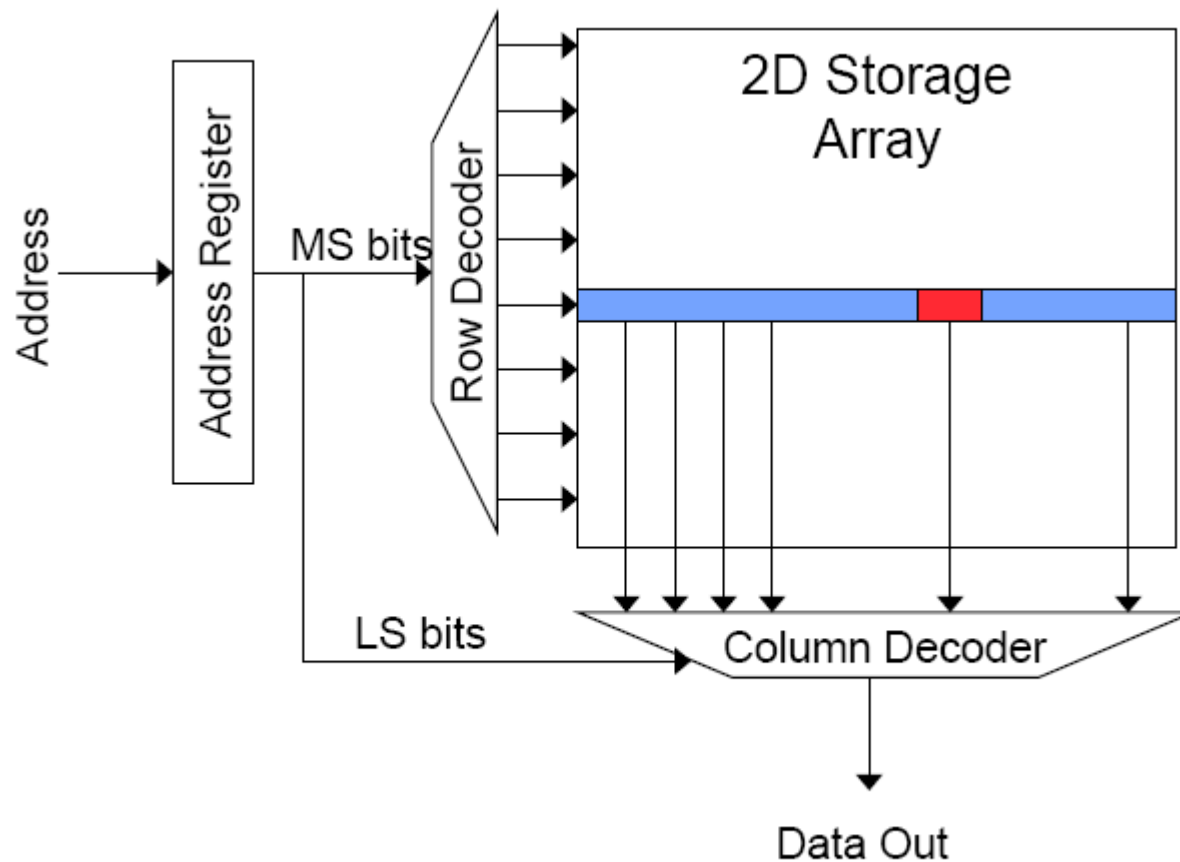


Memory Technology: SRAM

- Static random access memory
- Two cross coupled inverters store a single bit
 - ❑ Feedback path enables the stored value to persist in the “cell”
 - ❑ 4 transistors for storage
 - ❑ 2 transistors for access



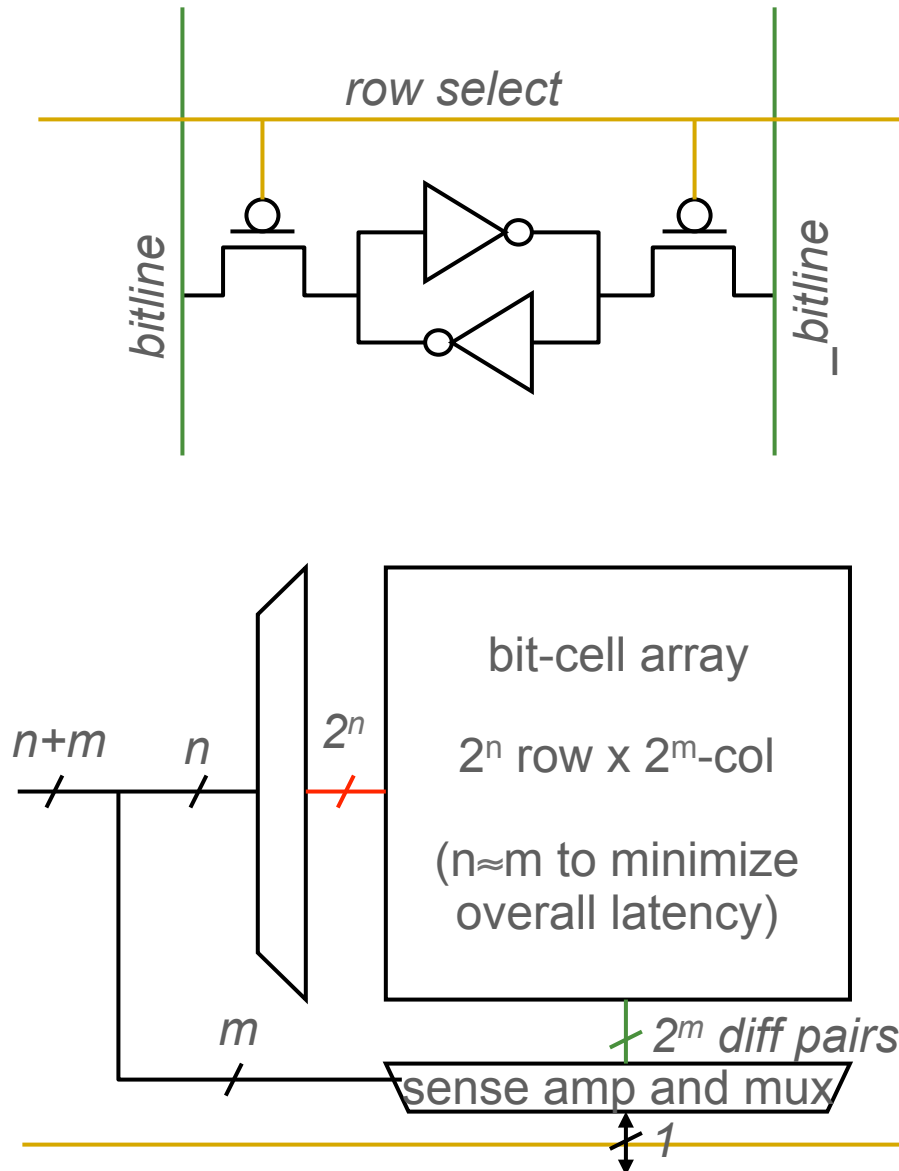
Memory Bank Organization and Operation



■ Read access sequence:

1. Decode row address & drive word-lines
2. Selected bits drive bit-lines
 - Entire row read
3. Amplify row data
4. Decode column address & select subset of row
 - Send to output
5. Precharge bit-lines
 - For next access

SRAM (Static Random Access Memory)



Read Sequence

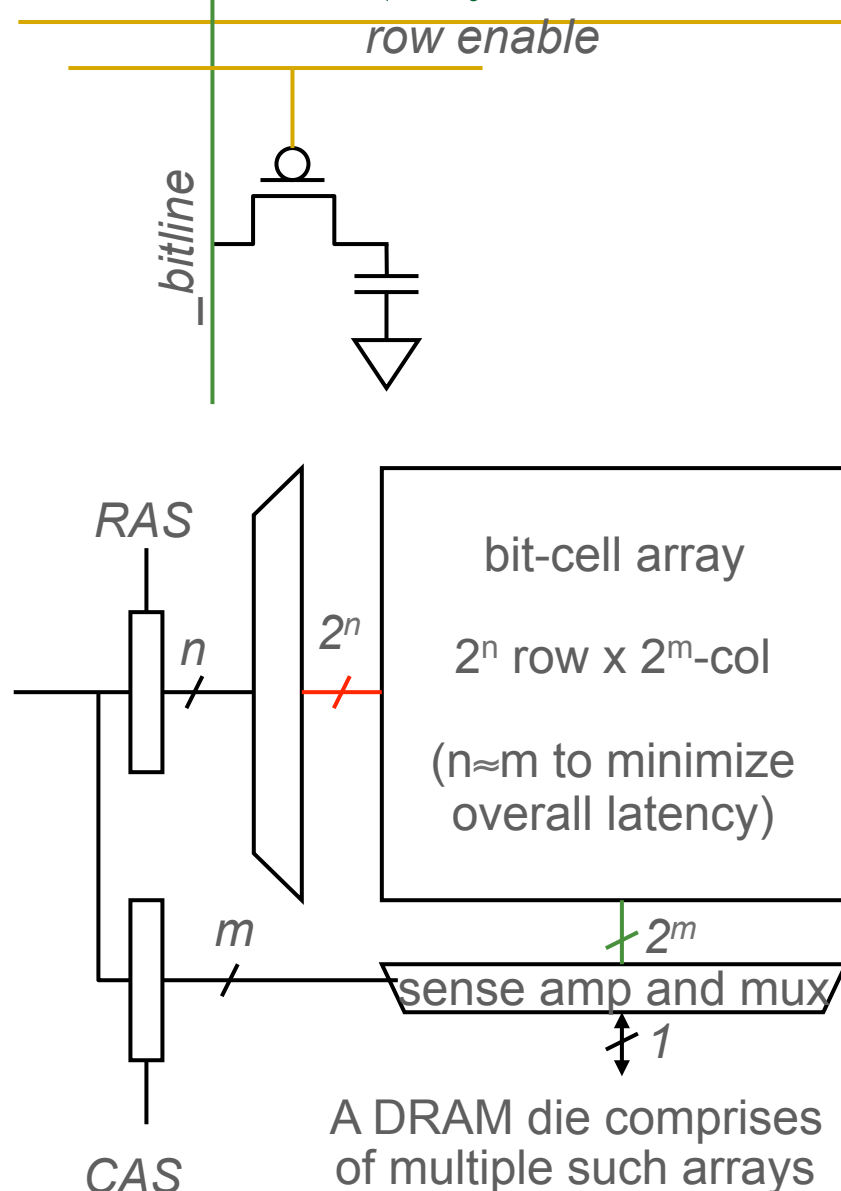
1. address decode
2. drive row select
3. selected bit-cells drive bitlines
(entire row is read together)
4. differential sensing and column select
(data is ready)
5. precharge all bitlines
(for next read or write)

Access latency dominated by steps 2 and 3

Cycling time dominated by steps 2, 3 and 5

- step 2 proportional to 2^m
- step 3 and 5 proportional to 2^n

DRAM (Dynamic Random Access Memory)



Bits stored as charges on node capacitance (non-restorative)

- bit cell loses charge when read
- bit cell loses charge over time

Read Sequence

1~3 same as SRAM

4. a “flip-flopping” sense amp amplifies and regenerates the bitline, data bit is mux’ed out

5. precharge all bitlines

Destructive reads

Charge loss over time

Refresh: A DRAM controller must periodically read each row within the allowed refresh time (10s of ms) such that charge is restored

DRAM vs. SRAM

■ DRAM

- ❑ Slower access (capacitor)
- ❑ Higher density (1T 1C cell)
- ❑ Lower cost
- ❑ Requires refresh (power, performance, circuitry)
- ❑ Manufacturing requires putting capacitor and logic together

■ SRAM

- ❑ Faster access (no capacitor)
- ❑ Lower density (6T cell)
- ❑ Higher cost
- ❑ No need for refresh
- ❑ Manufacturing compatible with logic process (no capacitor)

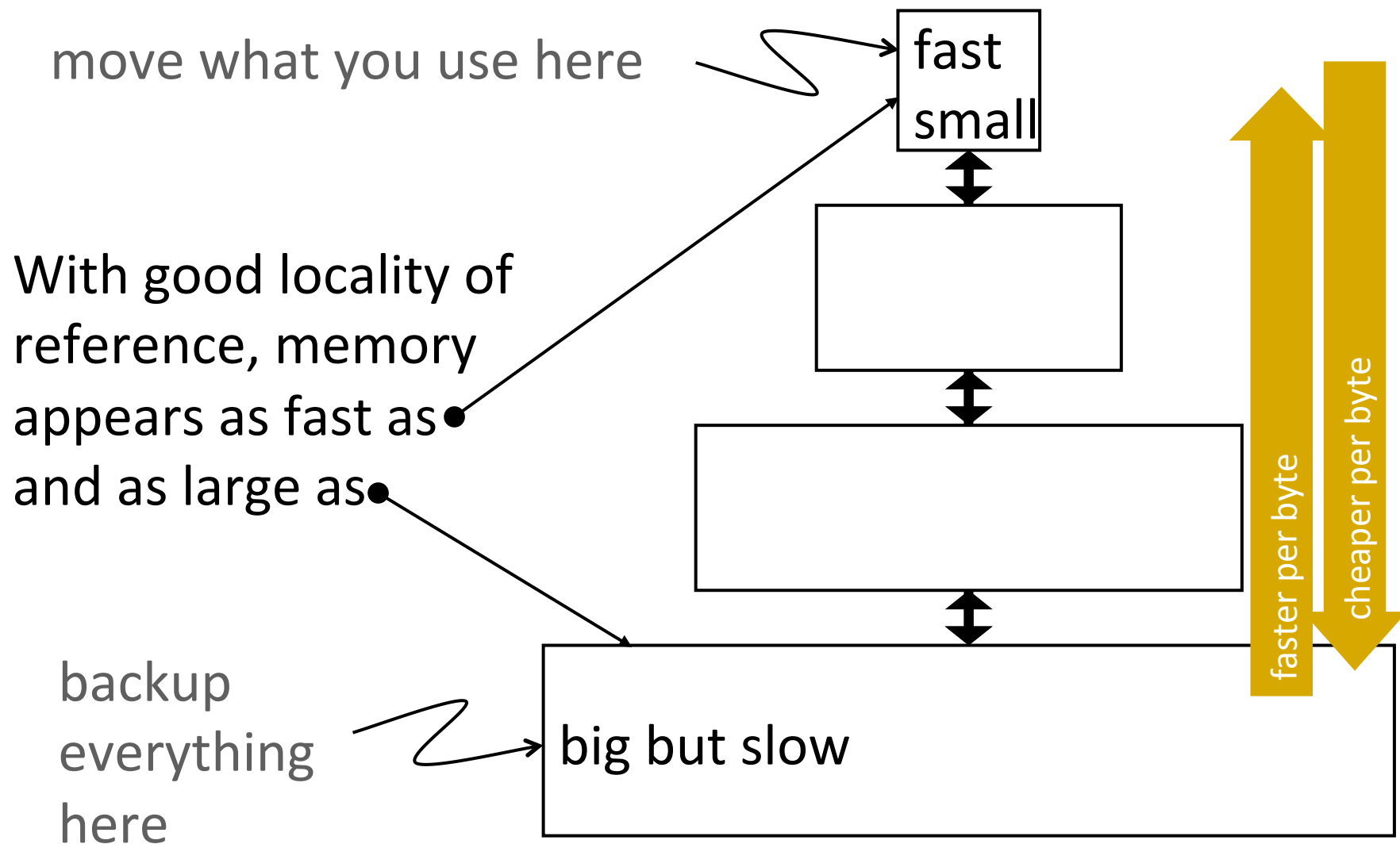
The Problem

- Bigger is slower
 - ❑ SRAM, 512 Bytes, sub-nanosec
 - ❑ SRAM, KByte~MByte, ~nanosec
 - ❑ DRAM, Gigabyte, ~50 nanosec
 - ❑ Hard Disk, Terabyte, ~10 millisec
- Faster is more expensive (dollars and chip area)
 - ❑ SRAM, < 10\$ per Megabyte
 - ❑ DRAM, < 1\$ per Megabyte
 - ❑ Hard Disk < 1\$ per Gigabyte
 - ❑ These sample values (circa ~2011) scale with time
- Other technologies have their place as well
 - ❑ Flash memory, PC-RAM, MRAM, RRAM (not mature yet)

Why Memory Hierarchy?

