# CS3350B Computer Organization Introduction

Alex Brandt

Department of Computer Science University of Western Ontario, Canada

Tuesday January 8, 2019

# Outline

### 1 Highlights of Hardware History

2 Modern Computer Architectures

3 System and Hardware Abstractions



Konrad Zuse's Z3 electro-mechanical computer (1941, Germany). Turing complete, though conditional jumps were missing.

Alex Brandt

CS3350B Computer Organization Introduction



Colossus (UK, 1941) was the world's first totally electronic programmable computing device. But not Turing complete.



Harvard Mark I – IBM ASCC (1944, US). Electro-mechanical computer. No conditional jumps and not Turing complete. It could store 72 numbers, each 23 decimal digits long. It could do three additions or subtractions in a second. A multiplication took six seconds, a division took 15.3 seconds, and a logarithm or a trigonometric function took over one minute.



Electronic Numerical Integrator And Computer (ENIAC, 1945). The first general-purpose, electronic computer. It was a Turing-complete, digital computer capable of being reprogrammed and was running at 5,000 cycles per second for operations on the 10-digit numbers.



The IBM Personal Computer (IBM PC) (Introduced on August 12, 1981).

# Outline

#### 1 Highlights of Hardware History

#### 2 Modern Computer Architectures

3 System and Hardware Abstractions

### Von Neumann Architecture



The Von Neumann architecture saw a shared memory for instructions and data. Modern computers, especially w.r.t what software "sees", use his model.

Alex Brandt	CS3350B Computer Organization Introduction	Tuesday January 8, 2019
-------------	--	-------------------------

9 / 28

### Multi-core Architecture



Modern multi-core processors see a varying hierarchy of independent and shared memories.

# Memory Hierarchy



11 / 28

# Cache Memory Statistics

L1 Data Cache									
Size	Line Size	e Latency Associa							
32 KB	64 bytes	3 cycles	les 8-way						
L1 Instruction Cache									
Size	Line Size	Latency	Associativty						
32 KB	64 bytes	3 cycles	8-way						
L2 Cache									
Size	Line Size	Latency	Associativty						
6 MB	64 bytes	14 cycles	24-way						

Typical cache specifications of a multicore in 2008.

# The CPU-Memory Gap

# The increasing gap between DRAM, disk, and CPU speeds.



The **Processor-Memory Gap** is a key contributor to the **Memory Wall** – the point where a program's performance is totally determined by memory speed. Faster processors will not make programs run faster!

# Outline

- 1 Highlights of Hardware History
- 2 Modern Computer Architectures
- 3 System and Hardware Abstractions

# Components of a computer



Same components for all kinds of computer

# Layers of Software



#### Application software

- → Written in a high-level language (HLL)
- → Javascript, Python, Java.

### System software

- → Compiler: translates HLL code to machine code
- - Handling input/output
  - Managing memory and storage
  - Scheduling tasks & sharing resources

#### Hardware

# Levels of Program Code

### High-level language

- ↓ Level of abstraction closer to problem domain
- Provides productivity and portability

#### Assembly language

→ Textual representation of machine instructions

#### Hardware representation

⇒ Binary encoding of instructions and data



# Layers of Abstraction in Computer Systems



After a brief look at processors and memory, we will start down at circuit design and build up to ISA.

# Complicated Systems are Inevitable

Think economic specialization.

Each layer is separate and can be optimized independently in the pursuit of performance.

Eight great ideas in computer architecture (Patterson)

- 1 Use abstraction to simplify design
- 2 Design for Moore's Law
- 3 Make the common case fast
- 4 Performance via parallelism
- 5 Performance via pipelining
- 6 Performance via prediction
- 7 Hierarchy of memories
- 8 Dependability via redundancy



### Great Idea #1: Abstraction



temp = v[k]; v[k] = v[k+1]; v[k+1] = temp;

lw \$t0, 0(\$2)
lw \$t1, 4(\$2)
sw \$t1, 0(\$2)
sw \$t0, 4(\$2)

# Anything can be represented as a
# number, i.e., data or instructions
0000 1001 1100 0110 1010 1111 0101 1000
1010 1111 0101 1000 0000 1001 1100 0110
1100 0110 1010 1111 0101 1000 0000 1001
0101 1000 0000 1001 1100 0110 1010 1111

### Great idea #2: Design for Moore's Law



# Great idea #4: Performance via Parallelism



Parallelism can be used to tackle the processor-memory gap.

# Great idea #5: Performance via Pipelining

	Time I	Time 2	Time 3	Time 4	Time 5	Time 6	Time 7	Time 8
instruction	Instruction fetch	Operand fetch	Execute	Operand store				
instruction 2		Instruction fetch	Operand fetch	Execute	Operand store			
instruction 3			Instruction fetch	Operand fetch	Execute	Operand store		
			-					
instruction 4				Instruction fetch	Operand fetch	Execute	Operand store	
	ln time	slot 3, instruc	tion I					
instruction 5	is being is in the and ins from m	executed, inst e operand fetc truction 3 is b emor <b>y</b>	ruction 2 h phase, eing fetched		Instruction fetch	Operand fetch	Execute	Operand store

# Great idea #7: Memory Hierarchy (Principle of Locality)



# Great Idea #8: Dependability via Redundancy

Redundancy so that a failing piece doesn't make the whole system fail



Increasing transistor density reduces the cost of redundancy

# Great Idea #8: Dependability via redundancy

Applies to everything from physical hardware components to the data encoded in binary.

- Redundant data centers so that can lose 1 datacenter but Internet service stays online
- Redundant disks so that can lose 1 disk but not lose data (Redundant Arrays of Independent Disks/RAID)
- Redundant memory bits of so that can lose 1 bit but no data (Error Correcting Code/ECC Memory)



# In This Course...

#### Check the syllabus!

- Course topics
- Evaluation
- Schedule
- Textbook (recommended, not required)
- Missed quiz policies
- TAs, Office hours

Check OWL for continued updates and announcements. It's your fault if you miss something that was announced.

The lecture slides of this course are adapted from previous years' slides by Marc Moreno Maza which in turn have been adapted from text book accompaniments and teaching materials posted on the internet by other Computer Architecture instructors.

# CS3350B Computer Organization Chapter 1: CPU and Memory Part 1: The CPU

Alex Brandt

Department of Computer Science University of Western Ontario, Canada

Thursday January 10, 2019

# Outline

#### 1 The Basics

2 Clock Cycles per Instruction (CPI)

3 Power, Trends, Limitations

4 Benchmarks and Profiling

2 / 34

### Components of a computer



**Micro-architecture**: the internals of a CPU (control and datapath) will come later in the course. For now, we look at the CPU as a whole.

CPU and memory are highly coupled, especially when we look at performance. This will be seen by the end of this chapter.

# Programmer's View of CPU Performance

At a basic level, the running time of some program on a CPU is determined by:

- The clock rate of the CPU (e.g. 3.4 GHz),
- The type of instructions being performed,
  - → Addition/Subtraction faster than Multiplication/Division, etc.
  - → Affects the *average* clock cycles per instruction (CPI)
- Memory access time.
  - L→ Recall processor-memory gap

Assuming CPU clock rate is fixed, programmers can influence program performance by changing the type of instructions they use as well as how their code accesses memory.

# Aside: Changing Instructions for Performance

On 6th-generation Intel CPUs

32-bit integer division ~26 clock cycles vs. logical bit shift ~1 clock cycle

int i = 1234567; i /= 2;
 int i = 1234567; i >>= 1;

Floating point division ~14 clock cycles vs. multiplication ~5 clock cycles

↓ float x = 1.2f; x /= 2.0f;
↓ float x = 1.2f; x \*= 0.5f;

Fast Inverse Square Root thanks to *Quake*: https://en.wikipedia.org/wiki/Fast\_inverse\_square\_root

# Understanding and Analyzing Performance

- Algorithmic analysis: Estimating the complexity of algorithms in some abstract, idealized way. In reality, two different O(n<sup>2</sup>) algorithms can have wildly different running times.
- Programming language, compiler, architecture:
  - → Programming language and corresponding compiler determine the actual machine instructions the CPU will perform.
  - Gesulting number and type of machine instructions vary by compiler and language and therefore resulting performance varies.
- Processor and Memory: Determines how fast instructions are executed and how fast data moves to and from the processor.
- I/O system (including OS): Determines how fast I/O operations are executed

### Need for Performance Metrics

- Purchasing Perspective. For a collection of machine which one has the
  - $\mapsto$  "best" cost?
  - → "best" cost relative to performance?
- Design perspective. Given many design options and directions which one has the
  - → "best" performance improvement?
  - └→ "best" cost relative to performance improvement?

In either case we need: (i) a basis for comparison, (ii) metrics for evaluation.

Our goal is to understand what factors in the architecture contribute to the overall system performance and the relative importance (and cost) of these factors.

# CPU Performance: Latency

CPU performance is largely measured by **latency**, **throughput**, **clock frequency**.

