

# Memory Performance Gap

The graph illustrates the performance gap between processors and memory over time. The Y-axis represents Performance on a logarithmic scale from 1 to 100,000. The X-axis represents the Year from 1980 to 2015. The Processor performance (black line with dots) shows exponential growth, reaching approximately 10,000 by 2005 and plateauing thereafter. The Memory performance (grey line with dots) shows much slower, linear growth on the log scale, reaching approximately 10 by 2015.


Year	Processor Performance	Memory Performance
1980	1	1
1985	2	1.2
1990	10	1.5
1995	100	2
2000	1,000	3
2005	10,000	5
2010	10,000	7
2015	10,000	10

Copyright © 2013, Elsevier Inc. All rights reserved.

# Memory Hierarchy Design

- Memory hierarchy design becomes more crucial with recent multi-core processors:
  - Aggregate peak bandwidth grows with # cores:
    - Intel Core i7 can generate two references per core per clock
    - Four cores and 3.2 GHz clock
      - 25.6 billion 64-bit data references/second +
      - 12.8 billion 128-bit instruction references/second
      - = 409.6 GB/s!
  - DRAM bandwidth is only 8% of this (34.1 GB/s)
  - Requires:
    - Multi-port, pipelined caches
    - Two levels of cache per core
    - Shared third-level cache on chip

Introduction

 Elsevier

Copyright © 2019, Elsevier Inc. All rights reserved.

5

# Performance and Power

- High-end microprocessors have >10 MB on-chip cache
  - Consumes large amount of area and power budget

Introduction

MK  
McGraw-Hill Education

Copyright © 2019, Elsevier Inc. All rights Reserved

6

Memory Hierarchy Basics

Introduction

- When a word is not found in the cache, a *miss* occurs:
  - Fetch word from lower level in hierarchy, requiring a higher latency reference
  - Lower level may be another cache or the main memory
  - Also fetch the other words contained within the *block*
    - Takes advantage of spatial locality
  - Place block into cache in any location within its *set*, determined by address
    - block address MOD number of sets in cache

MK

Copyright © 2019, Elsevier Inc. All rights Reserved

7

Memory Hierarchy Basics

Introduction

$$\frac{\text{Misses}}{\text{Instruction}} = \frac{\text{Missrate} \times \text{Memory accesses}}{\text{Instruction count}} = \text{Miss rate} \times \frac{\text{Memory accesses}}{\text{Instruction}}$$
$$\text{Average memory access time} = \text{Hit time} + \text{Miss rate} \times \text{Miss penalty}$$

- Speculative and multithreaded processors may execute other instructions during a miss
  - Reduces performance impact of misses

MK

Copyright © 2019, Elsevier Inc. All rights Reserved

10

Memory Hierarchy Basics

Introduction

- $n$  sets => *n-way set associative*
  - Direct-mapped cache* => one block per set
  - Fully associative* => one set
- Writing to cache: two strategies
  - Write-through*
    - Immediately update lower levels of hierarchy
  - Write-back*
    - Only update lower levels of hierarchy when an updated block is replaced
  - Both strategies use *write buffer* to make writes asynchronous

MK

Copyright © 2019, Elsevier Inc. All rights Reserved

8

Memory Hierarchy Basics

Introduction

- Six basic cache optimizations:
  - Larger block size*
    - Reduces compulsory misses
    - Increases capacity and conflict misses, increases miss penalty
  - Larger total cache capacity to reduce miss rate*
    - Increases hit time, increases power consumption
  - Higher associativity*
    - Reduces conflict misses
    - Increases hit time, increases power consumption
  - Higher number of cache levels*
    - Reduces overall memory access time
  - Giving priority to read misses over writes*
    - Reduces miss penalty
  - Avoiding address translation in cache indexing*
    - Reduces hit time

MK

Copyright © 2019, Elsevier Inc. All rights Reserved

11

Memory Hierarchy Basics

Introduction

- Miss rate
  - Fraction of cache access that result in a miss
- Causes of misses
  - Compulsory*
    - First reference to a block
  - Capacity*
    - Blocks discarded and later retrieved
  - Conflict*
    - Program makes repeated references to multiple addresses from different blocks that map to the same location in the cache

MK

Copyright © 2019, Elsevier Inc. All rights Reserved

9

Memory Technology and Optimizations

Memory Technology and Optimizations

- Performance metrics
  - Latency is concern of cache
  - Bandwidth is concern of multiprocessors and I/O
  - Access time
    - Time between read request and when desired word arrives
  - Cycle time
    - Minimum time between unrelated requests to memory
- SRAM memory has low latency, use for cache
- Organize DRAM chips into many banks for high bandwidth, use for main memory