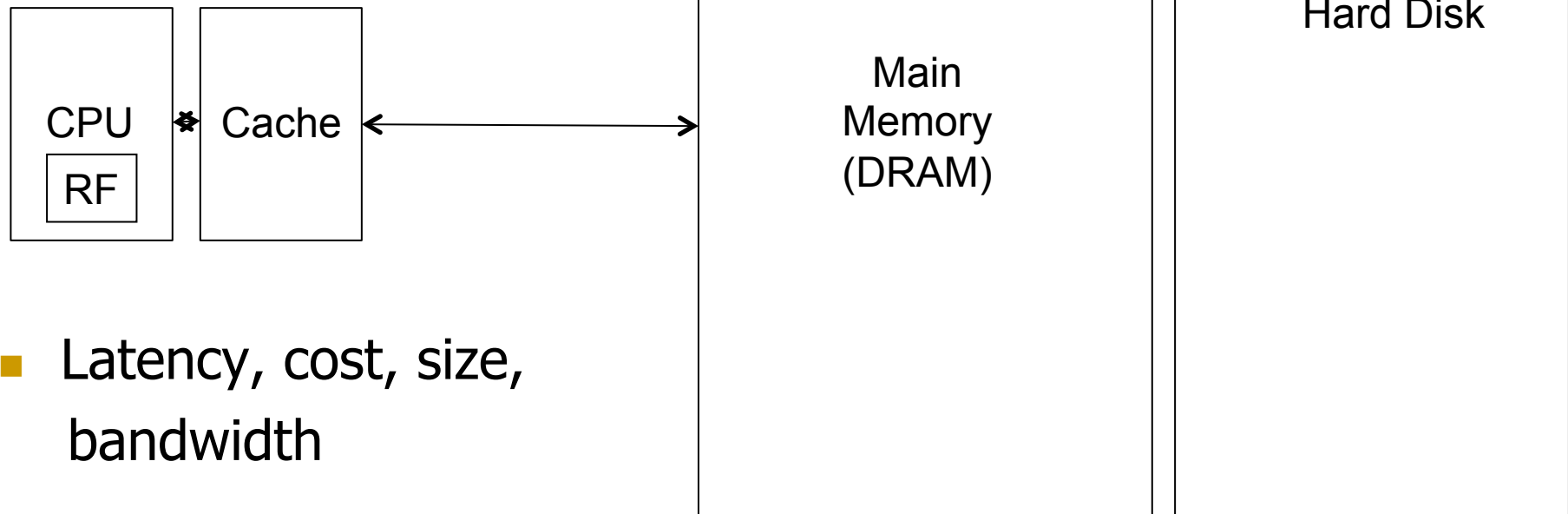
- We want both fast and large
- But we cannot achieve both with a single level of memory
- Idea: Have multiple levels of storage (progressively bigger and slower as the levels are farther from the processor) and ensure most of the data the processor needs is kept in the fast(er) level(s)

The Memory Hierarchy



Memory Hierarchy

- Fundamental tradeoff
 - Fast memory: small
 - Large memory: slow
- Idea: **Memory hierarchy**



- Latency, cost, size, bandwidth

Locality

- One's recent past is a very good predictor of his/her near future.
- **Temporal Locality**: If you just did something, it is very likely that you will do the same thing again soon
 - since you are here today, there is a good chance you will be here again and again regularly
- **Spatial Locality**: If you did something, it is very likely you will do something similar/related (in space)
 - every time I find you in this room, you are probably sitting close to the same people

Memory Locality

- A “typical” program has a lot of locality in memory references
 - typical programs are composed of “loops”
- **Temporal**: A program tends to reference the same memory location many times and all within a small window of time
- **Spatial**: A program tends to reference a cluster of memory locations at a time
 - most notable examples:
 - 1. instruction memory references
 - 2. array/data structure references

Caching Basics: Exploit Temporal Locality

- Idea: Store recently accessed data in automatically managed fast memory (called cache)
- Anticipation: the data will be accessed again soon
- Temporal locality principle
 - Recently accessed data will be again accessed in the near future
 - This is what Maurice Wilkes had in mind:
 - Wilkes, “Slave Memories and Dynamic Storage Allocation,” IEEE Trans. On Electronic Computers, 1965.
 - “The use is discussed of a fast core memory of, say 32000 words as a slave to a slower core memory of, say, one million words in such a way that in practical cases the effective access time is nearer that of the fast memory than that of the slow memory.”

Caching Basics: Exploit Spatial Locality

- Idea: Store addresses adjacent to the recently accessed one in automatically managed fast memory
 - Logically divide memory into equal size blocks
 - Fetch to cache the accessed block in its entirety
- Anticipation: nearby data will be accessed soon
- Spatial locality principle
 - Nearby data in memory will be accessed in the near future
 - E.g., sequential instruction access, array traversal
 - This is what IBM 360/85 implemented
 - 16 Kbyte cache with 64 byte blocks
 - Liptay, “Structural aspects of the System/360 Model 85 II: the cache,” IBM Systems Journal, 1968.

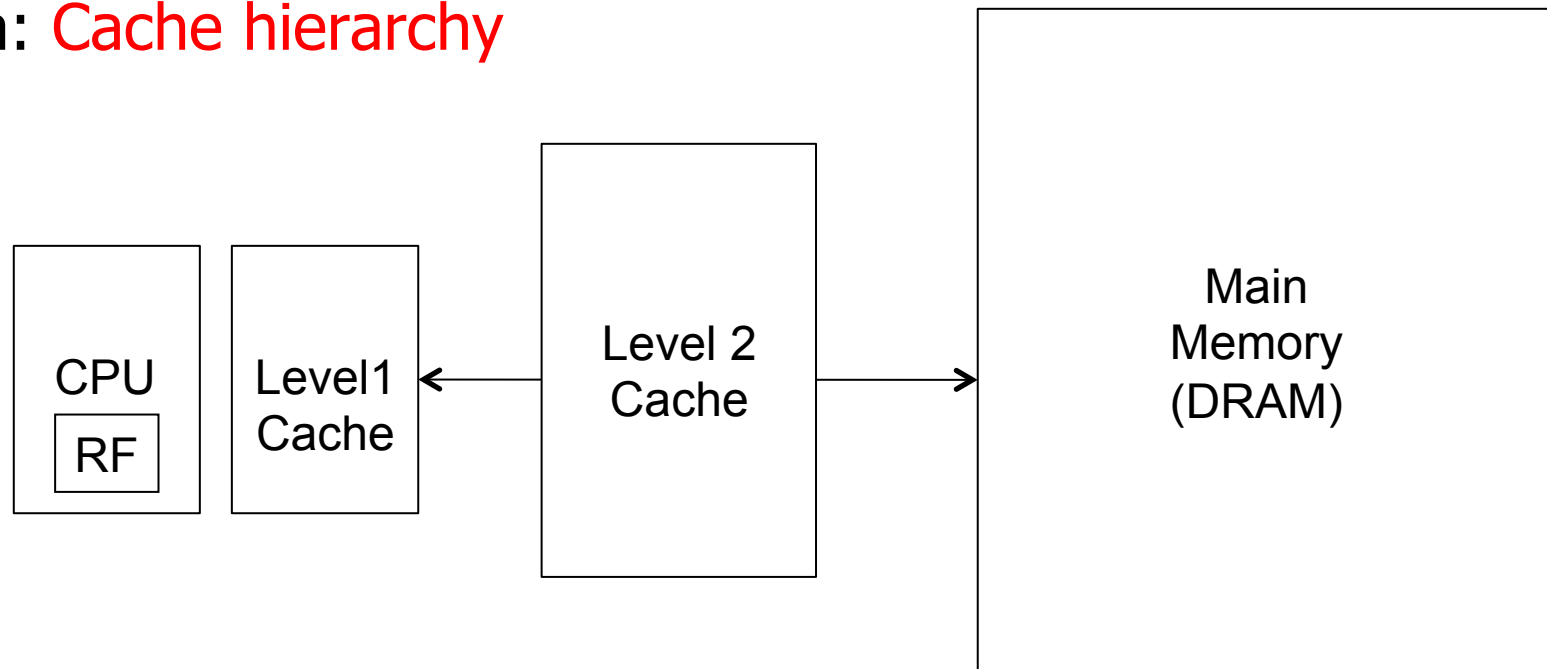
The Bookshelf Analogy

- Book in your hand
- Desk
- Bookshelf
- Boxes at home
- Boxes in storage

- Recently-used books tend to stay on desk
 - Comp Arch books, books for classes you are currently taking
 - Until the desk gets full
- Adjacent books in the shelf needed around the same time
 - If I have organized/categorized my books well in the shelf

Caching in a Pipelined Design

- The cache needs to be tightly integrated into the pipeline
 - Ideally, access in 1-cycle so that dependent operations do not stall
- High frequency pipeline → Cannot make the cache large
 - But, we want a large cache AND a pipelined design
- Idea: **Cache hierarchy**



A Note on Manual vs. Automatic Management

- **Manual:** Programmer manages data movement across levels
 - too painful for programmers on substantial programs
 - “core” vs “drum” memory in the 50’s
 - still done in some embedded processors (on-chip scratch pad SRAM in lieu of a cache) and GPUs (called “shared memory”)

- **Automatic:** Hardware manages data movement across levels, transparently to the programmer
 - ++ programmer’s life is easier
 - the average programmer doesn’t need to know about it
 - You don’t need to know how big the cache is and how it works to write a “correct” program! (What if you want a “fast” program?)

Automatic Management in Memory Hierarchy

- Wilkes, “**Slave Memories and Dynamic Storage Allocation**,” IEEE Trans. On Electronic Computers, 1965.

Slave Memories and Dynamic Storage Allocation

M. V. WILKES

SUMMARY

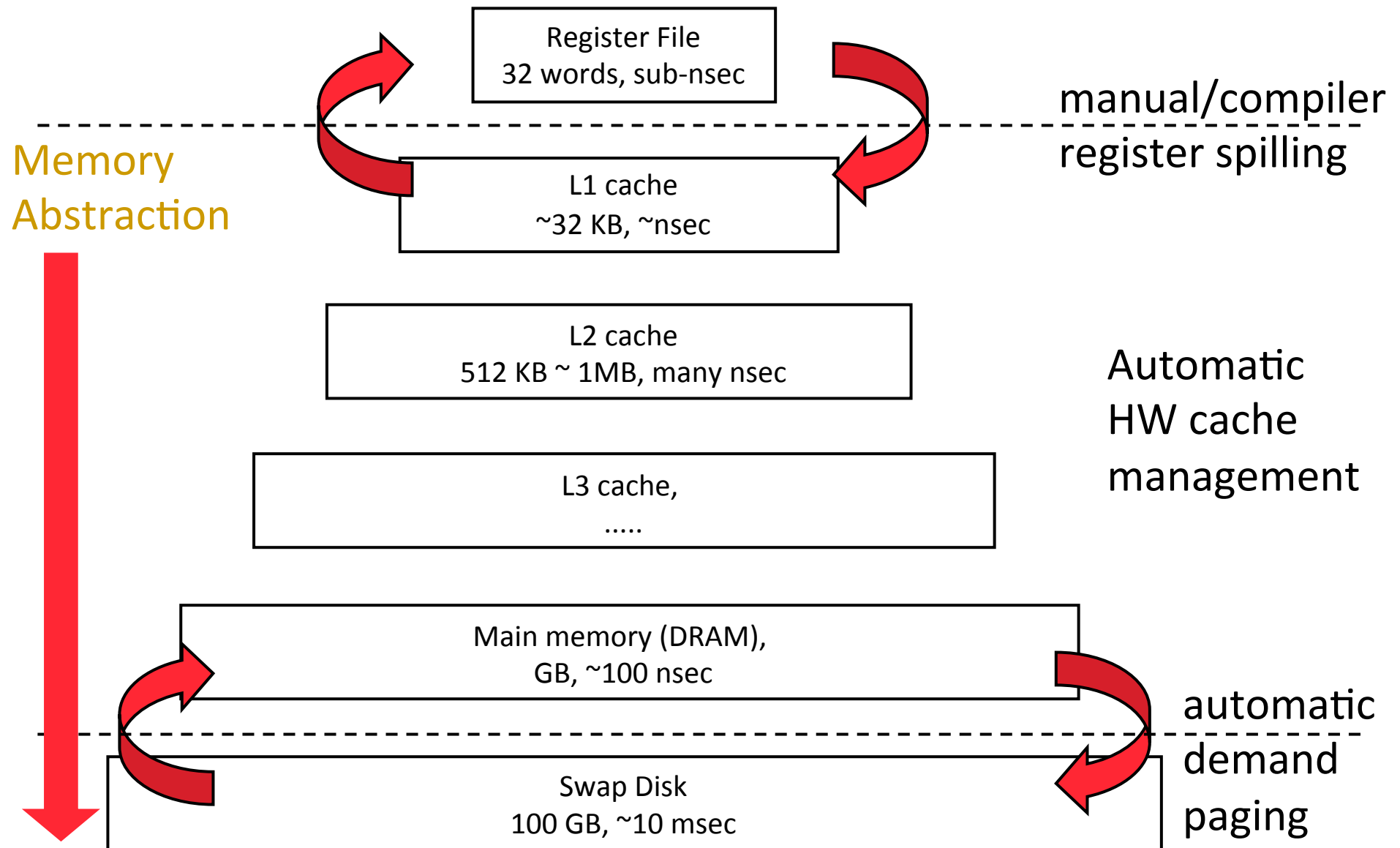
The use is discussed of a fast core memory of, say, 32 000 words as a slave to a slower core memory of, say, one million words in such a way that in practical cases the effective access time is nearer that of the fast memory than that of the slow memory.

- “By a slave memory I mean one which **automatically accumulates to itself words** that come from a slower main memory, and keeps them available for subsequent use without it being necessary for the penalty of main memory access to be incurred again.”

Historical Aside: Other Cache Papers

- Fotheringham, “Dynamic Storage Allocation in the Atlas Computer, Including an Automatic Use of a Backing Store,” CACM 1961.
 - <http://dl.acm.org/citation.cfm?id=366800>
- Bloom, Cohen, Porter, “Considerations in the Design of a Computer with High Logic-to-Memory Speed Ratio,” AIEE Gigacycle Computing Systems Winter Meeting, Jan. 1962.