- Want reduced response time (aka execution time, aka latency) the time between the start and the completion of a task.
  - Important to general PC users.
- To maximize performance of some code segment (program) *X*, we need to minimize execution time.

 $performance_X = 1/execution\_time_X$ 

If X is n times faster than Y, then

 $\frac{\text{performance}_X}{\text{performance}_Y} = \frac{\text{execution\_time}_Y}{\text{execution\_time}_X} = n$ 

Note: Can also compare latency of the same single instruction on different CPUs.
# CPU Performance: Throughput

- Want increased throughput the total amount of work done in a given unit of time.
  - Important to users like data center managers.
- Again, throughput depends on the code segment being executed. Different instructions result in different throughput measures.
- Decreasing response time usually improves throughput, but other factors are important (task scheduling, memory bandwidth, etc.).

# CPU Performance: Clock Frequency

Clock Frequency is a code-agnostic measure of CPU performance.

- Typically, a faster clock (higher frequency) yields a higher performing CPU.
- **But** the *micro-architecture* and the *instruction set architecture* play a large role.
  - $\downarrow$  They influence how much work the CPU does per cycle (i.e. efficiency)

#### Example:

CPU A runs at 3 GHz and a division takes 20 cycles. CPU B runs at 2 GHz and a division takes 10 cycles.

20 cycles / 3 GHz = 20 / 300000000 Hz = 6.66ns10 cycles / 2 GHz = 10 / 200000000 Hz = 5.00ns

# **CPU** Clocking

Synchronous digital systems (e.g. a CPU) are governed and controlled by a clock. This clock synchronizes internal circuits, memory states, and data movement.



Length of clock period determined by internal circuits and micro-architecture design. We will look at this with CPU datapaths and pipelining.

Alex Brandt	Chapter 1: CPU and Memory, Part 1: The CPU	Thursday Janu
-------------	--	---------------

Thursday January 10, 2019

# **CPU** Clocking



Clock period (cycle): duration of a clock cycle (CC)

- determines the speed of a computer processor
- Caveat: again, not necessarily latency or throughput though

e.g., 250ps = 0.25ns = 
$$250 \times 10^{-12} s$$

Clock frequency or rate (CR): cycles per second

the inverse of the clock period

e.g., 
$$3.0$$
GHz =  $3000$ MHz =  $3.0 \times 10^9$ Hz

CR = 1 / CC.

# CPU Time

- It is important to distinguish *elapsed time* and the *time spent on your* task.
  - → Wall time vs. CPU time

  - CPU execution time = #CPU clock cycles × clock cycle for a program for a program

or

- CPU execution time = #CPU clock cycles / clock rate for a program for a program
- We can improve performance by reducing either the *length of the clock cycle* or the **number of clock cycles required for a program**.

# Outline

#### 1 The Basics

### 2 Clock Cycles per Instruction (CPI)

3 Power, Trends, Limitations

4 Benchmarks and Profiling

### Instruction Performance

#CPU clock cycles = #Instructions × Average # of clock cycles for a program for a program per instruction

Clock cycles per instruction (CPI) - the average number of clock cycles each instruction takes to execute.

- □→ Different instructions may take different amounts of time depending on what they do.
- → Calculated by a simple averaging of each instruction type in a program and their corresponding number of cycles.

### The Classic Performance Equation

CPUtime	=	$Instruction\_count \times CPI \times clock\_cycle$
	or	
CPUtime	=	Instruction_count $\times$ CPI/clock_rate

- Keep in mind that the only complete and reliable measure of computer performance is time.
- For example, redesigning the hardware implementation of an instruction set to lower the instruction count may lead to an organization with

  - $\vdash$  higher CPI,

that offsets the improvement in instruction count.

■ *Note:* CPI depends on the type of instruction executed, so the code which executes the fewest instructions may not be the fastest.

A Simple Example (1/2)

Overall effective CPI = 
$$\sum_{i=1}^{n} (CPI_i \times IC_i)$$

Ор	Freq	$CPI_i$	$Freq \times CPI_i$	(1)
ALU	50%	1	.5	.5
Load	20%	5	1.0	
Store	10%	3	.3	.3
Branch	20%	2	.4	.4
			$\Sigma = 2.2$	

(1) How much faster would the machine be if a better data cache reduced the average load time to 2 cycles?

A Simple Example (1/2)

Overall effective CPI = 
$$\sum_{i=1}^{n} (CPI_i \times IC_i)$$

Ор	Freq	$CPI_i$	$Freq \times CPI_i$	(1)
ALU	50%	1	.5	.5
Load	20%	5	1.0	.4
Store	10%	3	.3	.3
Branch	20%	2	.4	.4
			$\Sigma = 2.2$	1.6

(1) How much faster would the machine be if a better data cache reduced the average load time to 2 cycles? CPU time new =  $1.6 \times IC \times CC$ ; so 2.2 versus 1.6 which means 37.5% faster

A Simple Example (2/2)

Overall effective CPI = 
$$\sum_{i=1}^{n} (CPI_i \times IC_i)$$

Ор	Freq	$CPI_i$	$Freq \times CPI_i$	(2)	(3)
ALU	50%	1	.5	.5	
Load	20%	5	1.0	1.0	1.0
Store	10%	3	.3	.3	.3
Branch	20%	2	.4		.4
			$\Sigma = 2.2$		

- (2) How does this CPI compare with using branch prediction to save a cycle off the branch time?
- (3) What if two ALU instructions could be executed at once?

A Simple Example (2/2)

Overall effective CPI = 
$$\sum_{i=1}^{n} (CPI_i \times IC_i)$$

Ор	Freq	$CPI_i$	$Freq \times CPI_i$	(2)	(3)
ALU	50%	1	.5	.5	
Load	20%	5	1.0	1.0	1.0
Store	10%	3	.3	.3	.3
Branch	20%	2	.4	.2	.4
			$\Sigma = 2.2$	2.0	

(2) How does this CPI compare with using branch prediction to save a cycle off the branch time?
 CPU time new = 2.0 × IC × CC so 2.2 versus 2.0 means 10% faster

(3) What if two ALU instructions could be executed at once?

A Simple Example (2/2)

Overall effective CPI = 
$$\sum_{i=1}^{n} (CPI_i \times IC_i)$$

Ор	Freq	$CPI_i$	$Freq \times CPI_i$	(2)	(3)
ALU	50%	1	.5	.5	.25
Load	20%	5	1.0	1.0	1.0
Store	10%	3	.3	.3	.3
Branch	20%	2	.4	.2	.4
			$\Sigma = 2.2$	2.0	1.95

(2) How does this CPI compare with using branch prediction to save a cycle off the branch time?
 CPU time new = 2.0 × IC × CC so 2.2 versus 2.0 means 10% faster

(3) What if two ALU instructions could be executed at once? CPU time new =  $1.95 \times IC \times CC$  so 2.2 versus 1.95 means 12.8% faster

# Understanding Program Performance

#### CPU Time = Instruction\_count × CPI × clock\_cycle

 $\label{eq:cpu_relation} \operatorname{CPUTime} = \frac{\operatorname{Instructions}}{\operatorname{Program}} \times \frac{\operatorname{Clock}\operatorname{cycles}}{\operatorname{Instruction}} \times \frac{\operatorname{Seconds}}{\operatorname{Clock}\operatorname{cycle}}$ 

The performance of a program depends on the algorithm, the language, the compiler, the architecture, and the actual hardware.

	Instruction_count	CPI	clock_cycle
Algorithm	Х	Х	
Programming language	Х	Х	
Compiler	Х	Х	
ISA	Х	Х	Х
Processor organization		Х	Х

# Check Yourself

A given application written in Java runs 15 seconds on a desktop processor. A new Java compiler is released that requires only 0.6 as many instructions as the old compiler. Unfortunately, it increases the CPI by 1.1. How fast can we expect the application to run using this new compiler? Pick the right answer from the three choices below:

a. 
$$\frac{15 \times 0.6}{1.1} = 8.2 \text{ sec}$$

b. 
$$15 \times 0.6 \times 1.1 = 9.9$$
 sec

c. 
$$\frac{15 \times 1.1}{0.6} = 27.5 \text{ sec}$$

# Check Yourself

A given application written in Java runs 15 seconds on a desktop processor. A new Java compiler is released that requires only 0.6 as many instructions as the old compiler. Unfortunately, it increases the CPI by 1.1. How fast can we expect the application to run using this new compiler? Pick the right answer from the three choices below:

a. 
$$\frac{15 \times 0.6}{1.1} = 8.2 \text{ sec}$$

b. 
$$15 \times 0.6 \times 1.1 = 9.9$$
 sec

c. 
$$\frac{15 \times 1.1}{0.6} = 27.5 \text{ sec}$$

#### b!

$$Time_1 = IC_1 \times CPI_1 \times CC_1$$
  

$$Time_2 = IC_2 \times CPI_2 \times CC_2$$
  

$$= (IC_1 \times 0.6) \times (CPI_1 \times 1.1) \times CC_1$$
  

$$= Time_1 \times 0.6 \times 1.1$$
  

$$= 15 \times 0.6 \times 1.1 = 9.9$$

# Outline

#### 1 The Basics

#### 2 Clock Cycles per Instruction (CPI)

- 3 Power, Trends, Limitations
- 4 Benchmarks and Profiling

# **CPU Power Usage**

Depending on an architect's design goals they may want to look at metrics different from latency, throughput, time, or clock frequency.