MK

Copyright © 2019, Elsevier Inc. All rights Reserved

12

Memory Technology and Optimizations

Memory Technology

- SRAM
  - Requires low power to retain bit
  - Requires 6 transistors/bit
- DRAM
  - Must be re-written after being read
  - Must also be periodically refreshed
    - Every ~ 8 ms (roughly 5% of time)
    - Each row can be refreshed simultaneously
  - One transistor/bit
  - Address lines are multiplexed:
    - Upper half of address: row access strobe (RAS)
    - Lower half of address: column access strobe (CAS)

MK

Copyright © 2019, Elsevier Inc. All rights Reserved

13

Memory Technology and Optimizations

Memory Optimizations

Production year	Chip size	DRAM type	Best case access time (no precharge)			Precharge needed
			RAS time (ns)	CAS time (ns)	Total (ns)	Total (ns)
2000	256M bit	DDR1	21	21	42	63
2002	512M bit	DDR1	15	15	30	45
2004	1G bit	DDR2	15	15	30	45
2006	2G bit	DDR2	10	10	20	30
2010	4G bit	DDR3	13	13	26	39
2016	8G bit	DDR4	13	13	26	39

MK

Copyright © 2019, Elsevier Inc. All rights Reserved

16

Memory Technology and Optimizations

Internal Organization of DRAM

MK

Copyright © 2019, Elsevier Inc. All rights Reserved

14

Memory Technology and Optimizations

Memory Optimizations

Standard	I/O clock rate	M transfers/s	DRAM name	MiB/s/DIMM	DIMM name
DDR1	133	266	DDR266	2128	PC2100
DDR1	150	300	DDR300	2400	PC2400
DDR1	200	400	DDR400	3200	PC3200
DDR2	266	533	DDR2-533	4264	PC4300
DDR2	333	667	DDR2-667	5336	PC5300
DDR2	400	800	DDR2-800	6400	PC6400
DDR3	533	1066	DDR3-1066	8528	PC8500
DDR3	666	1333	DDR3-1333	10,664	PC10700
DDR3	800	1600	DDR3-1600	12,800	PC12800
DDR4	1333	2666	DDR4-2666	21,300	PC21300

MK

Copyright © 2019, Elsevier Inc. All rights Reserved

17

Memory Technology and Optimizations

Memory Technology

- Amdahl:
  - Memory capacity should grow linearly with processor speed
  - Unfortunately, memory capacity and speed has not kept pace with processors
- Some optimizations:
  - Multiple accesses to same row
  - Synchronous DRAM
    - Added clock to DRAM interface
    - Burst mode with critical word first
  - Wider interfaces
  - Double data rate (DDR)
  - Multiple banks on each DRAM device

MK

Copyright © 2019, Elsevier Inc. All rights Reserved

15

Memory Technology and Optimizations

Memory Optimizations

- DDR:
  - DDR2
    - Lower power (2.5 V -> 1.8 V)
    - Higher clock rates (266 MHz, 333 MHz, 400 MHz)
  - DDR3
    - 1.5 V
    - 800 MHz
  - DDR4
    - 1-1.2 V
    - 1333 MHz
- GDDR5 is graphics memory based on DDR3

MK

Copyright © 2019, Elsevier Inc. All rights Reserved

18

Memory Technology and Optimizations

### Memory Optimizations

- Reducing power in SDRAMs:
  - Lower voltage
  - Low power mode (ignores clock, continues to refresh)
- Graphics memory:
  - Achieve 2-5 X bandwidth per DRAM vs. DDR3
    - Wider interfaces (32 vs. 16 bit)
    - Higher clock rate
      - Possible because they are attached via soldering instead of socketed DIMM modules

MK

Copyright © 2019, Elsevier Inc. All rights Reserved

19

Memory Technology and Optimizations

### Flash Memory

- Type of EEPROM
- Types: NAND (denser) and NOR (faster)
- NAND Flash:
  - Reads are sequential, reads entire page (.5 to 4 KiB)
  - 25 us for first byte, 40 MiB/s for subsequent bytes
  - SDRAM: 40 ns for first byte, 4.8 GB/s for subsequent bytes
  - 2 KiB transfer: 75 uS vs 500 ns for SDRAM, 150X slower
  - 300 to 500X faster than magnetic disk

MK

Copyright © 2019, Elsevier Inc. All rights Reserved

22

Memory Technology and Optimizations

### Memory Power Consumption

Usage Mode	Read, write, terminate power (mW)	Activate power (mW)	Background power (mW)	Total Power (mW)
Low power mode	0	0	~100	~100
Typical usage	~150	~150	~100	~400
Fully active	~250	~250	~100	~600

MK

Copyright © 2019, Elsevier Inc. All rights Reserved

20

Memory Technology and Optimizations

### NAND Flash Memory

- Must be erased (in blocks) before being overwritten
- Nonvolatile, can use as little as zero power
- Limited number of write cycles (~100,000)
- \$2/GiB, compared to \$20-40/GiB for SDRAM and \$0.09 GiB for magnetic disk
- Phase-Change/Memristor Memory
  - Possibly 10X improvement in write performance and 2X improvement in read performance