A Modern Memory Hierarchy



Hierarchical Latency Analysis

- For a given memory hierarchy level i it has a technology-intrinsic access time of t_i . The perceived access time T_i is longer than t_i
- Except for the outer-most hierarchy, when looking for a given address there is
 - a chance (hit-rate h_i) you “hit” and access time is t_i
 - a chance (miss-rate m_i) you “miss” and access time $t_i + T_{i+1}$
 - $h_i + m_i = 1$
- Thus

$$T_i = h_i \cdot t_i + m_i \cdot (t_i + T_{i+1})$$

$$T_i = t_i + m_i \cdot T_{i+1}$$

h_i and m_i are defined to be the hit-rate and miss-rate of just the references that missed at L_{i-1}

Hierarchy Design Considerations

- Recursive latency equation

$$T_i = t_i + m_i \cdot T_{i+1}$$

- The goal: achieve desired T_1 within allowed cost
- $T_i \approx t_i$ is desirable
- Keep m_i low
 - increasing capacity C_i lowers m_i , but beware of increasing t_i
 - lower m_i by smarter management (replacement::anticipate what you don't need, prefetching::anticipate what you will need)
- Keep T_{i+1} low
 - faster lower hierarchies, but beware of increasing cost
 - introduce intermediate hierarchies as a compromise

Intel Pentium 4 Example

- 90nm P4, 3.6 GHz
 - L1 D-cache
 - $C_1 = 16K$
 - $t_1 = 4 \text{ cyc int} / 9 \text{ cycle fp}$
 - L2 D-cache
 - $C_2 = 1024 \text{ KB}$
 - $t_2 = 18 \text{ cyc int} / 18 \text{ cyc fp}$
 - Main memory
 - $t_3 = \sim 50\text{ns or } 180 \text{ cyc}$
 - Notice
 - best case latency is not 1
 - worst case access latencies are into 500+ cycles
- if $m_1=0.1, m_2=0.1$
 $T_1=7.6, T_2=36$

if $m_1=0.01, m_2=0.01$
 $T_1=4.2, T_2=19.8$

if $m_1=0.05, m_2=0.01$
 $T_1=5.00, T_2=19.8$

if $m_1=0.01, m_2=0.50$
 $T_1=5.08, T_2=108$
-

Cache Basics and Operation

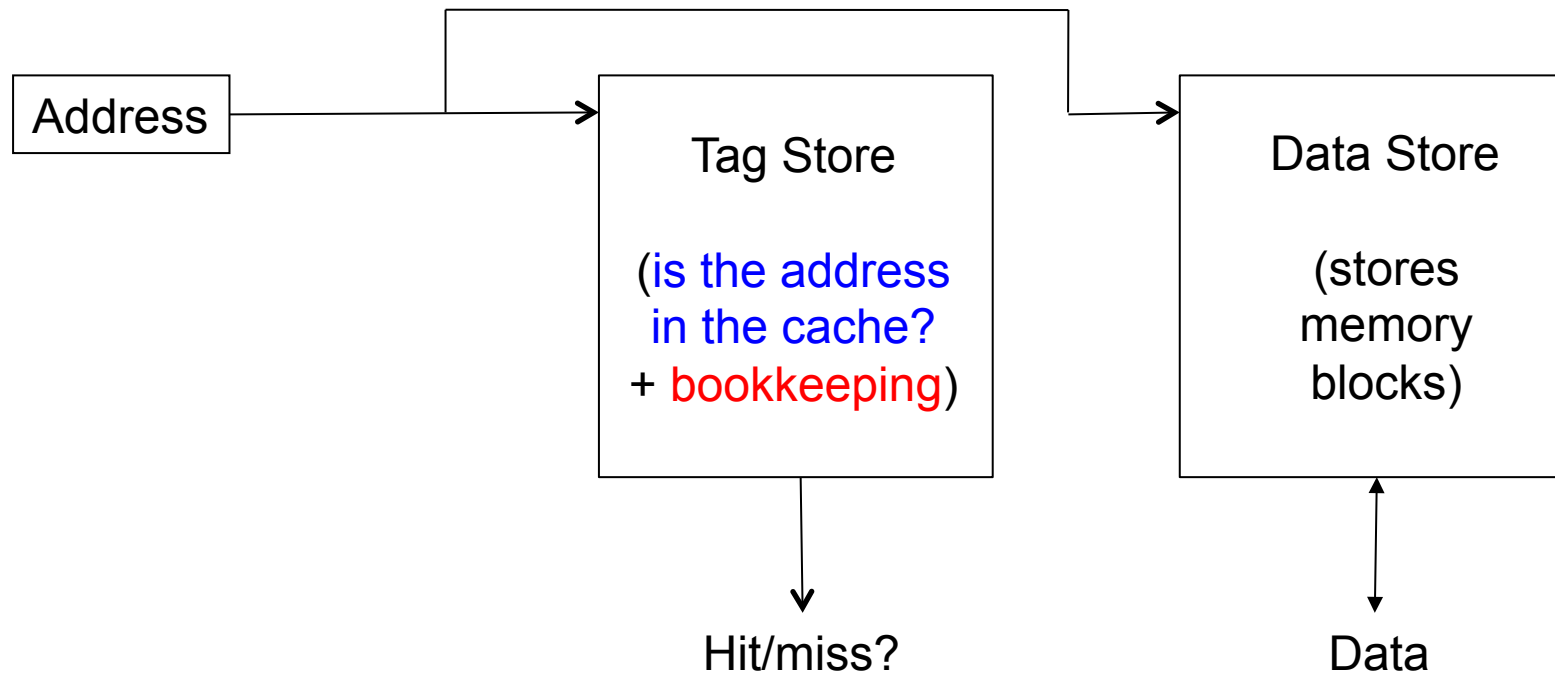
Cache

- Generically, any structure that “memoizes” frequently used results to avoid repeating the long-latency operations required to reproduce the results from scratch, e.g. a web cache
- Most commonly in the on-die context: an automatically-managed memory hierarchy based on SRAM
 - memoize in SRAM the most frequently accessed DRAM memory locations to avoid repeatedly paying for the DRAM access latency

Caching Basics

- **Block (line):** Unit of storage in the cache
 - Memory is logically divided into cache blocks that map to locations in the cache
- On a reference:
 - **HIT:** If in cache, use cached data instead of accessing memory
 - **MISS:** If not in cache, bring block into cache
 - Maybe have to kick something else out to do it
- Some important cache design decisions
 - **Placement:** where and how to place/find a block in cache?
 - **Replacement:** what data to remove to make room in cache?
 - **Granularity of management:** large or small blocks? Subblocks?
 - **Write policy:** what do we do about writes?
 - **Instructions/data:** do we treat them separately?

Cache Abstraction and Metrics



- Cache hit rate = $(\# \text{ hits}) / (\# \text{ hits} + \# \text{ misses}) = (\# \text{ hits}) / (\# \text{ accesses})$
- Average memory access time (AMAT)
= $(\text{hit-rate} * \text{hit-latency}) + (\text{miss-rate} * \text{miss-latency})$
- Aside: *Can reducing AMAT reduce performance?*

A Basic Hardware Cache Design

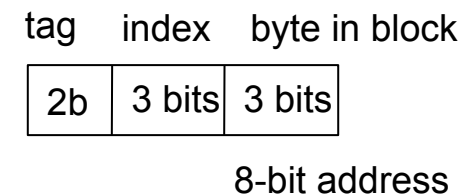
- We will start with a basic hardware cache design
- Then, we will examine a multitude of ideas to make it better

Blocks and Addressing the Cache

- Memory is logically divided into fixed-size blocks

- Each block maps to a location in the cache, determined by the **index bits** in the address

- used to index into the tag and data stores



- Cache access:

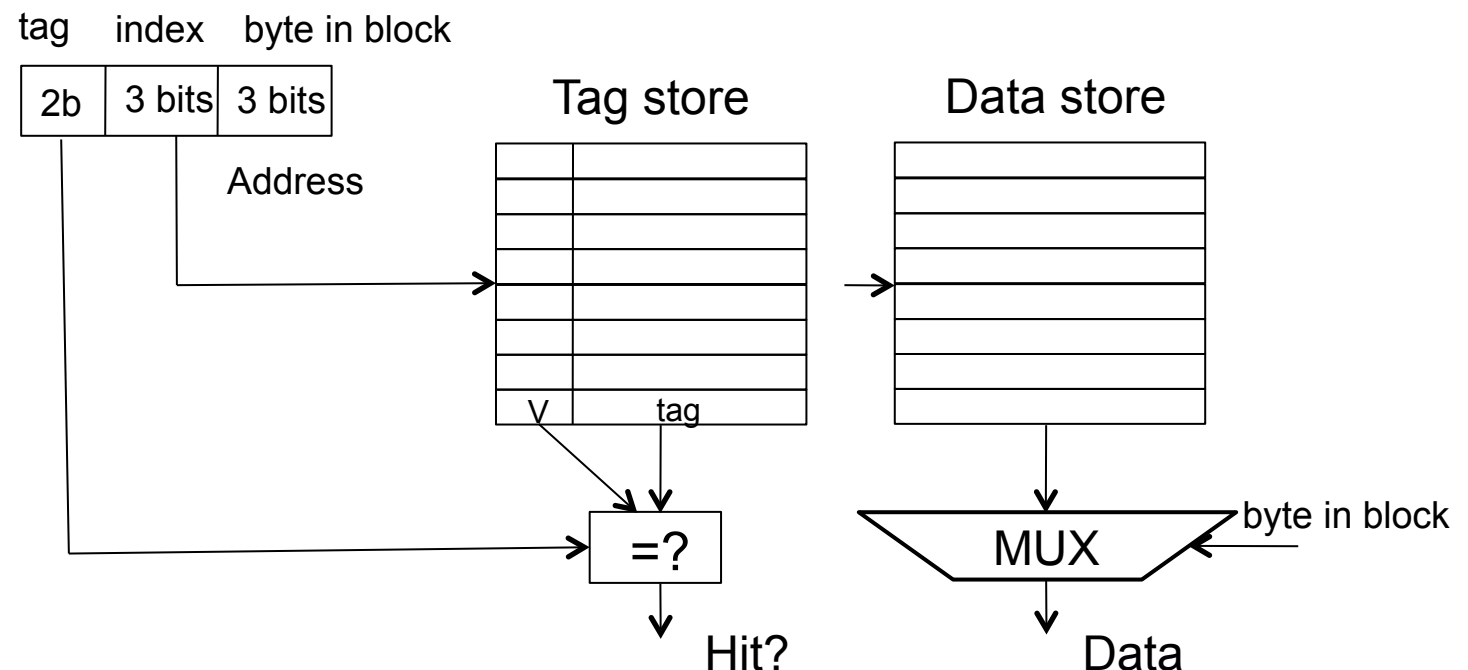
- 1) index into the tag and data stores with index bits in address
- 2) check valid bit in tag store
- 3) compare tag bits in address with the stored tag in tag store

- If a block is in the cache (cache hit), **the stored tag should be valid and match the tag of the block**

Direct-Mapped Cache: Placement and Access



- Assume byte-addressable memory:
256 bytes, 8-byte blocks → 32 blocks
- Assume cache: 64 bytes, 8 blocks
 - Direct-mapped: A block can go to only one location



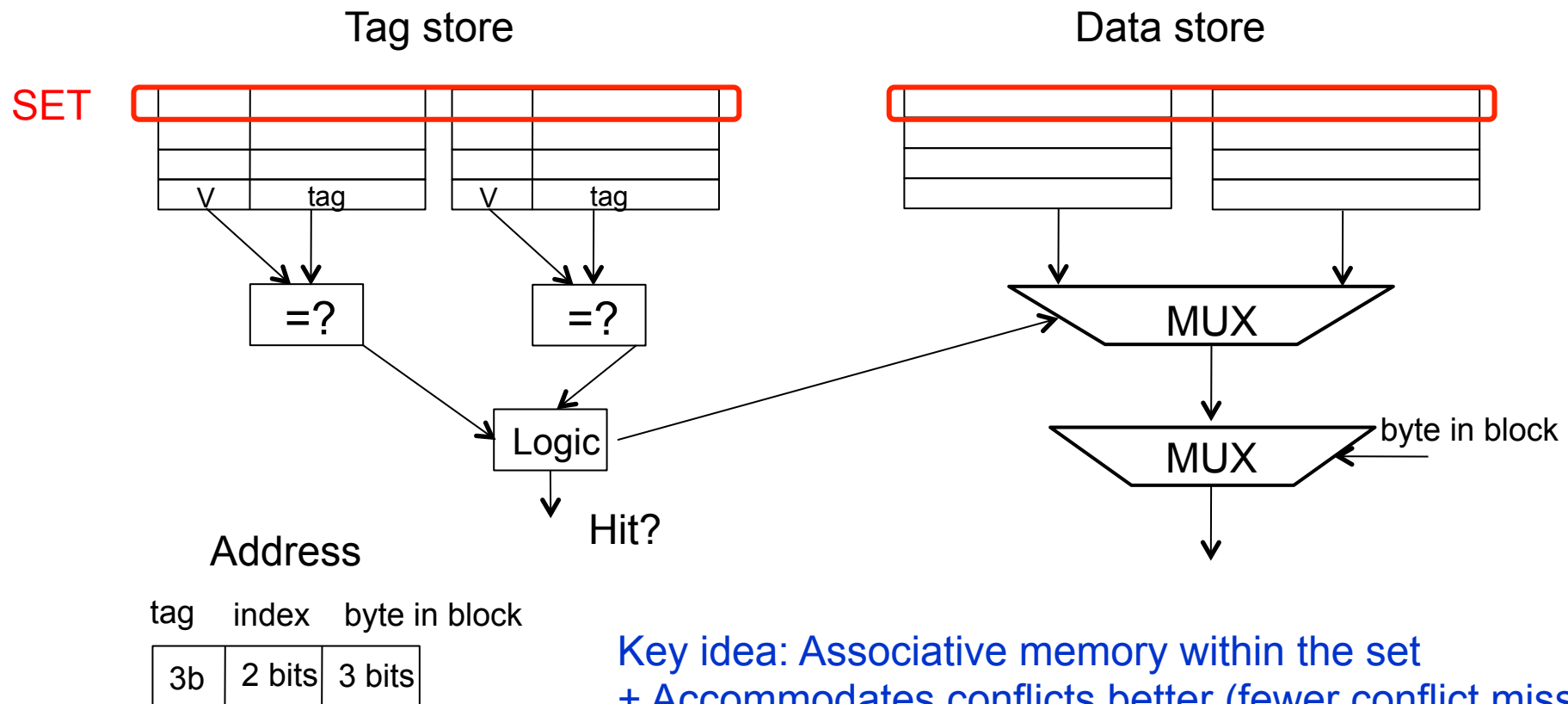
- Addresses with same index contend for the same location
 - Cause conflict misses

Direct-Mapped Caches

- **Direct-mapped cache:** Two blocks in memory that map to the same index in the cache cannot be present in the cache at the same time
 - One index → one entry
- Can lead to 0% hit rate if more than one block accessed in an interleaved manner map to the same index
 - Assume addresses A and B have the same index bits but different tag bits
 - A, B, A, B, A, B, A, B, ... → conflict in the cache index
 - All accesses are **conflict misses**