 $\begin{array}{rcl} \mbox{Power Usage} & \Longrightarrow & \mbox{Temperature} \\ \mbox{Power Usage} & \Longrightarrow & \mbox{Battery Life} \end{array}$ 

Ultrabooks are so (not) hot right now



### Power Trends



Complementary metal oxide semiconductor (CMOS) integrated circuits, the technology which implements the physical circuity inside modern CPUs, has a (simplified) power equation:

Power = Capacitive load × Voltage<sup>2</sup> × Frequency switched  
(×30) 
$$(5V \rightarrow 1V)$$
 (×1000)

# Reducing Power

#### Suppose a new CPU has

- $\, {\scriptstyle {\scriptstyle \vdash}}\,$  85% of capacitive load of old CPU
- $\, {\scriptstyle {\scriptstyle \vdash}} \,$  15% voltage and 15% frequency reduction

$$\frac{P_{new}}{P_{old}} = \frac{C_{old} \times 0.85 \times (V_{old} \times 0.85)^2 \times F_{old} \times 0.85}{C_{old} \times V_{old}^2 \times F_{old}} = 0.85^4 = 0.52$$

#### The Power Wall

- □→ Required input voltage may decrease but more transistors use more power.

### CPU Performance: Relative Performance vs VAX-11



Alex Brandt

Thursday January 10, 2019

### Multi-Processors to the Rescue

Moore's Law failing? Great idea #4: Performance via Parallelism.

- Multi-core processors!
  - ${\,\sqsubseteq\,}$  More than one processor per chip
  - $\, {\scriptstyle {\scriptstyle \vdash}}\,$  Huge benefits to OS and multiple processes.
- Within a single process requires explicit parallel programming
  - - Hardware executes multiple instructions at once (pipelining/multi-issue)
    - Hidden from the programmer
  - $\vdash$  Hard to do:
    - Programming for performance
    - Thread management, Load Balancing
    - Optimizing communication and synchronization

# Outline

#### 1 The Basics

- 2 Clock Cycles per Instruction (CPI)
- 3 Power, Trends, Limitations
- 4 Benchmarks and Profiling

# SPEC CPU Benchmark

Programs used to measure performance

- $\, {\scriptstyle {\scriptstyle {\scriptstyle \vdash}}} \,$  Supposedly typical of actual workload
- Standard Performance Evaluation Corp (SPEC)
- SPEC CPU2006
  - - Negligible I/O, so focuses on CPU performance
  - → Normalize relative to reference machine
  - $\, {\scriptstyle {\scriptstyle {\scriptstyle \leftarrow}}}\,$  Summarize as geometric mean of performance ratios
    - CINT2006 (integer) and CFP2006 (floating-point)

$$\sqrt[n]{\prod_{i=1}^{n} \text{Execution time ratio}_{i}}$$

# CINT2006 for Intel Core i7 920

Description	Name	Instruction Count x 10 <sup>9</sup>	CPI	Clock cycle time (seconds x 10 <sup>-9</sup> )	Execution Time (seconds)	Reference Time (seconds)	SPECratio
Interpreted string processing	peri	2252	0.60	0.376	508	9770	19.2
Block-sorting compression	bzip2	2390	0.70	0.376	629	9650	15.4
GNU C compiler	gcc	794	1.20	0.376	358	8050	22.5
Combinatorial optimization	mcf	221	2.66	0.376	221	9120	41.2
Go game (AI)	go	1274	1.10	0.376	527	10490	19.9
Search gene sequence	hmmer	2616	0.60	0.376	590	9330	15.8
Chess game (AI)	sjeng	1948	0.80	0.376	586	12100	20.7
Quantum computer	libquantum	659	0.44	0.376	109	20720	190.0
simulation							
Video compression	h264avc	3793	0.50	0.376	713	22130	31.0
Discrete event simulation library	omnetpp	367	2.10	0.376	290	6250	21.5
Games/path finding	astar	1250	1.00	0.376	470	7020	14.9
XML parsing	xalancbmk	1045	0.70	0.376	275	6900	25.1
Geometric mean	-	-	-	-	-	-	25.7

Alex Brandt

# **Profiling Tools**

- Many profiling tools
  - $\, \, \downarrow \, \,$  gprof (static instrumentation)
  - □ → cachegrind, Dtrace (dynamic instrumentation)
  - $\vdash$  perf (performance counters)

perf in linux-tools, based on event sampling

- ↓ Keep a list of where "interesting events" (cycle, cache miss, etc) happen
- $\, {\scriptstyle {\scriptstyle {\scriptstyle \leftarrow}}} \,$  CPU Feature: Counters for hundreds of events
  - Performance: Cache misses, branch misses, instructions per cycle, ...
- Intel®64 and IA-32 Architectures Software Developer's Manual: Appendix A lists all counters http:

//www.intel.com/products/processor/manuals/index.html

 $\rightarrow$  perf user guide:

https://perf.wiki.kernel.org/index.php/Tutorial

### Exercise 1

copymatrix1 vs copymatrix2

- $\lor$  What do they do?
- $\lor$  What is the difference?
- $\, \, \downarrow \, \,$  Which one performs better? Why?

```
perf stat -e cycles -e cache-misses ./copymatrix1
perf stat -e cycles -e cache-misses ./copymatrix2
```

- $\, \, \downarrow \, \,$  What does the output like?
- $\vdash$  How to interpret it?
- $\, \, \downarrow \, \,$  Which program performs better?

### Exercise 1: Results

abrandt5@node04:~/test-code\$ gcc -o copy.bin copymatrix.c abrandt5@node04:~/test-code\$ perf stat -e cycles -e cache-misses ./copy.bin

Performance counter stats for './copy.bin':

557,458,379 cycles 2,071,981 cache-misses

0.210788972 seconds time elapsed

```
abrandt5@node04:~/test-code$ vi copymatrix.c
abrandt5@node04:~/test-code$ gcc -o copy.bin copymatrix.c
abrandt5@node04:~/test-code$ perf stat -e cycles -e cache-misses ./copy.bin
```

Performance counter stats for './copy.bin':

1,418,374,698 cycles 3,980,687 cache-misses

0.534273226 seconds time elapsed

### Exercise 2

#### lower1 vs lower2

- $\lor$  What do they do?
- $\lor$  What is the difference?
- $\, \, \downarrow \, \,$  Which one performs better? Why?

```
perf stat -e cycles -e cache-misses ./lower1
perf stat -e cycles -e cache-misses ./lower2
```

- $\, \, \downarrow \, \,$  What does the output like?
- $\vdash$  How to interpret it?
- $\, \, \downarrow \, \,$  Which program performs better?

### Exercise 2: Results

```
abrandt5@node04:~/test-code$ gcc -o lower.bin lower.c
abrandt5@node04:~/test-code$ perf stat -e cycles -e cache-misses ./lower.bin
 Performance counter stats for './lower.bin':
  131,658,606,345 cvcles
           19,799 cache-misses
     43.195108384 seconds time elapsed
abrandt5@node04:~/test-code$ vi lower.c
abrandt5@node04:~/test-code$ gcc -o lower.bin lower.c
abrandt5@node04:~/test-code$ perf stat -e cycles -e cache-misses ./lower.bin
 Performance counter stats for './lower.bin':
       24,846,861 cycles
           12,013 cache-misses
      0.010991759 seconds time elapsed
```

# CS3350B Computer Organization Chapter 1: CPU and Memory Part 2: The Memory Hierarchy

#### Alex Brandt

Department of Computer Science University of Western Ontario, Canada

Tuesday January 15, 2019

# Recap: CPU Time

 $CPUTime = Instruction\_count \times CPI \times clock\_cycle$ 

	Instruction_count	CPI	clock_cycle
Algorithm	X	Х	
Programming language	Х	Х	
Compiler	Х	Х	
ISA	X	Х	Х
Processor organization		Х	Х

From a programmer's point of view CPU performance depends on:

- CPU Frequency
- The type of instructions performed
- Memory access time

### Recap: Processor-Memory Gap

# The CPU-Memory Gap

# The increasing gap between DRAM, disk, and CPU speeds.



The **Processor-Memory Gap** is a key contributor to the **Memory Wall** – the point where a program's performance is totally determined by memory speed. Faster processors will not make programs run faster!