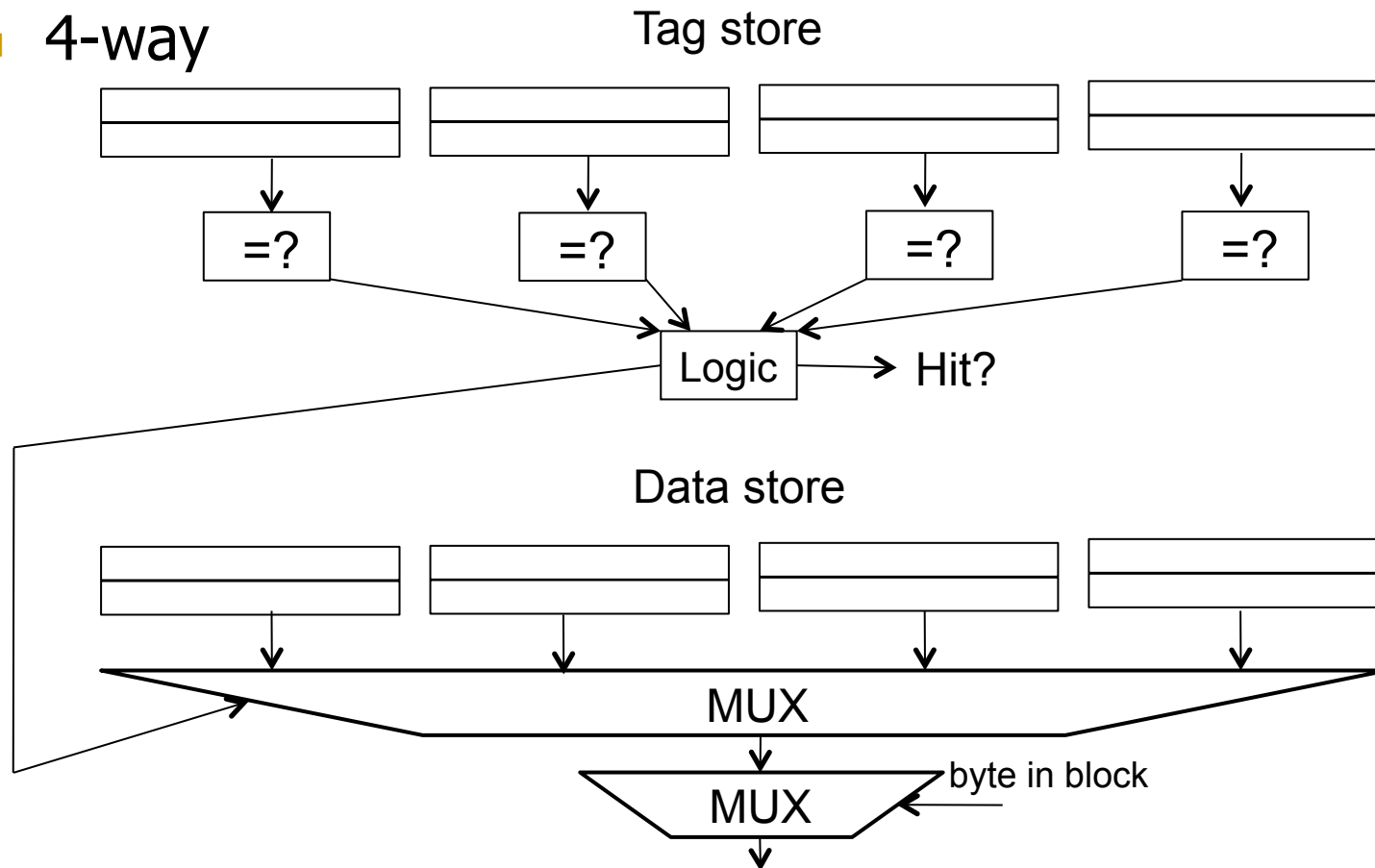
Set Associativity

- Addresses 0 and 8 always conflict in direct mapped cache
- Instead of having one column of 8, have 2 columns of 4 blocks



Higher Associativity

■ 4-way

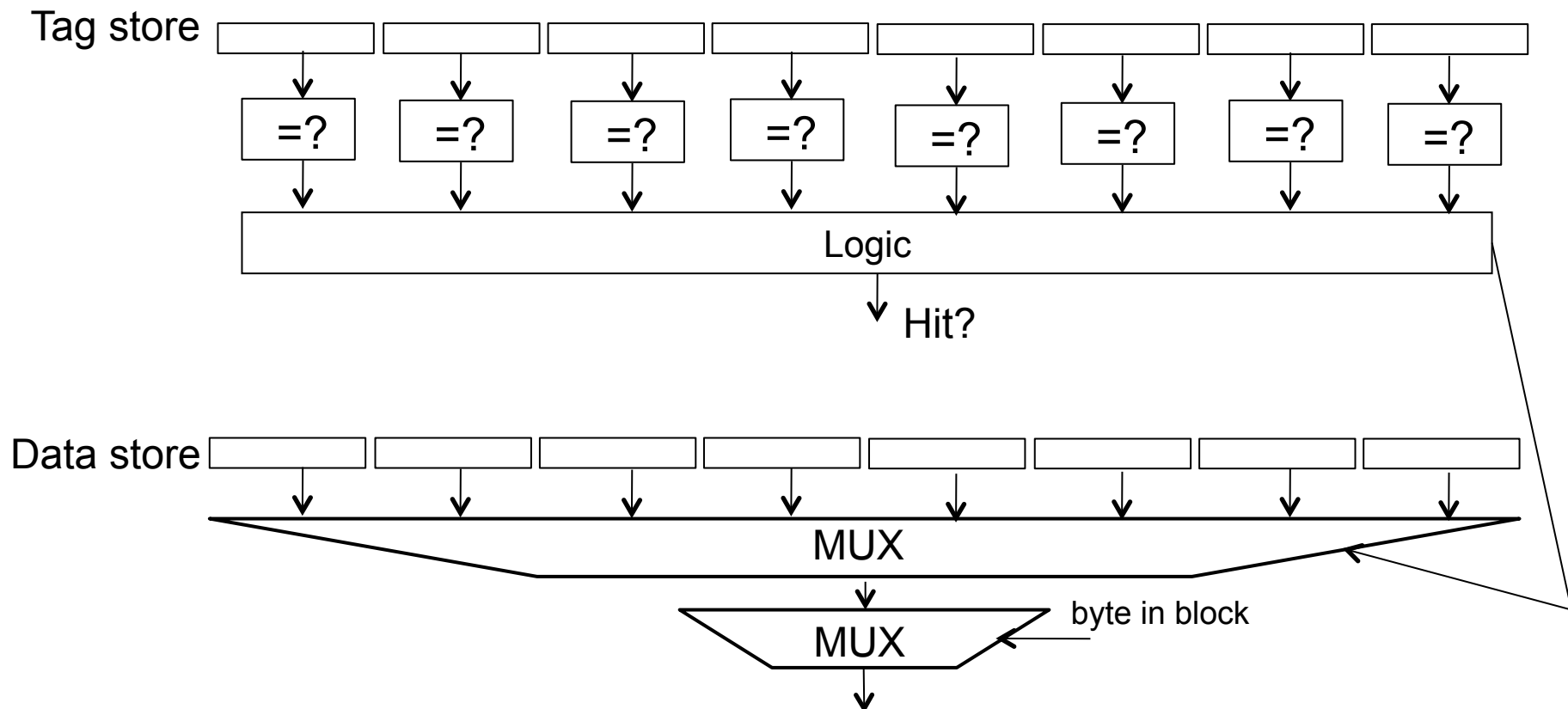


+ Likelihood of conflict misses even lower

-- More tag comparators and wider data mux; larger tags

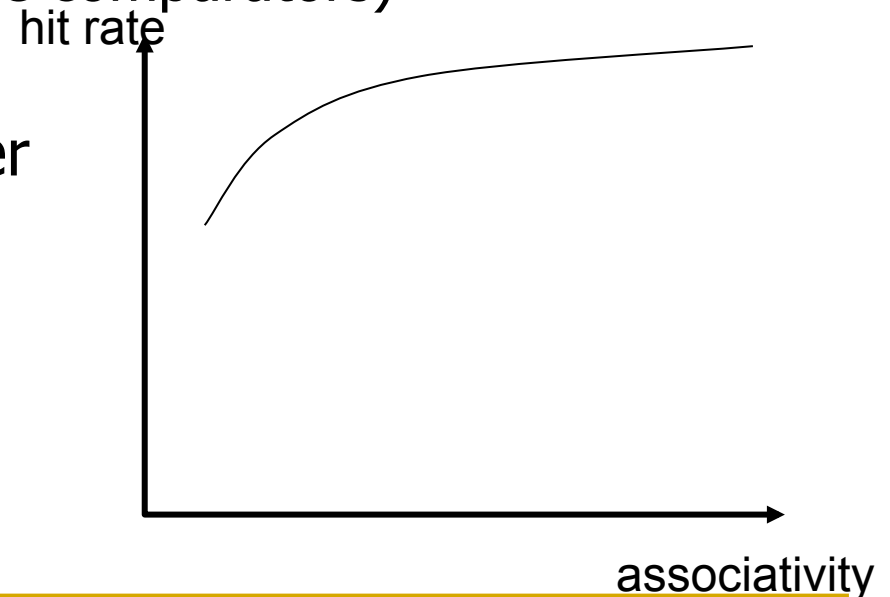
Full Associativity

- Fully associative cache
 - A block can be placed in **any** cache location



Associativity (and Tradeoffs)

- **Degree of associativity**: How many blocks can map to the same index (or set)?
- Higher associativity
 - ++ Higher hit rate
 - Slower cache access time (hit latency and data access latency)
 - More expensive hardware (more comparators)
- Diminishing returns from higher associativity



Computer Architecture

Lecture 2: Fundamentals, Memory Hierarchy, Caches

Prof. Onur Mutlu

ETH Zurich

Fall 2017

21 September 2017

We did not cover the following slides in lecture.
These are for your preparation for the next lecture.

Issues in Set-Associative Caches

- Think of each block in a set having a “priority”
 - Indicating how important it is to keep the block in the cache
- Key issue: How do you determine/adjust block priorities?
- There are three key decisions in a set:
 - Insertion, promotion, eviction (replacement)
- Insertion: What happens to priorities on a cache fill?
 - Where to insert the incoming block, whether or not to insert the block
- Promotion: What happens to priorities on a cache hit?
 - Whether and how to change block priority
- Eviction/replacement: What happens to priorities on a cache miss?
 - Which block to evict and how to adjust priorities

Eviction/Replacement Policy

- Which block in the set to replace on a cache miss?
 - Any invalid block first
 - If all are valid, consult the replacement policy
 - Random
 - FIFO
 - Least recently used (how to implement?)
 - Not most recently used
 - Least frequently used?
 - Least costly to re-fetch?
 - Why would memory accesses have different cost?
 - Hybrid replacement policies
 - Optimal replacement policy?

Implementing LRU

- Idea: Evict the least recently accessed block
- Problem: Need to keep track of access ordering of blocks

- Question: 2-way set associative cache:
 - What do you need to implement LRU perfectly?

- Question: 4-way set associative cache:
 - What do you need to implement LRU perfectly?
 - How many different orderings possible for the 4 blocks in the set?
 - How many bits needed to encode the LRU order of a block?
 - What is the logic needed to determine the LRU victim?

Approximations of LRU

- Most modern processors do not implement “true LRU” (also called “perfect LRU”) in highly-associative caches
- Why?
 - True LRU is complex
 - LRU is an approximation to predict locality anyway (i.e., not the best possible cache management policy)
- Examples:
 - Not MRU (not most recently used)
 - Hierarchical LRU: divide the N-way set into M “groups”, track the MRU group and the MRU way in each group
 - Victim-NextVictim Replacement: Only keep track of the victim and the next victim

Hierarchical LRU (not MRU)

- Divide a set into multiple groups
- Keep track of *only* the MRU group
- Keep track of *only* the MRU block in each group
- On replacement, select victim as:
 - A not-MRU block in one of the not-MRU groups (randomly pick one of such blocks/groups)

Hierarchical LRU (not MRU): Questions

- 16-way cache
- 2 8-way groups
- What is an access pattern that performs worse than true LRU?
- What is an access pattern that performs better than true LRU?

Victim/Next-Victim Policy

- Only 2 blocks' status tracked in each set:
 - victim (V), next victim (NV)
 - all other blocks denoted as (O) – Ordinary block

- On a cache miss
 - Replace V
 - Demote NV to V
 - Randomly pick an O block as NV

- On a cache hit to V
 - Demote NV to V
 - Randomly pick an O block as NV
 - Turn V to O

Victim/Next-Victim Policy (II)

- On a cache hit to NV
 - Randomly pick an O block as NV
 - Turn NV to O

- On a cache hit to O
 - Do nothing

Victim/Next-Victim Example

Example

	V	NV		V	NV
A	1	0		0	0
B	0	0		0	0
C	0	1		1	0
D	0	0		0	1

hit to A

randomly picked

time

Some questions as before

Cache Replacement Policy: LRU or Random

- LRU vs. Random: Which one is better?
 - Example: 4-way cache, cyclic references to A, B, C, D, E
 - 0% hit rate with LRU policy
- **Set thrashing:** When the “program working set” in a set is larger than set associativity
 - Random replacement policy is better when thrashing occurs
- In practice:
 - Depends on workload
 - Average hit rate of LRU and Random are similar
- Best of both Worlds: Hybrid of LRU and Random
 - How to choose between the two? **Set sampling**
 - See Qureshi et al., “**A Case for MLP-Aware Cache Replacement**,” ISCA 2006.

What Is the Optimal Replacement Policy?

- Belady's OPT
 - ❑ Replace the block that is going to be referenced furthest in the future by the program
 - ❑ Belady, “A study of replacement algorithms for a virtual-storage computer,” IBM Systems Journal, 1966.
 - ❑ How do we implement this? Simulate?

- Is this optimal for minimizing miss rate?
- Is this optimal for minimizing execution time?
 - ❑ No. Cache miss latency/cost varies from block to block!
 - ❑ Two reasons: Remote vs. local caches and miss overlapping
 - ❑ Qureshi et al. “A Case for MLP-Aware Cache Replacement,” ISCA 2006.

Aside: Cache versus Page Replacement

- Physical memory (DRAM) is a cache for disk
 - Usually managed by system software via the virtual memory subsystem
- Page replacement is similar to cache replacement
- Page table is the “tag store” for physical memory data store
- What is the difference?
 - Required speed of access to cache vs. physical memory
 - Number of blocks in a cache vs. physical memory
 - “Tolerable” amount of time to find a replacement candidate (disk versus memory access latency)
 - Role of hardware versus software

What's In A Tag Store Entry?

- Valid bit
- Tag
- Replacement policy bits
- Dirty bit?
 - Write back vs. write through caches

Handling Writes (I)

- When do we write the modified data in a cache to the next level?
 - Write through: At the time the write happens
 - Write back: When the block is evicted
- Write-back
 - + Can combine multiple writes to the same block before eviction
 - Potentially saves bandwidth between cache levels + saves energy
 - Need a bit in the tag store indicating the block is “dirty/modified”
- Write-through
 - + Simpler
 - + All levels are up to date. Consistency: Simpler cache coherence because no need to check lower-level caches
 - More bandwidth intensive; no combining of writes

Handling Writes (II)

- Do we allocate a cache block on a write miss?
 - Allocate on write miss: Yes
 - No-allocate on write miss: No

- Allocate on write miss
 - + Can combine writes instead of writing each of them individually to next level
 - + Simpler because write misses can be treated the same way as read misses
 - Requires (?) transfer of the whole cache block

- No-allocate
 - + Conserves cache space if locality of writes is low (potentially better cache hit rate)

Handling Writes (III)

- What if the processor writes to an entire block over a small amount of time?
- Is there any need to bring the block into the cache from memory in the first place?
- Ditto for a *portion* of the block, i.e., subblock
 - E.g., 4 bytes out of 64 bytes

Sectored Caches

- Idea: Divide a block into subblocks (or sectors)
 - Have separate valid and dirty bits for each sector
 - When is this useful? (Think writes...)
- ++ No need to transfer the entire cache block into the cache
(A write simply validates and updates a subblock)
- ++ More freedom in transferring subblocks into the cache (a cache block does not need to be in the cache fully)
(How many subblocks do you transfer on a read?)
- More complex design
- May not exploit spatial locality fully when used for reads



Instruction vs. Data Caches

- Separate or Unified?
- Unified:
 - + Dynamic sharing of cache space: no overprovisioning that might happen with static partitioning (i.e., split I and D caches)
 - Instructions and data can thrash each other (i.e., no guaranteed space for either)
 - I and D are accessed in different places in the pipeline. Where do we place the unified cache for fast access?
- First level caches are almost always split
 - Mainly for the last reason above
- Second and higher levels are almost always unified

Multi-level Caching in a Pipelined Design

- First-level caches (instruction and data)
 - Decisions very much affected by cycle time
 - Small, lower associativity
 - Tag store and data store accessed in parallel
- Second-level caches
 - Decisions need to balance hit rate and access latency
 - Usually large and highly associative; latency not as important
 - Tag store and data store accessed serially
- Serial vs. Parallel access of levels
 - Serial: Second level cache accessed only if first-level misses
 - Second level does not see the same accesses as the first
 - First level acts as a filter (filters some temporal and spatial locality)
 - Management policies are therefore different