# Recap: Memory Wall Example

```
void copymatrix1(int** src,
                                       void copymatrix2(int** src,
int** dst, int n) {
                                       int** dst, int n) {
  int i,j;
                                          int i,j;
  for (i = 0; i < n; i++)
                                          for (j = 0; j < n; j++)
    for (j = 0; j < n; j++)
                                            for (i = 0; i < n; i++)
       dst[i][j] = src[i][j];
                                              dst[i][j] = src[i][j];
}
                                       }
                                        Performance counter stats for './copy.bin':
Performance counter stats for './copy.bin':
                                                          cycles
                                           1,418,374,698
     557.458.379 cvcles
                                              3,980,687
                                                          cache-misses
      2.071.981 cache-misses
                                             0.534273226 seconds time elapsed
     0.210788972 seconds time elapsed
```

Cache misses are the reason for the performance drop.

# Outline

- 1 History, Trends, and Basics
- 2 The Principle of Locality
- 3 The Cache Hierarchy
- 4 Cache and Memory Performance
- 5 Locality and Cache

# Von Neumann Architecture



A simple shared memory for both instructions and data.

- From a software's point of view this model is still very much true.
- Hardware takes care of all of the sophistication.
#### Before the Hierarchy

In the 1950s, 60s, and 70s, software development was mostly mathematical and scientific.

- The heyday of FORTRAN.
- Computing resources were very limited.
  - $\, {\scriptstyle {\scriptstyle \vdash}}\,$  As late as the 1980s memory was only 64 mega bytes.
- Algorithm development focused on minimizing amount of memory a program used in order to solve larger problems.
  - $\, {\scriptstyle \downarrow} \,$  Sparse mathematics, linear algebra, symbolic computations.
  - Minimizing storage of numerical types: char/INT1 (1 byte), short/INT2 (2 bytes), int/INT4 (4 bytes), ...

## Current Memory Resources

- Historically, programmers where concerned with *how much* memory was used.
  - $\, {\scriptstyle {\scriptstyle {\scriptstyle \leftarrow}}}\,$  Speeds of processor and memory were about the same.
  - $\, {\scriptstyle {\scriptstyle {\scriptstyle \leftarrow}}}\,$  Memory access time was roughly same as arithmetic time.
- Today, we are for more concerned with *how* memory is used.
- Memory size is no longer limiting factor:

  - ${}_{\rightarrow}$  Compute clusters  $\implies$  hundreds of gigabytes *per processor*.
- Memory Wall (memory speed) is now the limiting factor.
- Memory Hierarchy introduced to tackle Memory Wall and limit the effects of the Processor-Memory Gap.

## The Memory Wall

- Processor vs DRAM speed disparity continues to grow
- Notice log scale!



Hennessy & Patterson, Computer Architecture: A Quantitative Approach

## Components of a Computer: Memory



At first, cache was a single, simple entity. Around the 1990s **multi-level** caches began.

## Memory Hierarchy



## But Why a Hierarchy?

Enduring properties of hardware and software:

- Fast storage technologies (SRAM) cost more per byte.
- Fast storage technologies take up more physical space per byte.
- Gap between CPU and main memory (DRAM) speed is widening.
- Principle of Locality.

Whole sections of memory are copied from lower levels of hierarchy to higher levels.

Assuming locality, this transfer is done infrequently and computations can occur continuously by only referencing data high in the hierarchy.

## Outline

- 1 History, Trends, and Basics
- 2 The Principle of Locality
- 3 The Cache Hierarchy
- 4 Cache and Memory Performance
- 5 Locality and Cache

## The Principle of Locality

"Programs tend to reuse data and instructions they have recently used."

→ Hennessy and Patterson, *Computer Architecture*.

During a short time range, a program is likely to access only a relatively small portion of:

- available address space,
- available program code.

This principle is a driving factor in computer architecture and design.

- Always in an effort to improve performance.
- Locality directs hardware design. If your program breaks this principle then it will not use the hardware properly.

# Locality in Time and Space

#### **Temporal Locality**

Recently accessed items are likely to be accessed again in the near future.

#### **Spatial Locality**

Recently accessed items are likely to have their adjacent (or near-by) items accessed in the near future.

Note: items can mean data memory addresses or instructions.

Appreciating locality for instructions is more difficult than for memory addresses. What programming constructs lead to locality for instructions?

Temporal?

Spatial?

Appreciating locality for instructions is more difficult than for memory addresses. What programming constructs lead to locality for instructions?

- Temporal?
  - ${\bf {\scriptstyle {\rm b}}} \ \ Loops$
- Spatial?

Appreciating locality for instructions is more difficult than for memory addresses. What programming constructs lead to locality for instructions?

- Temporal?
  - $\vdash$  Loops
  - $\, {\scriptstyle {\scriptstyle {\scriptstyle \vdash}}}\,$  Repeated calls to the same method
- Spatial?

Appreciating locality for instructions is more difficult than for memory addresses. What programming constructs lead to locality for instructions?

- Temporal?
  - $\, {\scriptstyle {\scriptstyle {\scriptstyle \vdash}}} \ \, Loops$
- Spatial?
  - $\, {\scriptstyle {\scriptstyle {\scriptstyle \vdash}}} \ \, Loops$

Appreciating locality for instructions is more difficult than for memory addresses. What programming constructs lead to locality for instructions?

- Temporal?
  - $\, {\scriptstyle {\scriptstyle {\scriptstyle \vdash}}} \ \, Loops$
  - $\, {\scriptstyle {\scriptstyle {\scriptstyle \vdash}}}\,$  Repeated calls to the same method
- Spatial?
  - $\, {\scriptstyle {\scriptstyle {\scriptstyle \vdash}}} \ \, Loops$
  - $\, {\scriptstyle {\scriptstyle {\scriptstyle \vdash}}}\,$  Repeated calls to the same method

Appreciating locality for instructions is more difficult than for memory addresses. What programming constructs lead to locality for instructions?

- Temporal?
  - $\, {\scriptstyle {\scriptstyle {\scriptstyle \vdash}}} \ \, Loops$
  - $\, {\scriptstyle {\scriptstyle {\scriptstyle \vdash}}}\,$  Repeated calls to the same method
- Spatial?
  - $\, {\scriptstyle {\scriptstyle {\scriptstyle \vdash}}} \ \, Loops$

  - → Sequential operation, no unconditional jumps!

Appreciating locality for instructions is more difficult than for memory addresses. What programming constructs lead to locality for instructions?

- Temporal?
  - $\, {\scriptstyle {\scriptstyle {\scriptstyle \vdash}}} \ \, Loops$
  - $\, {\scriptstyle {\scriptstyle {\scriptstyle \vdash}}}\,$  Repeated calls to the same method
- Spatial?
  - $\, {\scriptstyle {\scriptstyle {\scriptstyle \vdash}}} \ \, Loops$

  - → Sequential operation, no unconditional jumps!
    - Unless *inlined*, method calls can act like unconditional jumps.

## Locality in Data

Data locality easier to understand:

- Repeated access to the same variable.
- Sequential access to an array.
- Initializing a variable just before it is used.
- Re-use previous variables that are finished being useful instead of initializing new ones.
  - → When a new variable is created, it must be loaded into cache. Why not re-use existing variables?

Prime Example: Stride-1 access pattern

```
int arr[n];
for (int i = 0; i < n; ++i) {
    printf("%d", arr[i]);
}
```

#### Aside: Layout of C Arrays in Memory

C arrays allocated in row-major order.

$$\begin{bmatrix} 0 & 1 & 2 \\ 3 & 4 & 5 \\ 6 & 7 & 8 \end{bmatrix} \Longrightarrow \begin{bmatrix} 0, 1, 2, 3, 4, 5, 6, 7, 8 \end{bmatrix}$$

- FORTRAN, Matlab arrays allocated in column-major order.
  - $\, {\scriptstyle {\scriptstyle {\scriptstyle \leftarrow}}}\,$  Each column in contiguous memory locations.

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix} \Longrightarrow \begin{bmatrix} 1, 4, 7, 2, 5, 8, 3, 6, 9 \end{bmatrix}$$

Does this function in C have good locality? If yes, which type?

```
int sumArray(int* a, int n) {
    int sum = 0;
    for (i=0; i<n; i++)
        sum += a[i];
    return sum;
}</pre>
```

Does this function in C have good locality? If yes, which type?

```
int sumArray(int* a, int n) {
    int sum = 0;
    for (i=0; i<n; i++)
        sum += a[i];
    return sum;
}</pre>
```

■ stride-1 access to a ⇒ spatial locality

Does this function in C have good locality? If yes, which type?

```
int sumArray(int* a, int n) {
    int sum = 0;
    for (i=0; i<n; i++)
        sum += a[i];
    return sum;
}</pre>
```

- stride-1 access to a ⇒ spatial locality
- temporal locality in access to sum, i, n

Does this function in C have good locality? If yes, which type?

```
int sumArray(int* a, int n) {
    int sum = 0;
    for (i=0; i<n; i++)
        sum += a[i];
    return sum;
}</pre>
```

- stride-1 access to a ⇒ spatial locality
- temporal locality in access to sum, i, n

*Note:* Temporal locality in reference to the *pointer* a but *not* the array elements. a[i] is just syntactic sugar.