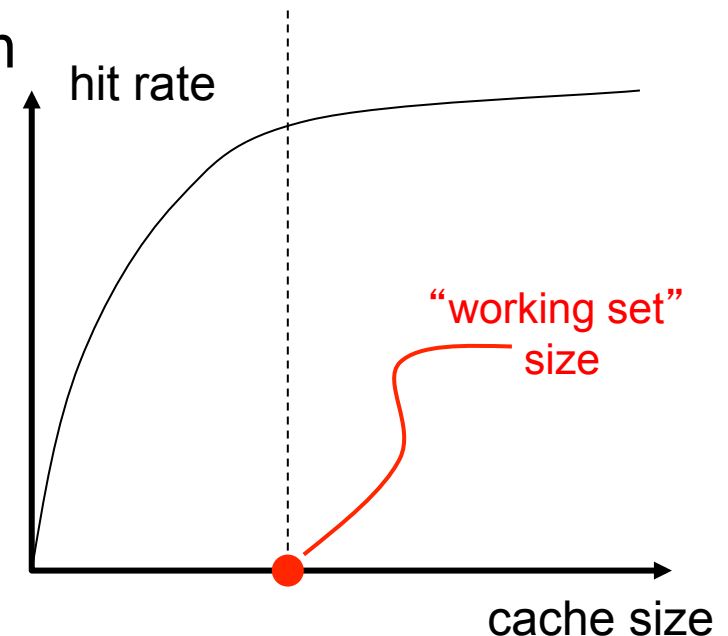
Cache Performance

Cache Parameters vs. Miss/Hit Rate

- Cache size
- Block size
- Associativity
- Replacement policy
- Insertion/Placement policy

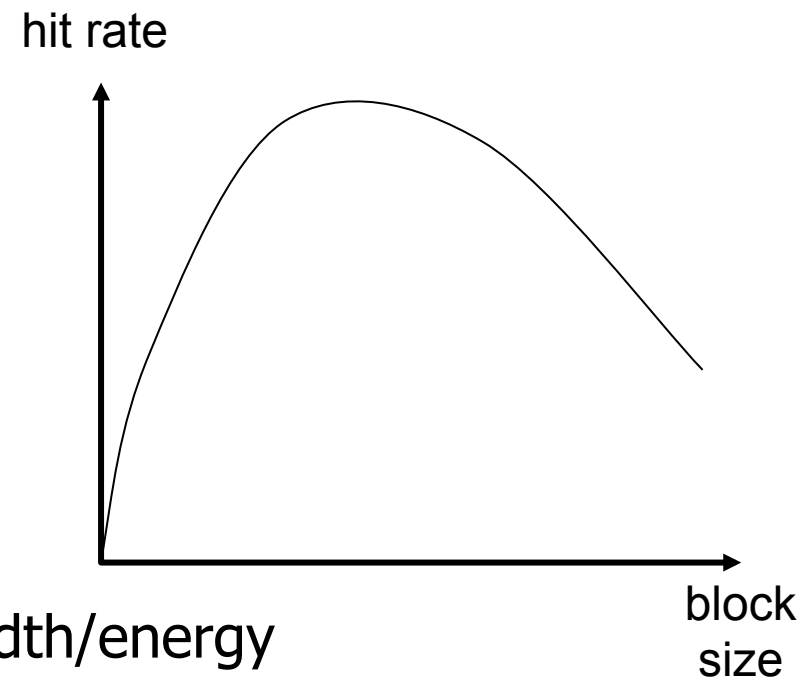
Cache Size

- Cache size: total data (not including tag) capacity
 - bigger can exploit temporal locality better
 - not ALWAYS better
- Too large a cache adversely affects hit and miss latency
 - smaller is faster => bigger is slower
 - access time may degrade critical path
- Too small a cache
 - doesn't exploit temporal locality well
 - useful data replaced often
- **Working set**: the whole set of data the executing application references
 - Within a time interval



Block Size

- Block size is the data that is associated with an address tag
 - not necessarily the unit of transfer between hierarchies
 - Sub-blocking: A block divided into multiple pieces (each with V bit)
 - Can improve “write” performance
- Too small blocks
 - don't exploit spatial locality well
 - have larger tag overhead
- Too large blocks
 - too few total # of blocks → less temporal locality exploitation
 - waste of cache space and bandwidth/energy if spatial locality is not high



Large Blocks: Critical-Word and Subblocking

- Large cache blocks can take a long time to fill into the cache
 - ❑ fill cache line **critical word first**
 - ❑ restart cache access before complete fill
- Large cache blocks can waste bus bandwidth
 - ❑ divide a block into subblocks
 - ❑ associate separate valid bits for each subblock
 - ❑ **When is this useful?**

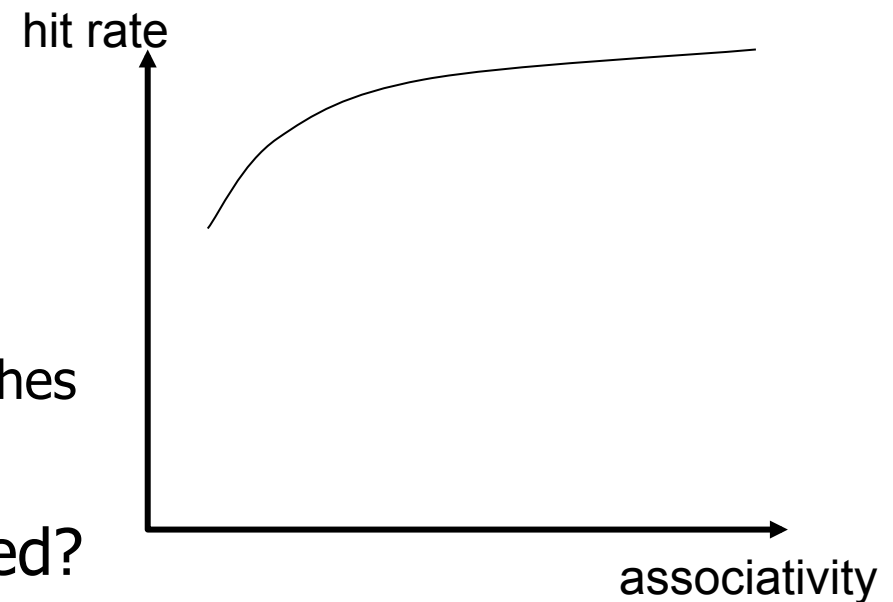


Associativity

- How many blocks can map to the same index (or set)?
- Larger associativity
 - ❑ lower miss rate (reduced conflicts)
 - ❑ higher hit latency and area cost (plus diminishing returns)

- Smaller associativity
 - ❑ lower cost
 - ❑ lower hit latency
 - Especially important for L1 caches

- Power of 2 associativity required?



Classification of Cache Misses

- Compulsory miss
 - ❑ first reference to an address (block) always results in a miss
 - ❑ subsequent references should hit unless the cache block is displaced for the reasons below
- Capacity miss
 - ❑ cache is too small to hold everything needed
 - ❑ defined as the misses that would occur even in a fully-associative cache (with optimal replacement) of the same capacity
- Conflict miss
 - ❑ defined as any miss that is neither a compulsory nor a capacity miss

How to Reduce Each Miss Type

- Compulsory
 - ❑ Caching cannot help
 - ❑ Prefetching can
- Conflict
 - ❑ More associativity
 - ❑ Other ways to get more associativity without making the cache associative
 - Victim cache
 - Better, randomized indexing
 - Software hints?
- Capacity
 - ❑ Utilize cache space better: keep blocks that will be referenced
 - ❑ Software management: divide working set such that each “phase” fits in cache

How to Improve Cache Performance

- Three fundamental goals
- Reducing miss rate
 - Caveat: reducing miss rate can reduce performance if more costly-to-refetch blocks are evicted
- Reducing miss latency or miss cost
- Reducing hit latency or hit cost

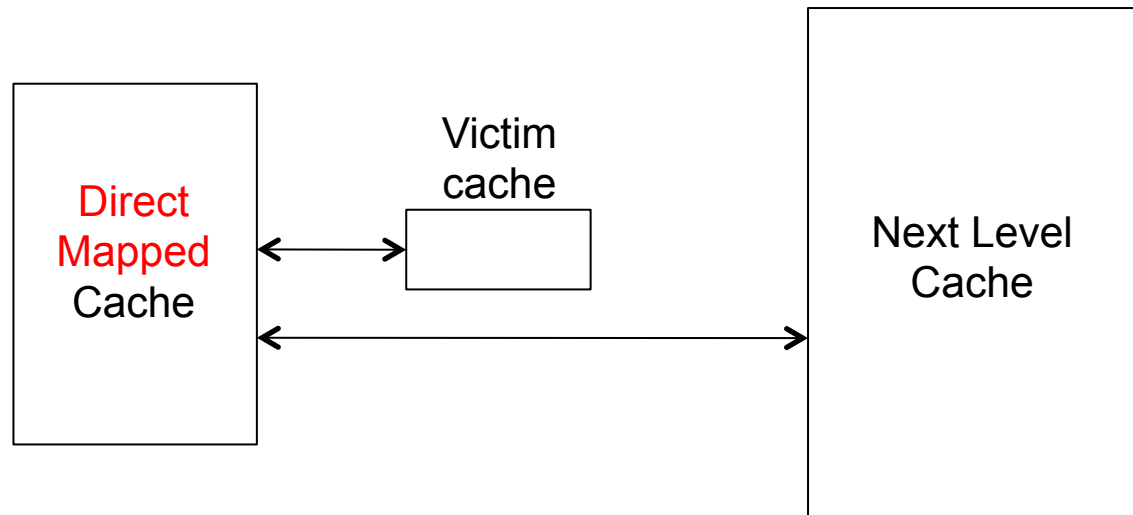
Improving Basic Cache Performance

- Reducing miss rate
 - More associativity
 - Alternatives/enhancements to associativity
 - Victim caches, hashing, pseudo-associativity, skewed associativity
 - Better replacement/insertion policies
 - Software approaches
- Reducing miss latency/cost
 - Multi-level caches
 - Critical word first
 - Subblocking/sectoring
 - Better replacement/insertion policies
 - Non-blocking caches (multiple cache misses in parallel)
 - Multiple accesses per cycle
 - Software approaches

Cheap Ways of Reducing Conflict Misses

- Instead of building highly-associative caches:
 - Victim Caches
 - Hashed/randomized Index Functions
 - Pseudo Associativity
 - Skewed Associative Caches
 - ...

Victim Cache: Reducing Conflict Misses



- Jouppi, “Improving Direct-Mapped Cache Performance by the Addition of a Small Fully-Associative Cache and Prefetch Buffers,” ISCA 1990.
- Idea: Use a small fully-associative buffer (victim cache) to store recently evicted blocks
 - + Can avoid ping ponging of cache blocks mapped to the same set (if two cache blocks continuously accessed in nearby time conflict with each other)
 - Increases miss latency if accessed serially with L2; adds complexity

Hashing and Pseudo-Associativity

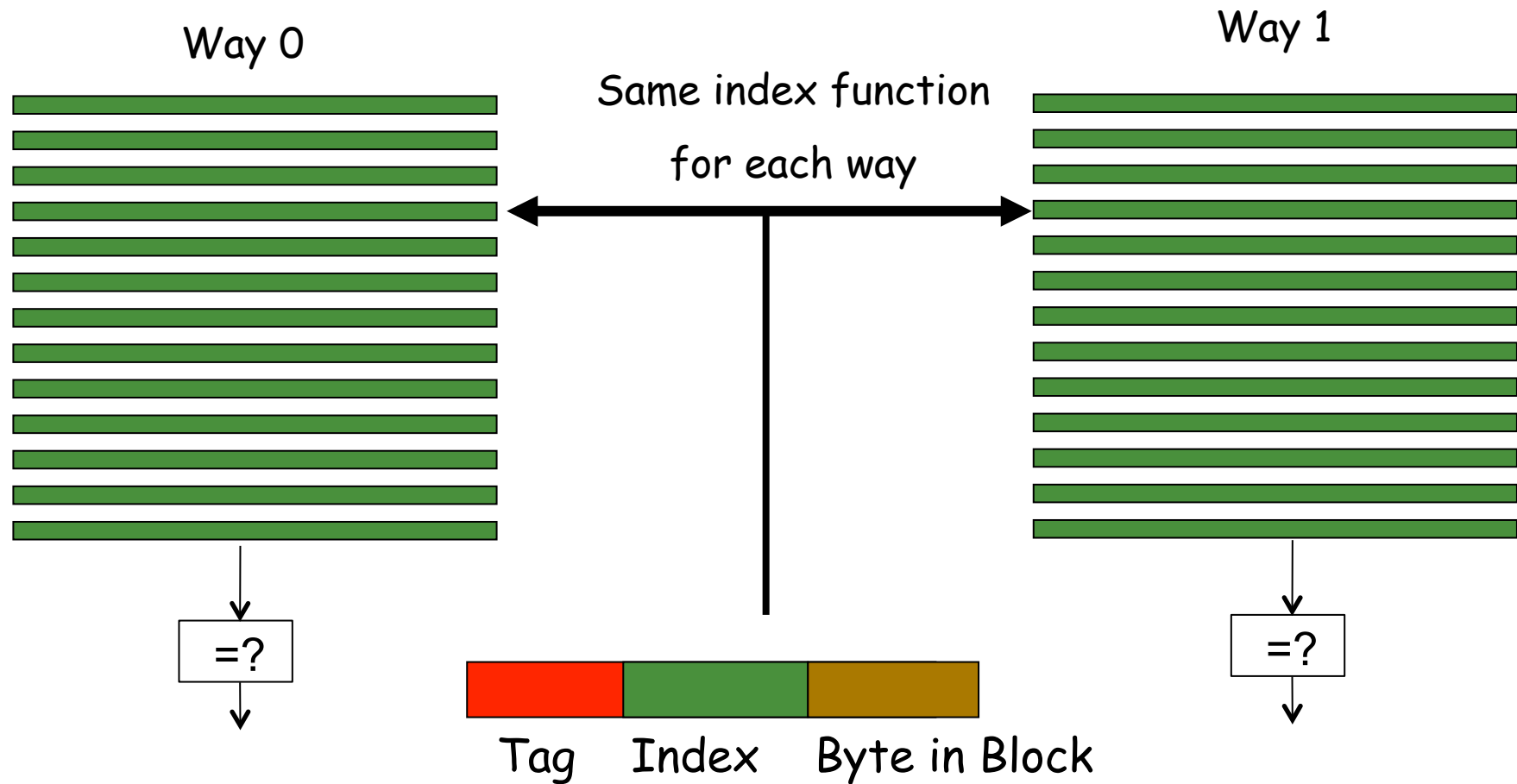
- Hashing: Use better “randomizing” index functions
 - + can reduce conflict misses
 - by distributing the accessed memory blocks more evenly to sets
 - Example of conflicting accesses: strided access pattern where stride value equals number of sets in cache
 - More complex to implement: can lengthen critical path
- Pseudo-associativity (Poor Man’s associative cache)
 - Serial lookup: On a miss, use a different index function and access cache again
 - Given a direct-mapped array with K cache blocks
 - Implement K/N sets
 - Given address Addr, sequentially look up: $\{0, \text{Addr}[\lg(K/N)-1: 0]\}$, $\{1, \text{Addr}[\lg(K/N)-1: 0]\}$, ... , $\{N-1, \text{Addr}[\lg(K/N)-1: 0]\}$

Skewed Associative Caches

- Idea: Reduce conflict misses by using different index functions for each cache way
- Seznec, “A Case for Two-Way Skewed-Associative Caches,” ISCA 1993.

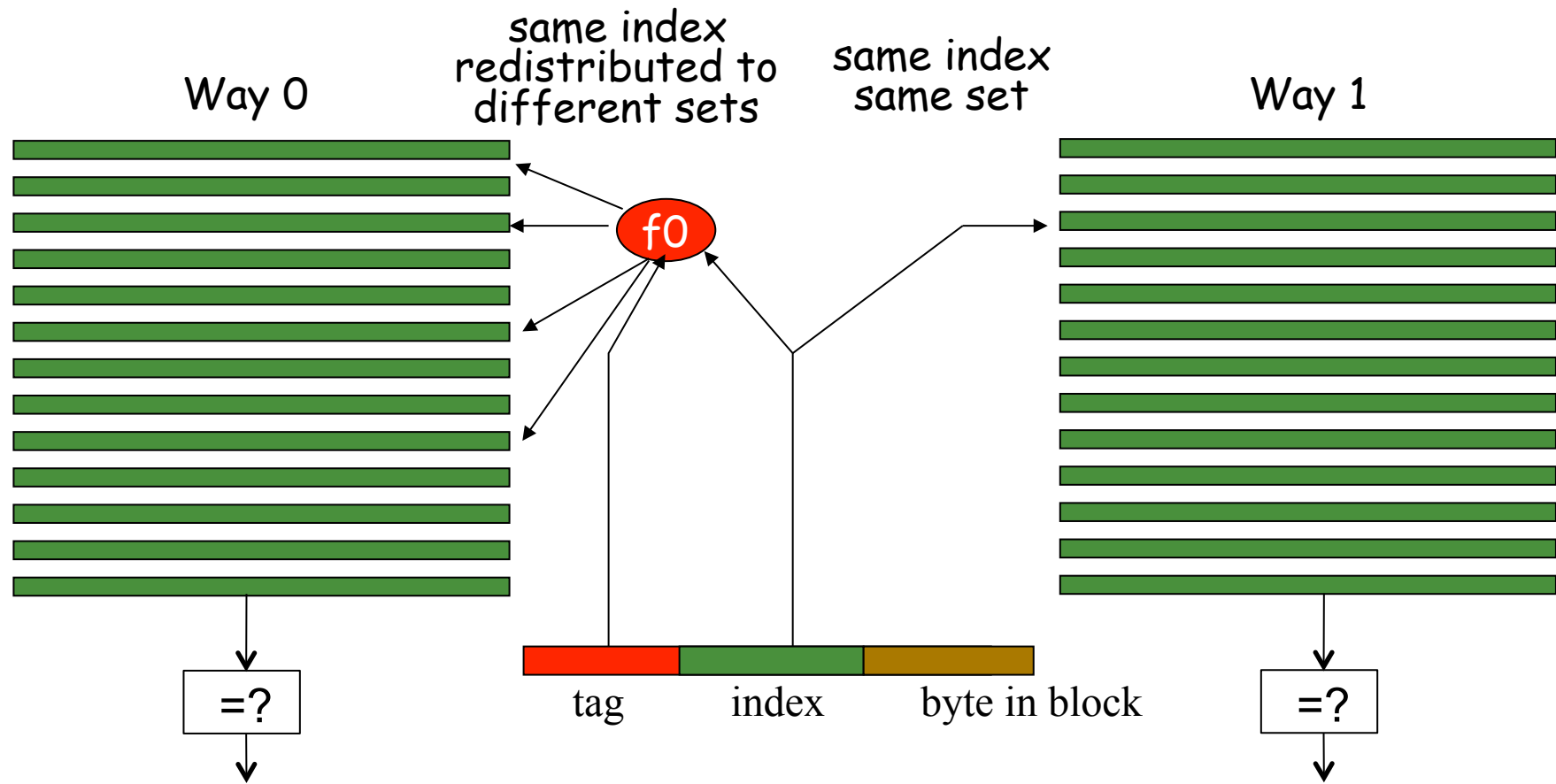
Skewed Associative Caches (I)

- Basic 2-way associative cache structure



Skewed Associative Caches (II)

- Skewed associative caches
 - Each bank has a different index function



Skewed Associative Caches (III)

- Idea: Reduce conflict misses by using **different index functions for each cache way**
- Benefit: indices are more randomized (memory blocks are better distributed across sets)
 - Less likely two blocks have same index (esp. with strided access)
 - Reduced conflict misses
- Cost: additional latency of hash function
- Seznec, “**A Case for Two-Way Skewed-Associative Caches,**” ISCA 1993.

Software Approaches for Higher Hit Rate

- Restructuring data access patterns
- Restructuring data layout

- Loop interchange
- Data structure separation/merging
- Blocking
- ...

Restructuring Data Access Patterns (I)

- Idea: Restructure data layout or data access patterns
- Example: If column-major
 - $x[i+1,j]$ follows $x[i,j]$ in memory
 - $x[i,j+1]$ is far away from $x[i,j]$

Poor code

```
for i = 1, rows
  for j = 1, columns
    sum = sum + x[i,j]
```

Better code

```
for j = 1, columns
  for i = 1, rows
    sum = sum + x[i,j]
```

- This is called **loop interchange**
- Other optimizations can also increase hit rate
 - Loop fusion, array merging, ...
- What if multiple arrays? Unknown array size at compile time?

Restructuring Data Access Patterns (II)

■ Blocking

- ❑ Divide loops operating on arrays into computation chunks so that each chunk can hold its data in the cache
- ❑ Avoids cache conflicts between different chunks of computation
- ❑ Essentially: Divide the working set so that each piece fits in the cache

■ But, there are still self-conflicts in a block

1. there can be conflicts among different arrays
2. array sizes may be unknown at compile/programming time

Restructuring Data Layout (I)

```
struct Node {  
    struct Node* node;  
    int key;  
    char [256] name;  
    char [256] school;  
}  
  
while (node) {  
    if (node→key == input-key) {  
        // access other fields of node  
    }  
    node = node→next;  
}
```

- Pointer based traversal (e.g., of a linked list)
- Assume a huge linked list (1M nodes) and unique keys
- Why does the code on the left have poor cache hit rate?
 - “Other fields” occupy most of the cache line even though rarely accessed!

Restructuring Data Layout (II)

```
struct Node {
    struct Node* node;
    int key;
    struct Node-data* node-data;
}

struct Node-data {
    char [256] name;
    char [256] school;
}

while (node) {
    if (node→key == input-key) {
        // access node→node-data
    }
    node = node→next;
}
```

- Idea: separate frequently-used fields of a data structure and pack them into a separate data structure
- Who should do this?
 - Programmer
 - Compiler
 - Profiling vs. dynamic
 - Hardware?
 - Who can determine what is frequently used?

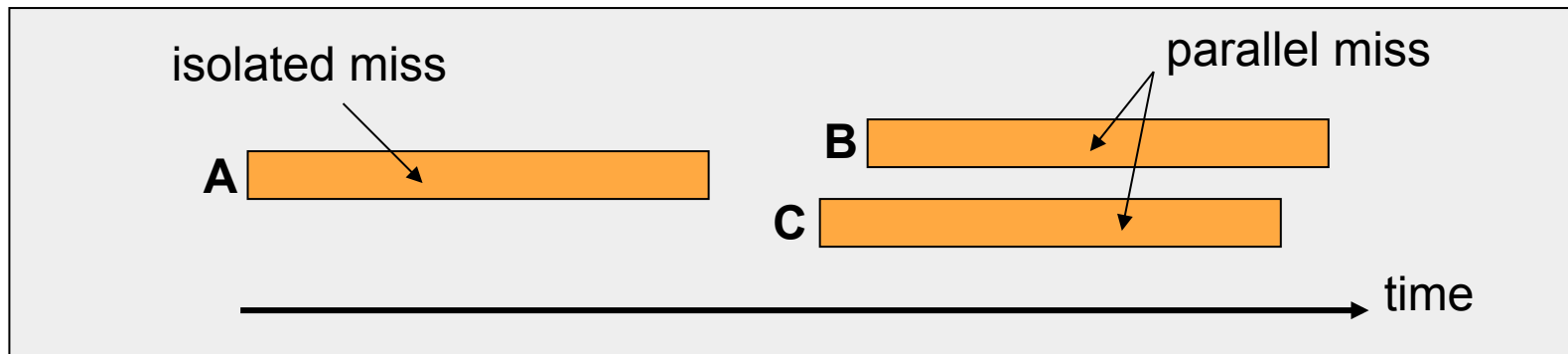
Improving Basic Cache Performance

- Reducing miss rate
 - More associativity
 - Alternatives/enhancements to associativity
 - Victim caches, hashing, pseudo-associativity, skewed associativity
 - Better replacement/insertion policies
 - Software approaches
- Reducing miss latency/cost
 - Multi-level caches
 - Critical word first
 - Subblocking/sectoring
 - Better replacement/insertion policies
 - Non-blocking caches (multiple cache misses in parallel)
 - Multiple accesses per cycle
 - Software approaches

Miss Latency/Cost

- What is miss latency or miss cost affected by?
 - Where does the miss get serviced from?
 - Local vs. remote memory
 - What level of cache in the hierarchy?
 - Row hit versus row miss
 - Queueing delays in the memory controller and the interconnect
 - ...
 - How much does the miss stall the processor?
 - Is it overlapped with other latencies?
 - Is the data immediately needed?
 - ...

Memory Level Parallelism (MLP)



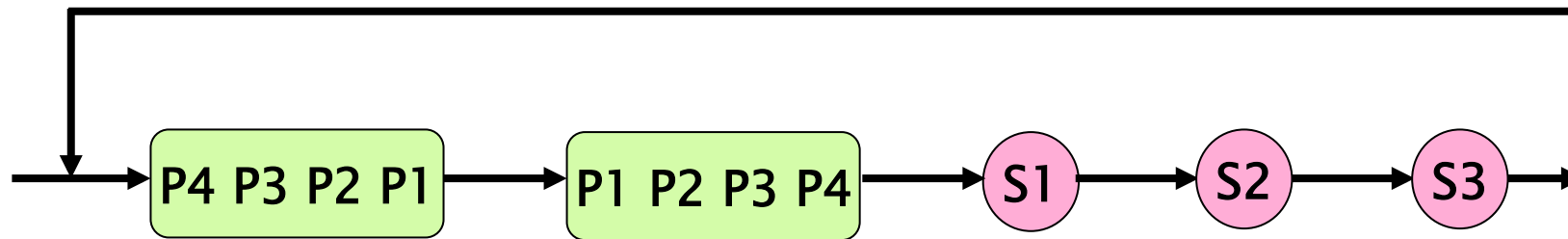
- ❑ Memory Level Parallelism (MLP) means generating and servicing multiple memory accesses in parallel [Glew' 98]
- ❑ Several techniques to improve MLP (e.g., out-of-order execution)
- ❑ MLP varies. Some misses are isolated and some parallel

How does this affect cache replacement?

Traditional Cache Replacement Policies

- ❑ Traditional cache replacement policies try to reduce miss count
- ❑ **Implicit assumption**: Reducing miss count reduces memory-related stall time
- ❑ Misses with varying cost/MLP **breaks** this assumption!
- ❑ Eliminating an isolated miss helps performance more than eliminating a parallel miss
- ❑ Eliminating a higher-latency miss could help performance more than eliminating a lower-latency miss

An Example



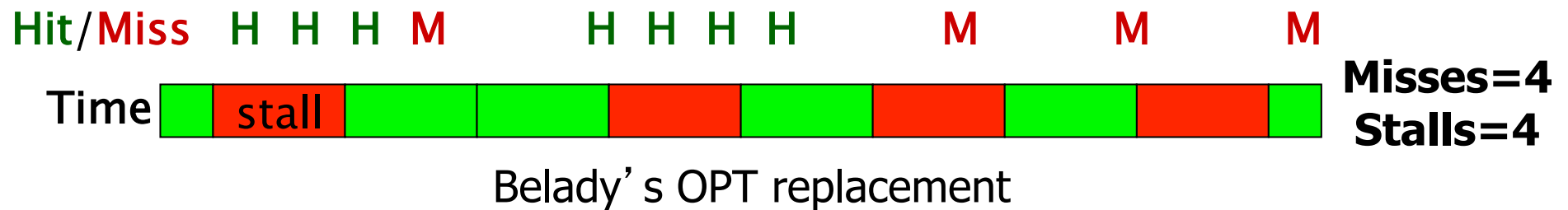
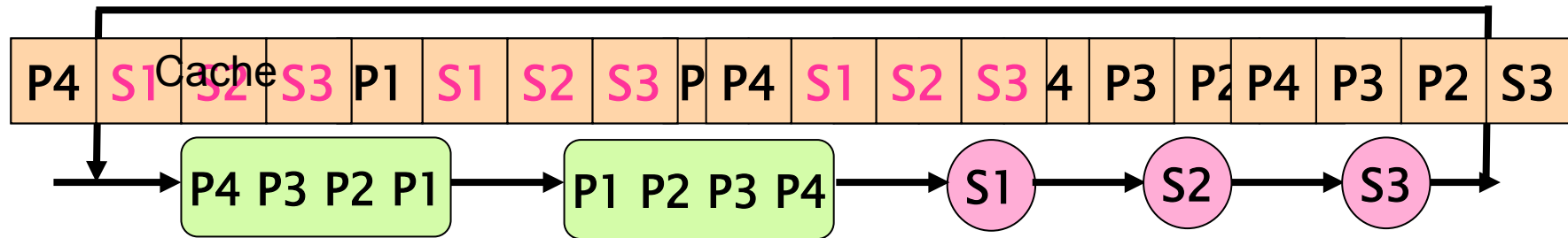
Misses to blocks P1, P2, P3, P4 can be parallel
Misses to blocks S1, S2, and S3 are isolated

Two replacement algorithms:

1. Minimizes miss count (Belady's OPT)
2. Reduces isolated miss (MLP-Aware)

For a fully associative cache containing 4 blocks

Fewest Misses \neq Best Performance



MLP-Aware Cache Replacement

- How do we incorporate MLP into replacement decisions?
- Qureshi et al., “A Case for MLP-Aware Cache Replacement,” ISCA 2006.
 - Reading for review

A Case for MLP-Aware Cache Replacement

Moinuddin K. Qureshi Daniel N. Lynch Onur Mutlu Yale N. Patt
Department of Electrical and Computer Engineering
The University of Texas at Austin
{moin, lynch, onur, patt}@hps.utexas.edu

Other Recommended Cache Papers (I)

- Qureshi et al., “Adaptive Insertion Policies for High Performance Caching,” ISCA 2007.

Adaptive Insertion Policies for High Performance Caching

Moinuddin K. Qureshi[†] Aamer Jaleel[§] Yale N. Patt[†] Simon C. Steely Jr.[§] Joel Emer[§]

[†]ECE Department
The University of Texas at Austin
{moin, patt}@hps.utexas.edu

[§]Intel Corporation, VSSAD
Hudson, MA
{aamer.jaleel, simon.c.steely.jr, joel.emer}@intel.com

Other Recommended Cache Papers (II)

- Seshadri et al., “The Evicted-Address Filter: A Unified Mechanism to Address Both Cache Pollution and Thrashing,” PACT 2012.

The Evicted-Address Filter: A Unified Mechanism to Address Both Cache Pollution and Thrashing

Vivek Seshadri[†]
vseshadr@cs.cmu.edu

Onur Mutlu[†]
onur@cmu.edu

Michael A Kozuch^{*}
michael.a.kozuch@intel.com

Todd C Mowry[†]
tcm@cs.cmu.edu

[†]Carnegie Mellon University

^{*}Intel Labs Pittsburgh

Enabling Multiple Outstanding Misses

Handling Multiple Outstanding Accesses

- Question: If the processor can generate multiple cache accesses, can the later accesses be handled while a previous miss is outstanding?
- Goal: Enable cache access when there is a pending miss
- Goal: Enable multiple misses in parallel
 - Memory-level parallelism (MLP)
- Solution: Non-blocking or lockup-free caches
 - Kroft, “Lockup-Free Instruction Fetch/Prefetch Cache Organization,” ISCA 1981.