 $\vdash a[i] \iff *(a + i)$ 

Where is locality present in this C function? For each specify the type of locality.

int sumarray(int[][] a, int M, int N) {
 int i, j, sum = 0;
 for (i = 0; i < M; i++)
 for (j = 0; j < N; j++)
 sum += a[i][j];
 return sum;
}</pre>

Where is locality present in this C function? For each specify the type of locality.

```
int sumarray(int[][] a, int M, int N) {
    int i, j, sum = 0;
    for (i = 0; i < M; i++)
        for (j = 0; j < N; j++)
            sum += a[i][j];
    return sum;
}</pre>
```

Spatial locality in access to a, *row-major order in C*.

Where is locality present in this C function? For each specify the type of locality.

```
int sumarray(int[][] a, int M, int N) {
    int i, j, sum = 0;
    for (i = 0; i < M; i++)
        for (j = 0; j < N; j++)
            sum += a[i][j];
    return sum;
}</pre>
```

- Spatial locality in access to a, *row-major order in C*.
- Temporal locality in access to sum, i, j, N.

Where is locality present in this C function? For each specify the type of locality.

```
int sumarray(int[][] a, int M, int N) {
    int i, j, sum = 0;
    for (j = 0; j < N; j++)
        for (i = 0; i < M; i++)
            sum += a[i][j];
    return sum;
}</pre>
```

Where is locality present in this C function? For each specify the type of locality.

```
int sumarray(int[][] a, int M, int N) {
    int i, j, sum = 0;
    for (j = 0; j < N; j++)
        for (i = 0; i < M; i++)
            sum += a[i][j];
    return sum;
}</pre>
```

Temporal locality in access to sum, i, j, M.

Where is locality present in this C function? For each specify the type of locality.

```
int sumarray(int[][] a, int M, int N) {
    int i, j, sum = 0;
    for (j = 0; j < N; j++)
        for (i = 0; i < M; i++)
            sum += a[i][j];
    return sum;
}</pre>
```

- Temporal locality in access to sum, i, j, M.
- No spacial locality in access to a.

Where is locality present in this C function? For each specify the type of locality.

```
int sumarray(int[][] a, int M, int N) {
    int i, j, sum = 0;
    for (j = 0; j < N; j++)
        for (i = 0; i < M; i++)
            sum += a[i][j];
    return sum;
}</pre>
```

- Temporal locality in access to sum, i, j, M.
- No spacial locality in access to a.
- If M small possibility for temporal locality in N.

Does this C function have good spatial locality? If not, permute the loops so that spatial locality is achieved by stride-1 access.

```
int sumarray3d(int[][] [] a, int N) {
    int i, j, k, sum = 0;
```

}

## Outline

- 1 History, Trends, and Basics
- 2 The Principle of Locality
- 3 The Cache Hierarchy
- 4 Cache and Memory Performance
- 5 Locality and Cache

## Characteristics of the Memory Hierarchy



CPU looks first for data in L1, then in L2, ..., then in main memory.

Modern computers usually have three levels of cache before main memory.

# Photo of a CPU: Nehalem (First-Generation i7) Die



#### Core Area Breakdown

32KB I\$ per core 32KB D\$ per core 512KB L2\$ per core Share one 2-24MB L3\$



L3 Cache

## What is a Cache?

**Cache**: A small and fast (usually SRAM) storage device that acts as a staging area between the CPU and the larger and slower main memory.

- Fundamental idea of a memory hierarchy:
  - $\vdash$  For each k, the fast and small device at level k serves as cache for the larger and slower device at level k + 1.
  - ightarrow The data in cache at level k is always a subset of the data available at level k + 1.

#### Why do memory hierarchies work?

- Programs tend to access data at level k more often than they access data at level k + 1 (locality!)
- Storage at level k + 1 can be slower, larger, and cheaper per byte.
- Net effect: Large pool of memory that costs as little as the cheap storage near the bottom, but that serves data to programs at ≈ rate of the fast storage near the top.

# Caching in a Memory Hierarchy



- Each level of memory is partitioned into blocks of consecutive bytes.
- Each block within a level of cache is the same size.
- The smaller, faster, more expensive storage-device at level k caches a subset of the data from level k + 1
- If block-size is same between two levels of cache, the subset of data is actually a subset of blocks.
- Data is copied between levels in block-sized transfer units.
## Cache Hits and Cache Misses



The processor requests a memory address contained in a block b.

#### Cache hit (at level k)

■ The Program finds *b* in the cache at level *k*. e.g., block 14

#### Cache miss (at level k)

- b is not at level k, so the level k cache must fetch it from level k + 1.
   e.g., block 12
- If level k cache is full, then some block must be replaced (evicted) by the newly requested block. Which one is the "victim"?

### **Cache Policies**



Which one is the "victim"?

- Placement (mapping) policy: where can the new block go?

  - $\, {\, \sqsubseteq \, }$  If all slots it fits into are full, overwrite whatever is in that slot.
  - $\vdash$  e.g.,  $b \mod 4$ , a "mod-4 mapping"
- Replacement policy: which block should be evicted? Replace a block based on some policy:
  - $\, \, \downarrow \, \,$  LRU (least recently used).
  - $\vdash$  FIFO (first in, first out).

# Cache Misses Explained (1/2)

### ■ Cold (compulsory) Miss (at level k)

- ightarrow A cold miss occurs at level k for a block b when this block is missing for the first time at level k cache (i.e. it is the first request for block b).
- → Cold misses always occur.

### ■ Capacity Miss (at level k)

- ightarrow Occurs when the set of active blocks (that is, the data set with which the program is actively working on) is larger than the size of the cache at level k.
- ↓ Essentially a special kind of Conflict Miss when the cache happens to be full.

# Cache Misses Explained (2/2)

### ■ Conflict Miss (at level k)

- → Under a mapping policy caches limit the positions a block can be placed to a small subset (sometimes a singleton) of the total available positions.

- → *Note:* thrashing can occur if two blocks are fighting for the same slot.
- → e.g. Requesting blocks 0, 8, 0, 8, 0, 8, ... would cause a miss every time under a mod-4 mapping.

## Multiple Cache Levels



*Caveat:* While cache levels are numbered 1, 2, etc. it is common to call a cache at level k + 1 a *lower level cache* with respect to a cache at level k.

## Outline

- 1 History, Trends, and Basics
- 2 The Principle of Locality
- 3 The Cache Hierarchy
- 4 Cache and Memory Performance
- 5 Locality and Cache

# CPU Time with Caching

### Simplified CPU Time:

CPU Time = Instruction\_count × CPI × clock\_cycle  $\downarrow$  Here, CPI is CPI<sub>ideal</sub>.

### **CPU Time with Memory:**

Average memory stall cycles = access\_count × miss\_rate × miss\_penalty

*Note:* It is generally assumed that cache hit time is included as part of CPI for load/store instructions.

## Quantifying Cache Hits and Misses

- Hit Rate: The percentage of memory accesses which are found in a given level of the memory hierarchy.
- Hit Time: The time to retrieve data during memory access at a given level of the memory hierarchy, consisting of:

Time to determine hit/miss + Time to access/transmit the block.

- Miss Rate: The percentage of memory accesses which are not found in a given level of the memory hierarchy, that is, 1 - (Hit Rate).
- Miss Penalty: The time required to search and retrieve the requested block from a lower (and slower) level of cache.

Time to determine hit/miss

- + Time to access the block in the lower level
- + Time to transmit that block back to the current level

36 / 52

- + Time to insert the block in that level
- + Time to pass the block to the requester

### Hit Time << Miss Penalty

### Impacts of Cache Performance

 $\mathrm{CPI}_{stall} = \mathrm{CPI}_{ideal} + \mathrm{Average\ memory\ stall\ cycles}$ 

Avg. mem. stall cycles = access\_count  $\times$  miss\_rate  $\times$  miss\_penalty

- The relative cost for a cache miss increases as processor performance increases (↑ clock rate, ↓ CPI).
- For a fixed CPU (and main memory), cannot change miss\_penalty.
- Stall cycles mainly affected by program and compiler:

  - $\rightarrow$  miss\_rate: locality!
- This calculation assumes an idealized cache: one level of cache between CPU and main memory.

# CPI Example (1/2)

A program running on a particular processor has a  $\rm CPI_{ideal}$  of 2, a 100 cycle miss penalty, 36% load/store instr's, and a 4% miss rate. What is the average memory stall cycles? What is  $\rm CPI_{stall}?$ 

- Memory-stall cycles =  $36\% \times 4\% \times 100 = 1.44$
- So  $CPI_{stall} = 2 + 1.44 = 3.44$

- Memory-stall cycles =  $36\% \times 4\% \times 100 = 1.44$
- So  $CPI_{stall} = 2 + 1.44 = 3.44$

What if the previous miss rate is broken down as a 2% instruction-cache miss rate and 4% data-cache miss rate?

Memory-stall cycles =  $36\% \times 4\% \times 100 = 1.44$ 

So 
$$CPI_{stall} = 2 + 1.44 = 3.44$$

What if the previous miss rate is broken down as a 2% instruction-cache miss rate and 4% data-cache miss rate?