Handling Multiple Outstanding Accesses

- Idea: Keep track of the status/data of misses that are being handled in Miss Status Handling Registers (MSHRs)
 - A cache access checks MSHRs to see if a miss to the same block is already *pending*.
 - If pending, a new request is not generated
 - If pending and the needed data available, data forwarded to later load
 - Requires buffering of outstanding miss requests

Miss Status Handling Register

- Also called “miss buffer”
- Keeps track of
 - Outstanding cache misses
 - Pending load/store accesses that refer to the missing cache block
- Fields of a single MSHR entry
 - Valid bit
 - Cache block address (to match incoming accesses)
 - Control/status bits (prefetch, issued to memory, which subblocks have arrived, etc)
 - Data for each subblock
 - For each pending load/store
 - Valid, type, data size, byte in block, destination register or store buffer entry address

Miss Status Handling Register Entry

1	27	1
Valid	Block Address	Issued

1	3	5	5	
Valid	Type	Block Offset	Destination	Load/store 0
Valid	Type	Block Offset	Destination	Load/store 1
Valid	Type	Block Offset	Destination	Load/store 2
Valid	Type	Block Offset	Destination	Load/store 3

MSHR Operation

- On a cache miss:
 - Search MSHRs for a pending access to the same block
 - Found: Allocate a load/store entry in the same MSHR entry
 - Not found: Allocate a new MSHR
 - No free entry: stall
- When a subblock returns from the next level in memory
 - Check which loads/stores waiting for it
 - Forward data to the load/store unit
 - Deallocate load/store entry in the MSHR entry
 - Write subblock in cache or MSHR
 - If last subblock, deallocate MSHR (after writing the block in cache)

Non-Blocking Cache Implementation

- When to access the MSHRs?
 - In parallel with the cache?
 - After cache access is complete?
- MSHRs need not be on the critical path of hit requests
 - Which one below is the common case?
 - Cache miss, MSHR hit
 - Cache hit

Enabling High Bandwidth Memories

Multiple Instructions per Cycle

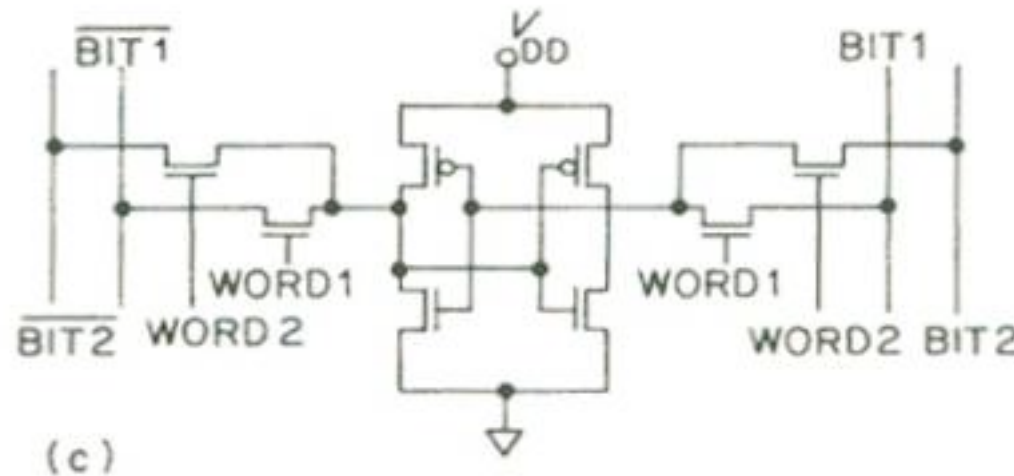
- Can generate multiple cache/memory accesses per cycle
- How do we ensure the cache/memory can handle multiple accesses in the same clock cycle?

- Solutions:
 - true multi-porting
 - virtual multi-porting (time sharing a port)
 - multiple cache copies
 - banking (interleaving)

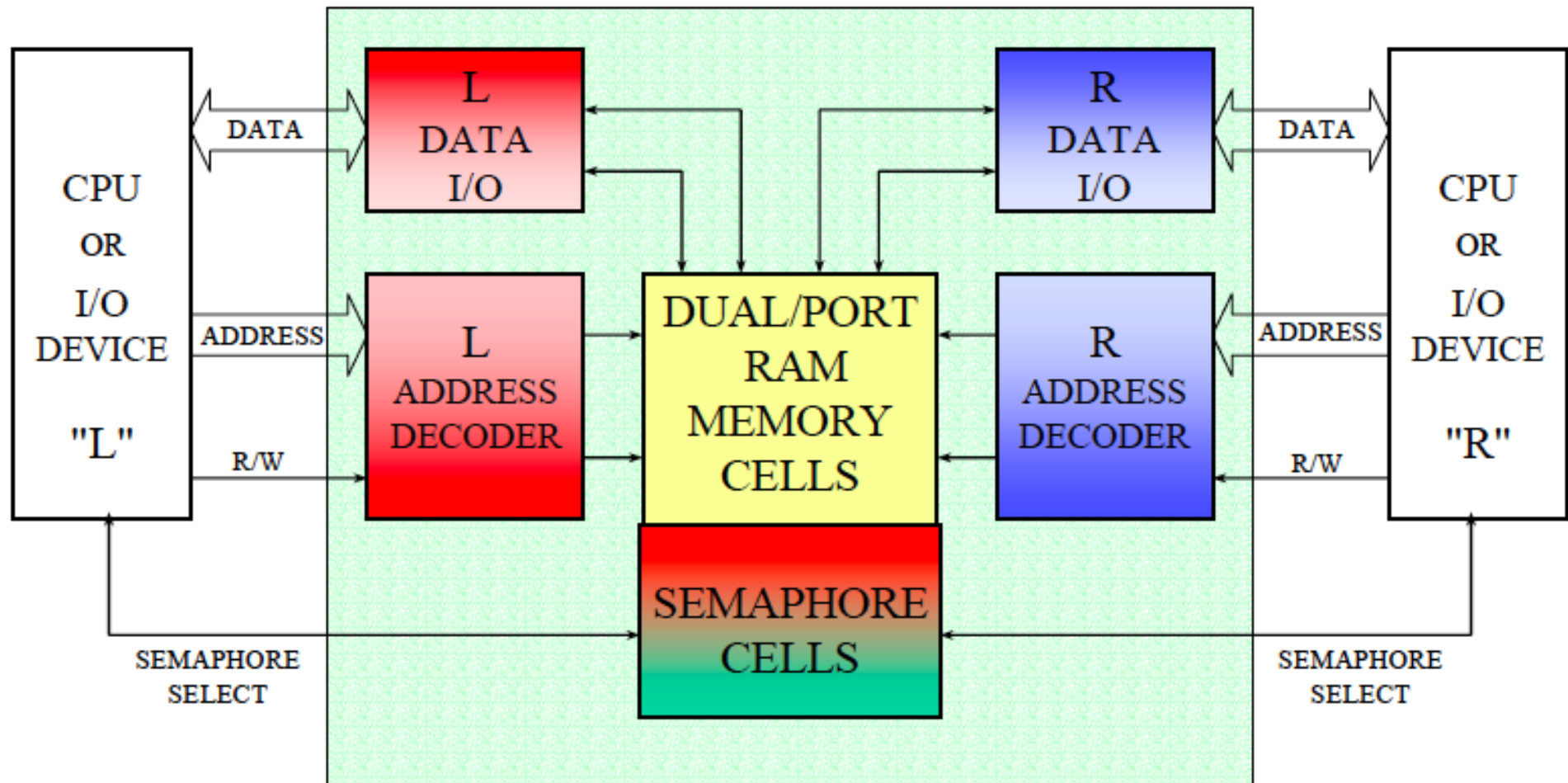
Handling Multiple Accesses per Cycle (I)

■ True multiporting

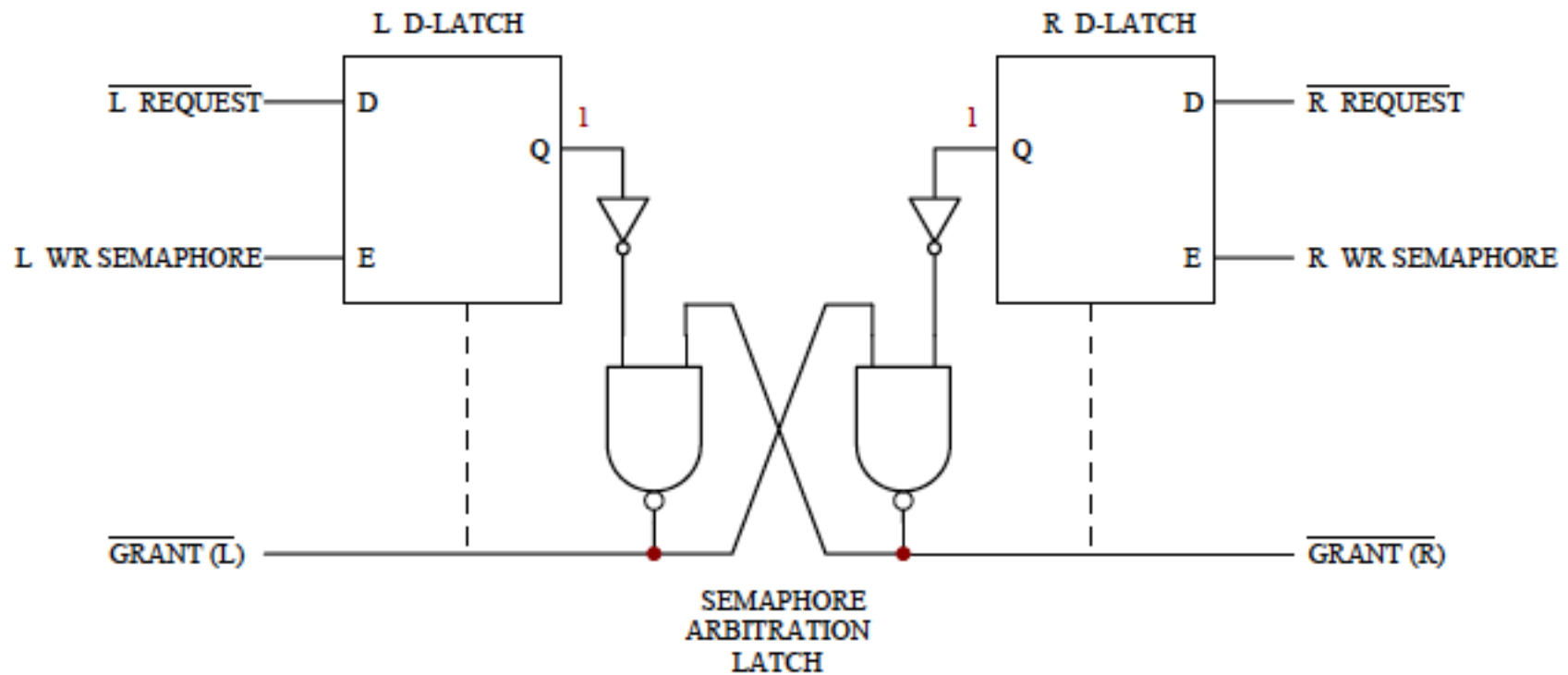
- Each memory cell has multiple read or write ports
- + Truly concurrent accesses (no conflicts on read accesses)
- Expensive in terms of latency, power, area
- What about read and write to the same location at the same time?
 - Peripheral logic needs to handle this



Peripheral Logic for True Multiporting



Peripheral Logic for True Multiporting



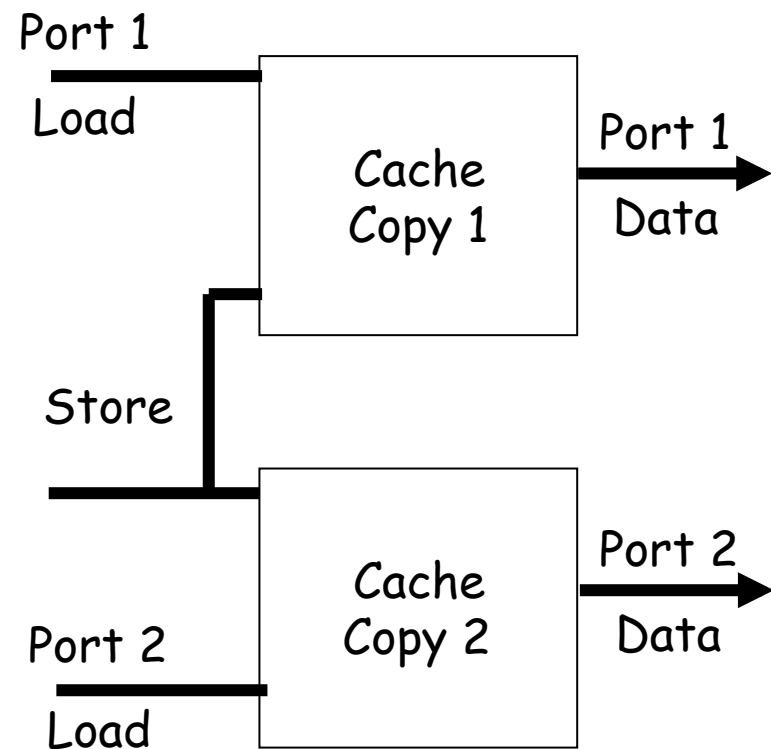
Handling Multiple Accesses per Cycle (II)

- Virtual multiporting

- Time-share a single port
- Each access needs to be (significantly) shorter than clock cycle
- Used in Alpha 21264
- Is this scalable?

Handling Multiple Accesses per Cycle (III)

- Multiple cache copies
 - ❑ Stores update both caches
 - ❑ Loads proceed in parallel
- Used in Alpha 21164
- Scalability?
 - ❑ Store operations cause a bottleneck
 - ❑ Area proportional to “ports”



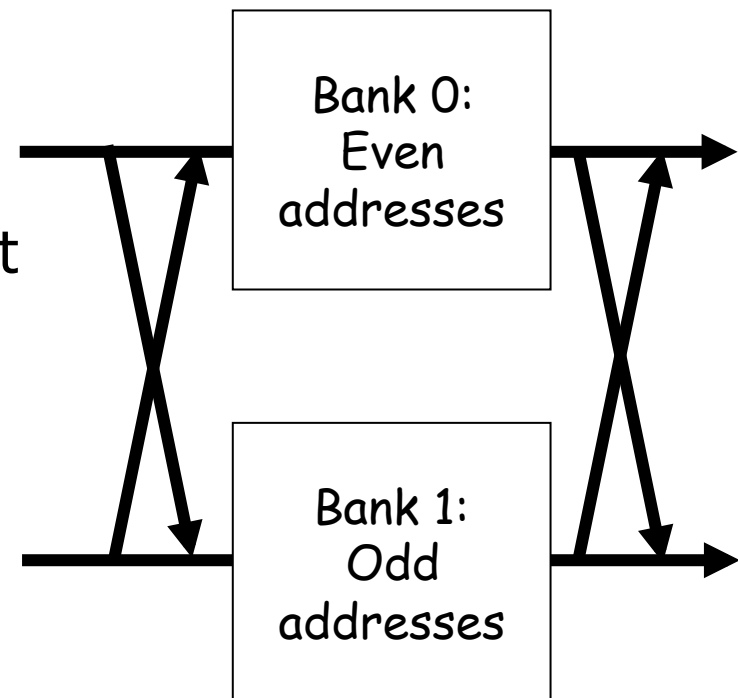
Handling Multiple Accesses per Cycle (III)

■ Banking (Interleaving)

- Bits in address determines which bank an address maps to
 - Address space partitioned into separate banks
 - Which bits to use for “bank address”?
- + No increase in data store area
- Cannot satisfy multiple accesses to the same bank
- Crossbar interconnect in input/output

■ Bank conflicts

- Two accesses are to the same bank
- How can these be reduced?
 - Hardware? Software?



General Principle: Interleaving

■ Interleaving (banking)

- **Problem:** a single monolithic memory array takes long to access and does not enable multiple accesses in parallel
- **Goal:** Reduce the latency of memory array access and enable multiple accesses in parallel
- **Idea:** Divide the array into multiple banks that can be accessed independently (in the same cycle or in consecutive cycles)
 - Each bank is smaller than the entire memory storage
 - Accesses to different banks can be overlapped
- **A Key Issue:** How do you map data to different banks? (i.e., how do you interleave data across banks?)