 $\blacksquare$  Memory-stall cycles = 1  $\times$  2%  $\times$  100 + 36%  $\times$  4%  $\times$  100 = 3.44

• So 
$$CPI_{stall} = 2 + 3.44 = 5.44$$

What if the  $\mathrm{CPI}_{\mathrm{ideal}}$  is reduced to 1?

What if the data-cache miss rate went up by 1%? (Instruction-cache miss rate still 2%)

What if the  $CPI_{ideal}$  is reduced to 1?

What if the data-cache miss rate went up by 1%? (Instruction-cache miss rate still 2%)

- Memory-stall cycles =  $1 \times 2\% \times 100 + 36\% \times 5\% \times 100 = 3.80$
- So  $CPI_{stall} = 2 + 3.80 = 5.80$

# Aside: Banked and Unified Caches

**Banked Cache:** A cache that is divided into two sections: one for instructions and one for data.

- Allows CPU to easily (and simultaneously) access both instructions and data.
- Gemember this when we talk about pipelining later.

**Unified Cache:** A cache where instructions and data are stored together and likely intermixed.

Core			
L1 D	L1 I		



12

## Another Memory Performance Metric: AMAT

Cache Miss Rate: number of cache misses total number of cache references (accesses) → Miss rate + Hit rate = 1.0 (100%)

**Miss Penalty**: the access time of a memory level one below the given memory level, where a requested block is missing.

Average Memory Access Time (AMAT) is the average time to access memory considering both hits and misses.

 $AMAT = Time for a Hit + Miss Rate \times Miss Penalty$ 

*Note:* In contrast to CPI<sub>stall</sub> this is measured in *time*.

What is the AMAT for a processor with a 200 ps clock, a miss penalty of 50 clock cycles, a miss rate of 0.02 misses per instruction and a cache access time of 1 clock cycle?

What is the AMAT for a processor with a 200 ps clock, a miss penalty of 50 clock cycles, a miss rate of 0.02 misses per instruction and a cache access time of 1 clock cycle?

 $AMAT = Time for a Hit + Miss Rate \times Miss Penalty$ 

What is the AMAT for a processor with a 200 ps clock, a miss penalty of 50 clock cycles, a miss rate of 0.02 misses per instruction and a cache access time of 1 clock cycle?

 $\begin{aligned} \mathsf{AMAT} &= \mathsf{Time for a Hit} + \mathsf{Miss Rate} \times \mathsf{Miss Penalty} \\ &= 1 + 0.02 \ ^* \ 50 = 2 \ \mathsf{clock cycles, or} \ 2 \ ^* \ 200 = 400 \ \mathsf{ps} \end{aligned}$ 

# Multiple Cache Levels (again)



## AMAT & Multiple Cache Levels

Given the memory hierarchy, speed at each level differs. Must calculate miss penalty on a *per level* basis.

Miss penalties defined per cache level:

:

L1 Miss Penalty = L2 Hit Time + L2 Miss Rate × L2 Miss Penalty L2 Miss Penalty = L3 Hit Time + L3 Miss Rate × L3 Miss Penalty L3 Miss Penalty = Main Memory Hit Time (usually)

#### New AMAT:

### New AMAT Example

Calculate AMAT given:

- 200 ps clock cycle,
- 1 cycle L1 hit time, 2% L1 miss rate,
- 5 cycle L2 hit time, 5% L2 miss rate,
- 100 cycle main memory access time.

Without L2 cache:

With L2 cache:

### New AMAT Example

Calculate AMAT given:

- 200 ps clock cycle,
- 1 cycle L1 hit time, 2% L1 miss rate,
- 5 cycle L2 hit time, 5% L2 miss rate,
- 100 cycle main memory access time.

Without L2 cache:

AMAT =  $1 + .02 \times 100 = 3$  cycles,  $3 \times 200 = 600$  ps

With L2 cache:

### New AMAT Example

Calculate AMAT given:

- 200 ps clock cycle,
- 1 cycle L1 hit time, 2% L1 miss rate,
- 5 cycle L2 hit time, 5% L2 miss rate,
- 100 cycle main memory access time.

Without L2 cache:

AMAT = 
$$1 + .02 \times 100 = 3$$
 cycles,  $3 \times 200 = 600$  ps

With L2 cache:

$$AMAT = 1 + .02 \times (5 + .05 \times 100) = 1.2$$
 cycles,  $1.2 \times 200 = 240$  ps

### Cache Memories & Relative Speeds

"Cache" Type	What Cached	Where Cached	Latency	Managed By
			(cycles)	
Registers	4/8-byte word	CPU registers	0.5	Compiler
TLB	Address	On-Chip TLB	0.5	Hardware
	translations			
L1 cache	32-byte block	On-Chip L1	1	Hardware
L2 cache	32-byte block	On/Off-Chip L2	10	Hardware
Buffer cache	Parts of files	Main memory	100	OS
Virtual Memory (Paging)	4-KB page	NVMe	100,000	Hardware+
		HDD	1,000,000	OS
Network buffer	Parts of files	Local disk	10,000,000	AFS/NFS
cache				client
Web cache	Web pages	Remote server	1,000,000,000	Web proxy
		disks		server

- The TLB (Translation lookaside buffer) "address-translation cache" caches virtual memory/physical memory mappings.
- The Andrew File System (AFS) and Network File System (NFS) are distributed file system protocols

## Outline

- 1 History, Trends, and Basics
- 2 The Principle of Locality
- 3 The Cache Hierarchy
- 4 Cache and Memory Performance
- 5 Locality and Cache

Advanced Temporal Locality

```
int sumarray(int[][] a, int M, int N) {
    int i, j, sum = 0;
    for (i = 0; i < M; i++)
        for (j = 0; j < N; j++)
            sum += a[i][j];
    return sum;
}</pre>
```

Recall: Temporal locality of M depends on size of N?

Depending on size of cache and on size of N, locality is possible.

- ↓ If sizeof(int) × (N + 6) ≤ cache size then M still in cache!
- → Assumes LRU cache policy and data-specific cache.

## Advanced Spatial Locality

In a row-major language:

Stepping through columns in one row:

```
for (i = 0; i < N; i++)
   sum += a[0][i];</pre>
```

- ightarrow If block size (B) > k bytes, exploit spatial locality Cold miss rate = k bytes / B.
- $\downarrow$  Typically k = 4 or 8 (32/64 bits) and B = 8k or B = 16k.

Stepping through rows in one column:

for (i = 0; i < n; i++)
 sum += a[i][0];</pre>

- $\downarrow$  Accesses distant elements (assuming number of columns in a is large).

49 / 52

## Advanced Spatial Locality Example

Calculate the number of cache hits and cache misses given:

- A row-major programming language with a having INT4 elements,
- Cache capacity of 4 Kilobytes, cache block is 32 bytes,
- All local variables are stored in registers, not cache.

```
int i, sum = 0, N = 16384;
for (i = 0; i < N; i++)
    sum += a[i];
```

## Advanced Spatial Locality Example

Calculate the number of cache hits and cache misses given:

- A row-major programming language with a having INT4 elements,
- Cache capacity of 4 Kilobytes, cache block is 32 bytes,
- All local variables are stored in registers, not cache.

```
int i, sum = 0, N = 16384;
for (i = 0; i < N; i++)
    sum += a[i];
```

- Stride-1 access to a.
- Since cache block is 32 bytes, each block holds 32/4 = 8 ints.
- One cold miss for each block.
- $(1/8) \times 16384 = 2048$  cache misses.
- 16384 2048 cache misses = **14336 cache hits**.

# Summary

### ■ We want a large, cheap, fast memory

### Solution: Memory Hierarchy

- $\, {\scriptstyle {\scriptstyle \vdash}}\,$  Successively lower levels contain "most used" data from next higher level.
- → Exploits temporal & spatial locality of programs.
- ightarrow Great Idea #3: Do the common case fast, worry less about the exceptions (RISC design principle).

### Challenges to programmer:

 $\, \, \downarrow \, \,$  Develop cache friendly (efficient) programs.

### Extra Reading

Being able to look at code and have a *qualitative sense of its locality* is a key skill for programmers.

Some projects driven by data locality (and other features):

- BLAS http://www.netlib.org/blas/
- FFTW, by Matteo Frigo and Steven G. Johnson http://www.fftw.org/
- M. Frigo, C. Leiserson, H. Prokop, and S. Ramachandran, Cache-Oblivious Algorithms, 1999.

Wikipedia is surprisingly well written on some of these topics.

- https://en.wikipedia.org/wiki/CPU\_cache
- https://en.wikipedia.org/wiki/Cache\_hierarchy https://en.wikipedia.org/wiki/Cache-oblivious\_algorithm

# CS3350B Computer Organization Chapter 1: CPU and Memory Part 3: Cache Implementations

#### Alex Brandt

Department of Computer Science University of Western Ontario, Canada

Thursday January 17, 2019

## Outline

### 1 The Basics

- 2 Cache Organization Schemes
- 3 Handling Cache Hits & Misses
- 4 Cache Design & Improvements

### How is the Hierarchy Managed?

*Recall:* CPU  $\leftarrow$  Registers  $\leftarrow$  L1  $\leftarrow \ldots \leftarrow$  Main Memory  $\leftarrow$  Disk

- Registers ↔ Cache Memory:
  - → The compiler (sometimes with the programmer's help) decides which values are stored in which registers.
- Caches ↔ Main Memory:
  - ↓ The cache controller (thus the hardware) handles cache memory movement.
- $\blacksquare Main Memory \leftrightarrow Disk:$ 
  - $\downarrow$  The operating system (which controls the virtual memory).
  - ↓ the TLB (thus the hardware) which assists the virtual-to-physical address mapping.
  - ↓ The programmer (who organizes the data into files), if data comes from file.
### Cache Design Questions

- Q1 How best to organize the memory blocks (a.k.a lines) inside the cache?
- Q2 To which block (line) of the cache does a given (main) memory address map?
  - → *Note:* since the cache is a subset of the main memory, multiple memory addresses can map to the same cache location.
- Q3 How do we know if a block of the main memory currently has a copy in cache?
- Q4 How do we quickly find a particular copy of main memory (memory address contents) in the cache?

# General Organization of a Cache Memory (1/2)

- Cache is an array of  $R = 2^s$  sets
- $\blacksquare$  Each set contains  $N \geq 1$  lines
- Each cache line (block) holds  $B = 2^b$  bytes of data



Cache size:  $C = B \times N \times R$  data bytes

# General Organization of a Cache Memory (2/2)

• Cache size = bytes per line × lines per set × number of sets  $C = B \times N \times R$ 

Everything is a power of 2.



Cache size:  $C = B \times N \times R$  data bytes

# Memory-Cache Mapping (Addressing Cache Memories)

- The data word at the *m*-bit address
  A is in the cache if the tag bits in one of the <valid> lines in set <set index> match <tag>
- The word contents begin at offset <block offset> bytes from the beginning of the block

#### Address Mapping:

block address =  $\langle tag \rangle || \langle set index \rangle$ set# = (block address) mod R  $\Box$  just take the "s bits" as set index.

*Note:* Main memory address space and cache size/implementation highly coupled.



### Outline

#### 1 The Basics

#### 2 Cache Organization Schemes

3 Handling Cache Hits & Misses

4 Cache Design & Improvements

# Types of Cache Organization

### **Direct-Mapped**

■ N = 1

- $\, \, \downarrow \, \,$  One line per set.

 $\bullet b = \log_2(B), \ R = C/B, \ s = \log_2(R), \ t \ (\text{tag size}) = m - s - b.$ 

**Fully Associative** 

- $\blacksquare$  R = 1 (allow a memory address to be mapped to any cache block).
- Tag is whole address except block offset.
- $b = \log_2(B)$ , N = C/B, s = 0, t = m b.

N-way set associative

- *N* is typically 2, 4, 8, or 16.
- A memory block maps to a specific set but can and can be placed in any way of that set (so there are N choices of mapping).

$$b = \log_2(B), R = C/(B \times N), s = \log_2(R), t = m - s - b.$$

### Direct-Mapped Cache Example

- Direct-Mapped  $\implies N = 1$ . Let B = 1, R = 4.
- Therefore we have a 2-bit set index. Assume a 2-bit tag  $\implies m = 4$ .
- Start with an empty cache blanks are considered invalid.

Consider the sequence of memory address accesses:



### Why Middle Bits For Set Index?

Consider a 4-line direct-mapped cache; sets are indexed as 00, 01, 10, 11:

#### 4-line Cache



- High-Order Bit Indexing
  - ⊢ Adjacent memory lines would map to same cache entry
  - $\vdash$  Poor use of locality
- Middle-Order Bit Indexing
  - to different cache lines
  - $\vdash$  Can hold C-bytes of contiguous memory in cache at one time

High-Order		N	∕lid
Bit In	dexing		Bit
<u>00</u> 00x		< <u>00</u> 00	C
<u>00</u> 01 <i>x</i>		(00 <u>01</u>	
<u>00</u> 10x		(00 <u>10</u>	(
<u>00</u> 11x		(00 <u>11</u>	
<u>01</u> 00x		(01 <u>00</u>	
<u>01</u> 01 <i>x</i>		(01 <u>01</u>	
<u>01</u> 10x		(01 <u>10</u>	
<u>01</u> 11x		(01 <u>11</u>	(
<u>10</u> 00x		(10 <u>00</u>	(
<u>10</u> 01x		(10 <u>01</u>	(
<u>10</u> 10x		(10 <u>10</u>	(
<u>10</u> 11 <i>x</i>		(10 <u>11</u>	C
<u>11</u> 00 <i>x</i>		(11 <u>00</u>	C
<u>11</u> 01 <i>x</i>		(11 <u>01</u>	
<u>11</u> 10x		(11 <u>10</u>	(

ldle-Order Indexina 1111)

In this example (0000x, 0001x, ...) with an *m*-bit address, x must be m-4 bits.

1111x

### Direct-Mapped Cache Circuit Logic

- One word (4 byte) cache lines (B = 4).
- Notice byte offset related to, but not exactly, block offset.
- $2^{10} = 1024$  sets  $\implies$  10 bits for set index.



### Direct-Mapped Cache Example 2: (B > 1)

- Let the cache line hold two words = bytes (B = 2).
- Direct-Mapped  $\implies N = 1$ . R = 2.
- Therefore 1-bit set index. Assume a 2-bit tag again.
- Start with an empty cache blanks are considered invalid.

Consider the sequence of memory address accesses:



		2 10	1155		3 11	IL	01		4 []	1155	
0	00	Mem(1)	Mem(0)	00	Mem(1)	Mem(0)		<b>V</b>	Mem(1)	Mem(Q)	4
1	00	Mem(3)	Mem(2)	00	Mem(3)	Mem(2)	[	00	Mem(3)	Mem(2)	



8 requests, 4 hits, 4 misses = 50% hit rate!

### Direct-Mapped Cache Circuit Logic 2

■ Four data words/block (B = 16).
 ■ 2<sup>8</sup> = 256 sets ⇒ 8 bits for set index.



### Block Size & Cache Performance



Miss rate goes up if the block size becomes a significant fraction of the cache size.

- $\vdash$  For a fixed cache size,  $\uparrow$  block size  $\implies \downarrow$  number of blocks.
- → Increases number of **capacity misses**.

### Direct-Mapped Cache Worst-Case

Direct-Mapped  $\implies N = 1$ . Let B = 1, R = 4.

Consider the sequence of memory addresses accesses:

 $0 = (0000), 4 = (0100), 0, 4, 0, 4, 0, 4, \dots$ 



Thrasing!

 $\, {\scriptstyle ij}$  conflict misses – two addresses map into the same cache block

16 / 37

### Thrashing Fix with Set Associativity

2-way associative cache  $\implies N = 2$ . Let B = 1, R = 2 now.

Consider the sequence of memory addresses accesses:

 $0 = (0000), 4 = (0100), 0, 4, 0, 4, 0, 4, \ldots$ 



8 requests, 2 misses = 75% hit rate

Solves thrashing of direct-mapped cache caused by conflict misses

 $\downarrow$  Now two memory locations that map to the same block can co-exist.

### Four-way Set Associative Cache Circuit Logic

•  $2^8 = 256$  sets each with four ways (each with one block of one word)



### Set Associativity Costs Extra

When a miss occurs, which way do we pick for replacement?

- Least Recently Used (LRU): the block replaced is the one that has been unused for the longest time.
  - $\vdash$  Hardware must keep track of when each way was last used relative to the other blocks in the set.
  - → For 2-way set associative, takes **one bit per set**.
  - $\downarrow$  Set the bit when a block is referenced (and reset the other way's bit).

#### First In First Out (FIFO)

 $\vdash$  Can be implemented as a circular buffer or counter.

N-way set associative cache costs:

- N comparators for LRU policy (delay and physical area on chip).
- MUX delay (selecting the proper way) before data is available.
- In a direct mapped cache, the cache block is available before the hit/miss decision/
  - → Impossible to assume a hit and recover later in set associative caches.

19 / 37

### Range of Set Associative Caches

For a **fixed size cache** (and fixed size of cache lines) each doubling in associativity (ways):

- Doubles number of blocks per set,
- Halves number of sets,
- Decreases size of set-index by 1 bit,
- Increases size of tag by 1 bit.



### Benefits of Set Associative Caches

The choice between *direct mapped* (1-way) and *set associative* depends on the cost of a miss versus the cost of implementation.



→ Largest gains are in going from direct mapped to 2-way (>20% reduction in miss rate).