Further Readings on Caching and MLP

- **Required:** Qureshi et al., “A Case for MLP-Aware Cache Replacement,” ISCA 2006.
- **One Pager:** Glew, “MLP Yes! ILP No!,” ASPLOS Wild and Crazy Ideas Session, 1998.
- Mutlu et al., “Runahead Execution: An Effective Alternative to Large Instruction Windows,” IEEE Micro 2003.

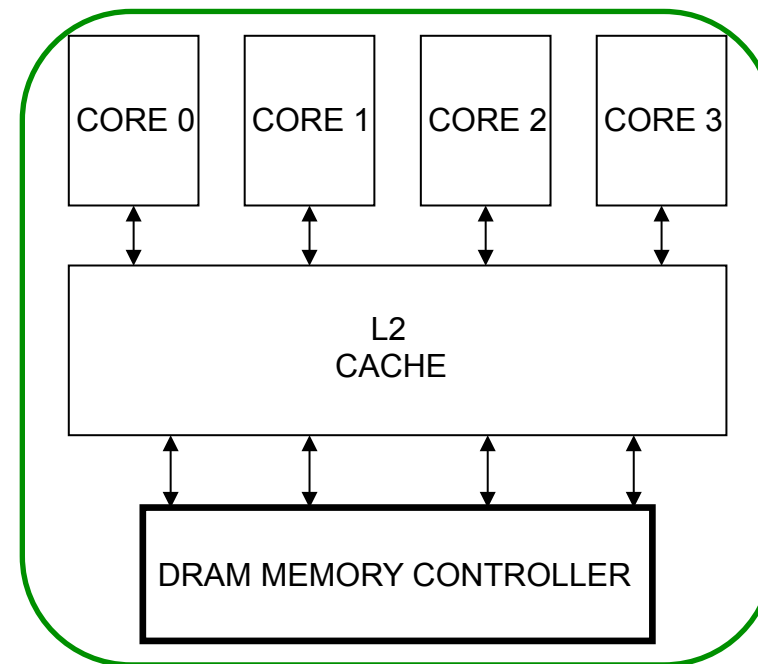
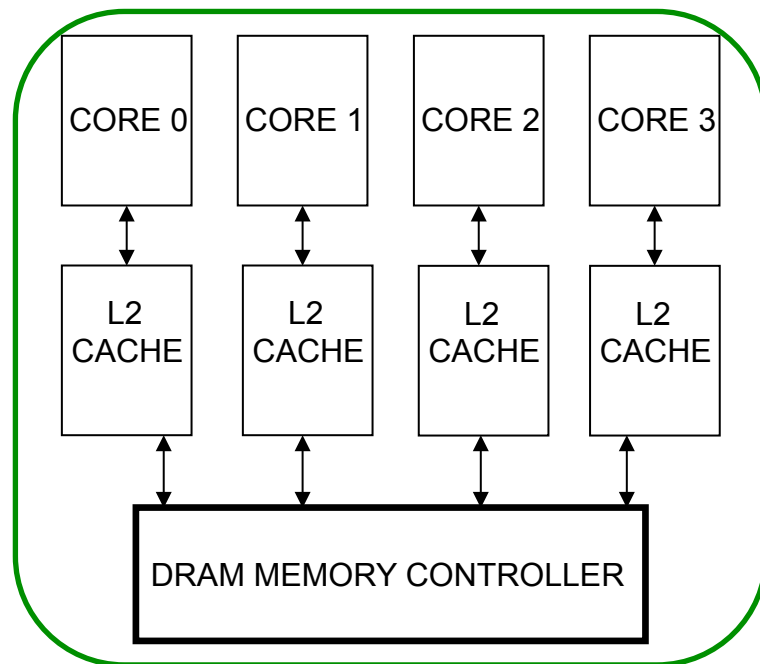
Multi-Core Issues in Caching

Caches in Multi-Core Systems

- Cache efficiency becomes even more important in a multi-core/multi-threaded system
 - ❑ Memory bandwidth is at premium
 - ❑ Cache space is a limited resource
- How do we design the caches in a multi-core system?
- Many decisions
 - ❑ Shared vs. private caches
 - ❑ How to maximize performance of the entire system?
 - ❑ How to provide QoS to different threads in a shared cache?
 - ❑ Should cache management algorithms be aware of threads?
 - ❑ How should space be allocated to threads in a shared cache?

Private vs. Shared Caches

- **Private** cache: Cache belongs to one core (a shared block can be in multiple caches)
- **Shared** cache: Cache is shared by multiple cores



Resource Sharing Concept and Advantages

- Idea: Instead of dedicating a hardware resource to a hardware context, allow multiple contexts to use it
 - Example resources: functional units, pipeline, caches, buses, memory
 - Why?
-
- + Resource sharing improves utilization/efficiency → throughput
 - When a resource is left idle by one thread, another thread can use it; no need to replicate shared data
 - + Reduces communication latency
 - For example, shared data kept in the same cache in multithreaded processors
 - + Compatible with the shared memory model

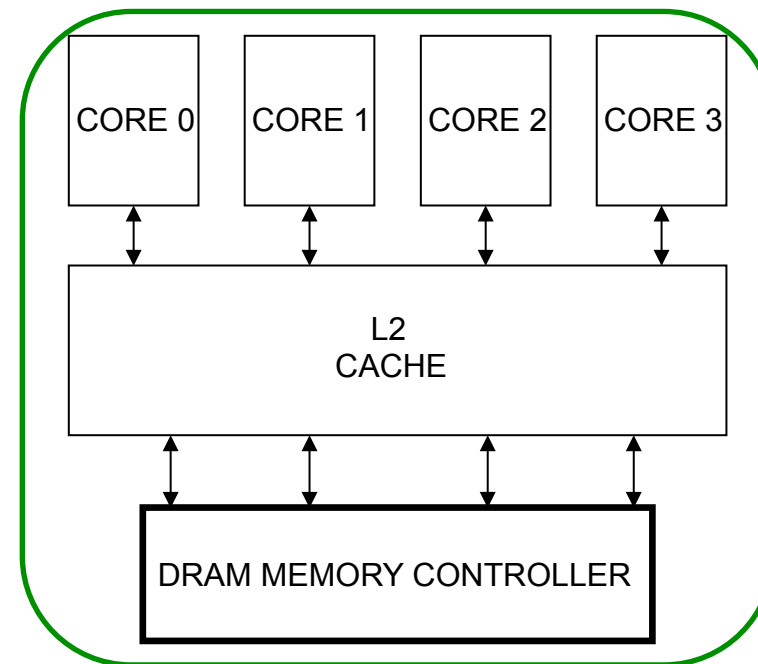
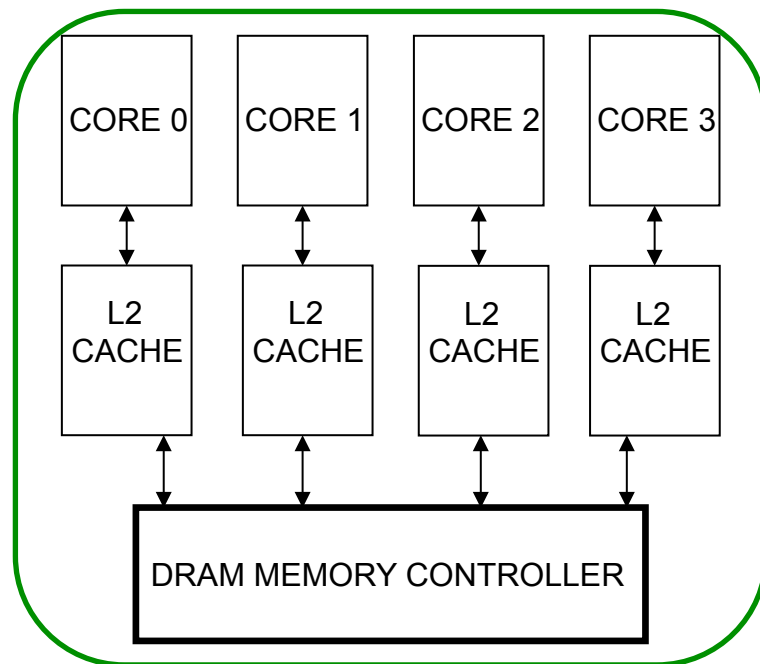
Resource Sharing Disadvantages

- Resource sharing results in **contention for resources**
 - When the resource is not idle, another thread cannot use it
 - If space is occupied by one thread, another thread needs to re-occupy it
- **Sometimes reduces each or some thread's performance**
 - Thread performance can be worse than when it is run alone
- **Eliminates performance isolation** → inconsistent performance across runs
 - Thread performance depends on co-executing threads
- Uncontrolled (free-for-all) sharing **degrades QoS**
 - Causes unfairness, starvation

Need to efficiently and fairly utilize shared resources

Private vs. Shared Caches

- **Private** cache: Cache belongs to one core (a shared block can be in multiple caches)
- **Shared** cache: Cache is shared by multiple cores



Shared Caches Between Cores

■ Advantages:

- High effective capacity
- **Dynamic partitioning** of available cache space
 - No fragmentation due to static partitioning
- **Easier to maintain coherence** (a cache block is in a single location)
- **Shared data and locks do not ping pong between caches**

■ Disadvantages

- Slower access
- Cores incur **conflict misses due to other cores' accesses**
 - Misses due to inter-core interference
 - Some cores can destroy the hit rate of other cores
- Guaranteeing a minimum level of service (or fairness) to each core is harder (how much space, how much bandwidth?)

Shared Caches: How to Share?

■ Free-for-all sharing

- ❑ Placement/replacement policies are the same as a single core system (usually LRU or pseudo-LRU)
- ❑ Not thread/application aware
- ❑ An incoming block evicts a block regardless of which threads the blocks belong to

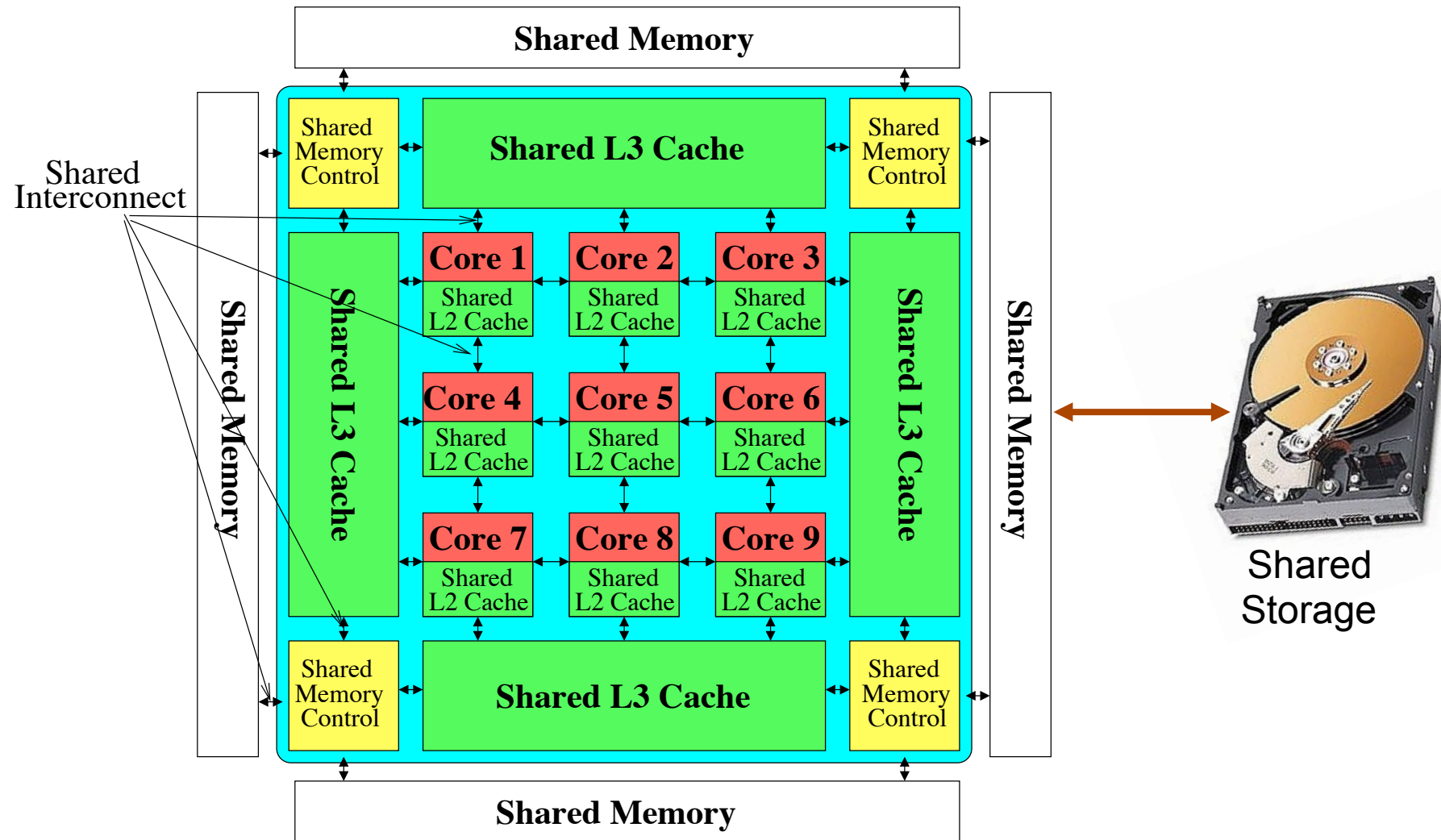
■ Problems

- ❑ Inefficient utilization of cache: LRU is not the best policy
- ❑ A cache-unfriendly application can destroy the performance of a cache friendly application
- ❑ Not all applications benefit equally from the same amount of cache: free-for-all might prioritize those that do not benefit
- ❑ Reduced performance, reduced fairness

Example: Utility Based Shared Cache Partitioning

- Goal: Maximize system throughput
- Observation: Not all threads/applications benefit equally from caching → simple LRU replacement not good for system throughput
- Idea: Allocate more cache space to applications that obtain the most benefit from more space
- The high-level idea can be applied to other shared resources as well.
- Qureshi and Patt, “Utility-Based Cache Partitioning: A Low-Overhead, High-Performance, Runtime Mechanism to Partition Shared Caches,” MICRO 2006.
- Suh et al., “A New Memory Monitoring Scheme for Memory-Aware Scheduling and Partitioning,” HPCA 2002.

The Multi-Core System: *A Shared Resource View*



Need for QoS and Shared Resource Mgmt.

- Why is unpredictable performance (or lack of QoS) bad?
- Makes programmer's life difficult
 - An optimized program can get low performance (and performance varies widely depending on co-runners)
- Causes discomfort to user
 - An important program can starve
 - Examples from shared software resources
- Makes system management difficult
 - How do we enforce a Service Level Agreement when hardware resources are sharing is uncontrollable?

Resource Sharing vs. Partitioning

- Sharing improves throughput
 - Better utilization of space
- Partitioning provides performance isolation (predictable performance)
 - Dedicated space
- Can we get the benefits of both?
- Idea: Design shared resources such that they are efficiently utilized, controllable and partitionable
 - No wasted resource + QoS mechanisms for threads

Shared Hardware Resources

- Memory subsystem (in both multithreaded and multi-core systems)
 - Non-private caches
 - Interconnects
 - Memory controllers, buses, banks
 - I/O subsystem (in both multithreaded and multi-core systems)
 - I/O, DMA controllers
 - Ethernet controllers
 - Processor (in multithreaded systems)
 - Pipeline resources
 - L1 caches
-

Computer Architecture

Lecture 2: Fundamentals, Memory Hierarchy, Caches

Prof. Onur Mutlu

ETH Zurich

Fall 2017

21 September 2017