### Cache Structure in Different Processors

#### What do we know so far?

	Intel Nehalem	AMD Barcelona	
L1 cache size &	32KB for each per core;	64KB for each per core;	
organization	64B blocks; Split I\$ and D\$	64B blocks; Split I\$ and D\$	
L1 associativity	4-way (I), 8-way (D) set	2-way set assoc.; LRU	
	assoc.; ~LRU replacement	replacement	
L1 write policy	write-back, write-allocate	write-back, write-allocate	
L2 cache size &	256KB per core;	512KB per core;	
organization	64B blocks; Unified	64B blocks; Unified	
L2 associativity	8-way set assoc.; ~LRU	16-way set assoc.; ~LRU	
L2 write policy	write-back, write-allocate	write-back, write-allocate	
L3 cache size &	8192KB (8MB) shared by	2048KB (2MB) shared by	
organization	cores; 64B blocks; Unified	by cores; 64B blocks; Unified	
L3 associativity	16-way set assoc.	32-way set assoc.; evict block	
		shared by fewest cores	
L3 write policy	write-back, write-allocate	write-back, write-allocate	

### Outline

#### 1 The Basics

- 2 Cache Organization Schemes
- 3 Handling Cache Hits & Misses
- 4 Cache Design & Improvements

# Handling Cache Hits

Read Hits (Instructions and Data)

• This is what we want!  $\implies$  Do nothing special.

Write Hits (Data only) has two policies:

- Write-through: Require the cache and backing memory to always be consistent.
  - → When writing to a cache, also pass the value to the next lower level cache (or main memory) to update its copy.
  - In naïve implementation this is very slow. Speed of cache writes limited by the speed of the next lower level.
  - $\, \, \downarrow \,$  Use a write buffer between levels and stall only if buffer is full.
- Write-back: Allow cache and memory to be inconsistent
  - → Write data into the cache block, this becomes a **dirty block**.
  - $\, \, \downarrow \,$  Only update the next lower level when a dirty block is evicted.
  - Gradient Structure → Requires two cycles one to check for evict/dirty and another to actually do write or again use a write buffer and use only one cycle.

# Handling Cache Misses

Read Misses (Instruction and Data)

Stall the execution, fetch block from the next lower level of memory, write ("install") it into cache, pass it to next higher level of memory (or the processor), and let processor resume.

Write Misses (Data only) stall execution and perform one of two policies:

#### Write allocate:

- $\, {\scriptstyle {\scriptstyle {\scriptstyle \vdash}}} \,$  Fetch the block from the next level in the memory hierarchy,
- ↓ Installing block and writing the updated word can be done simultaneously.
- No-Write Allocate: Skip fetching/installing block into cache, write directly into lower level of memory (or a write buffer).

then let the processor resume.

# Dealing With Cache Misses with Hardware Design

#### **Compulsory Misses:**

- Caused by cold starts, process migrations or very first references.
- Reduce impact by

#### **Capacity Misses:**

- Caused by the cache becoming full; it cannot hold all blocks referenced by the program.
- Reduce impact by

#### **Conflict Misses:**

- Caused by multiple addresses mapping to the same cache block.
- Reduce impact by

# Dealing With Cache Misses with Hardware Design

#### **Compulsory Misses:**

- Caused by **cold** starts, process migrations or very first references.
- Reduce impact by increasing block size. But this causes increased miss penalty and could increase miss rate.

#### **Capacity Misses:**

- Caused by the cache becoming full; it cannot hold all blocks referenced by the program.
- Reduce impact by increasing cache size. But this may increase access time.

#### **Conflict Misses:**

- Caused by multiple addresses mapping to the same cache block.
- Reduce impact by increase associativity or increasing cache/block size and/or increase associativity (may increase access time)
  - $\, {\scriptstyle {\scriptstyle \vdash}} \,$  Larger cache  $\implies$  more sets  $\implies$  fewer addresses map to same loc.

### Outline

#### 1 The Basics

- 2 Cache Organization Schemes
- 3 Handling Cache Hits & Misses
- 4 Cache Design & Improvements

### Reducing effects of Cache Miss

#### Reducing Cache Miss Rate $\implies$ increase cache size

With increasing technology (especially transistor size and density) there is more room for larger caches.

#### Reducing Cache Miss Penalty $\implies$ use more levels of cache

- L1 cache around for a *long* time.
- L2 cache first appeared with Intel's Celeron processors with only 128KB (1999).
- L3 was first seen in 2003 but prohibitively expensive. Became standard with Intel's Nehalem in 2008.

#### Recall: New AMAT Example

1 cycle L1 hit time, 2% L1 miss rate, 5 cycle L2 hit time, 5% L2 miss rate,100 cycle main memory access time

- Without L2 cache: AMAT = 1 + .02\*100 = 3
- With L2 cache: AMAT = 1 + .02\*(5 + .05\*100) = 1.2

### Intel Pentium Memory Hierarchy



・ロト ・ 雪 ト ・ ヨ ト

### Multilevel Cache Design

Design considerations for L1 and L2 caches are very different:

- Primary cache should focus on minimizing hit time in support of faster processor clock.
- Secondary cache(s) should focus on reducing miss penalty by reducing miss rate to main memory.
  - $\,\, \downarrow \,\,$  Larger capacity with larger block sizes.
  - → Higher levels of associativity.
- The miss penalty of the L1 cache is significantly reduced by the presence of an L2 cache.
- Hit time is less important for L2 cache than miss rate.

  - Presence of L3 cache greatly improves the situation, allowing L2 miss rate to be *slightly* less important.

It's all a balancing act.

Thursday January 17, 2019

30 / 37

# Improving Cache Performance (1/2)

#### AMAT = Time for a Hit + Miss Rate \* Miss Penalty

#### (1) Reduce the time to hit in the cache.

- → Smaller cache size, smaller block size.
- $\, \, \downarrow \, \, \text{Direct-Mapped}$
- ${\,\rightarrowtail\,}$  For writes, two possible strategies:
  - No-Write Allocate: no "hit" on cache, just write to write buffer. Makes subsequent reads tricky.
  - Write Allocate: avoid 2 cycles using write buffer to write to lower cache.

#### (2) Reduce the miss rate.

- $\downarrow$  Larger cache, Larger block size (16 64 bytes typical).
- → More flexible placement (increase associativity).
- ↓ Use a victim-cache a small buffer holding most recently discarded blocks.

# Improving Cache Performance (2/2)

#### AMAT = Time for a Hit + Miss Rate \* Miss Penalty

#### (3) Reduce the miss penalty

- $\, \downarrow \,$  Smaller blocks.
- $\, \, \downarrow \, \,$  Use a write-buffer.
- → Check the write-buffer (or the victim-cache) on a read miss: *luck!*

- $\vdash$  Faster backing store (= main memory).
- ↓ Improved memory bandwidth amount and speed of memory transfer between levels (e.g. wider buses).

# The Cache Design Space

#### It's all a balancing act.

- Several interacting dimensions
  - $\vdash$  cache size
  - $\vdash$  block size
  - $\, \, \downarrow \,$  associativity
  - $\, \vdash \, \mathsf{replacement} \, \mathsf{policy}$
  - $\, \downarrow \,$  write-through vs write-back
  - $\, \downarrow \,$  write allocation

#### The optimal choice is a compromise

- ↓ depends on access characteristics (what the program is doing).
- Simplicity often wins.



# Memory Hierarchy: Critical Aspects

#### The Principle of Locality:

- Program likely to access a relatively small portion of the address space at any instant of time.
  - → **Temporal locality:** Locality in Time.
  - → **Spatial locality:** Locality in Space.

Three major types of cache misses:

- Compulsory/Cold Misses: Sad facts of life. The *first* reference to an address/block.
- **Conflict misses**: Increase cache size and/or associativity. Thrashing!
- **Capacity misses**: It turns out size does matter :(

Cache Design Space

- Total size, Block size, Associativity (& Replacement policy).
- Write-hit policy: write-through, write-back
- Write-miss policy: (no)-write-allocate. Use write buffers.

Alex Brandt

34 / 37

- $\ensuremath{\mathbb{Q}1}$  Where can an entry be placed or found in cache?
  - → Cache Organization, Direct-Mapping, Associativity.
- Q2 Which entry should be replaced on a miss?
- Q3 What happens on a write?
  - $\, {\scriptstyle {\scriptstyle {\sqcup}}}\,$  Write hit/miss strategies: write-through, write-back, write-allocate.

**Typical Example:** Given list of memory references describe the cache's state after each reference.

# Q1: Where can an entry be placed?

	# of sets	Entries per set	
Direct mapped	# of cache lines	1	
Set associative	(# of cache lines) /	Associativity	
	associativity	(typically 2 to 16)	
Fully associative	1	# of entries	

	Location method	# of comparisons
Direct mapped	Index	1
Set associative	Index the set;	Degree of
	compare set's tags	associativity
Fully associative	Compare all entries' tags	# of entries

# Q2: Which entry should be replaced on a miss?

- Easy for direct mapped only one choice.
- Set associative or fully associative:
  - $\vdash$  Random
  - $\, \downarrow \,$  LRU (Least Recently Used)
  - $\, \downarrow \,$  FIFO (First In First Out)
- For 2-way set associative LRU is easy and ideal.
- Comparisons required for LRU can be too costly for high levels of associativity (> 4-way).