MK

Copyright © 2019, Elsevier Inc. All rights Reserved

23

Memory Technology and Optimizations

### Stacked/Embedded DRAMs

- Stacked DRAMs in same package as processor
  - High Bandwidth Memory (HBM)

The diagram illustrates two DRAM architectures. On the left, 'Vertical stacking (3D)' shows multiple layers of DRAM stacked vertically, with an xPU chip positioned below them. On the right, 'Interposer stacking (2.5D)' shows a stack of DRAM layers connected to an xPU chip via a central interposer.

MK

Copyright © 2019, Elsevier Inc. All rights Reserved

21

Memory Technology and Optimizations

### Memory Dependability

- Memory is susceptible to cosmic rays
- Soft errors: dynamic errors
  - Detected and fixed by error correcting codes (ECC)
- Hard errors: permanent errors
  - Use spare rows to replace defective rows
- Chipkill: a RAID-like error recovery technique

MK

Copyright © 2019, Elsevier Inc. All rights Reserved

24

Advanced Optimizations

- Reduce hit time
  - Small and simple first-level caches
  - Way prediction
- Increase bandwidth
  - Pipelined caches, multibanked caches, non-blocking caches
- Reduce miss penalty
  - Critical word first, merging write buffers
- Reduce miss rate
  - Compiler optimizations
- Reduce miss penalty or miss rate via parallelization
  - Hardware or compiler prefetching

MK

Copyright © 2019, Elsevier Inc. All rights Reserved

25

Way Prediction

- To improve hit time, predict the way to pre-set mux
  - Mis-prediction gives longer hit time
- Prediction accuracy
  - > 90% for two-way
  - > 80% for four-way
  - I-cache has better accuracy than D-cache
- First used on MIPS R10000 in mid-90s
- Used on ARM Cortex-A8
- Extend to predict block as well
  - “Way selection”
  - Increases mis-prediction penalty

MK

Copyright © 2019, Elsevier Inc. All rights Reserved

28

L1 Size and Associativity

Cache size	1-way	2-way	4-way	8-way
16 KB	1.0	1.2	1.6	2.9
32 KB	1.2	1.3	1.8	2.5
64 KB	1.8	1.5	1.9	2.6
128 KB	2.0	2.0	2.1	2.6
256 KB	2.5	2.5	2.5	2.9

Access time vs. size and associativity

MK

Copyright © 2019, Elsevier Inc. All rights Reserved

26

Pipelined Caches

- Pipeline cache access to improve bandwidth
  - Examples:
    - Pentium: 1 cycle
    - Pentium Pro – Pentium III: 2 cycles
    - Pentium 4 – Core i7: 4 cycles
- Increases branch mis-prediction penalty
- Makes it easier to increase associativity

MK

Copyright © 2019, Elsevier Inc. All rights Reserved

29

L1 Size and Associativity

Cache size	1-way	2-way	4-way	8-way
16 KB	1.0	1.8	3.2	6.5
32 KB	1.0	1.8	3.5	8.8
64 KB	1.2	2.0	3.8	8.8
128 KB	2.5	2.5	4.2	9.0
256 KB	4.5	4.5	5.5	9.5

Energy per read vs. size and associativity

MK

Copyright © 2019, Elsevier Inc. All rights Reserved

27

Multibanked Caches

- Organize cache as independent banks to support simultaneous access
  - ARM Cortex-A8 supports 1-4 banks for L2
  - Intel i7 supports 4 banks for L1 and 8 banks for L2
- Interleave banks according to block address

Block address	Bank 0	Block address	Bank 1	Block address	Bank 2	Block address	Bank 3
0		1		2		3	
4		5		6		7	
8		9		10		11	
12		13		14		15	

MK

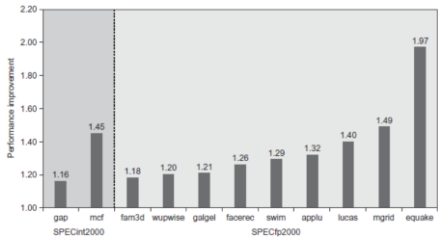
Copyright © 2019, Elsevier Inc. All rights Reserved

30



### Hardware Prefetching

- Fetch two blocks on miss (include next sequential block)



Pentium 4 Pre-fetching

Advanced Optimizations

### Use HBM to Extend Hierarchy

- Another approach (Alloy cache):
  - Mold tag and data together
  - Use direct mapped
- Both schemes require two DRAM accesses for misses
  - Two solutions:
    - Use map to keep track of blocks
    - Predict likely misses

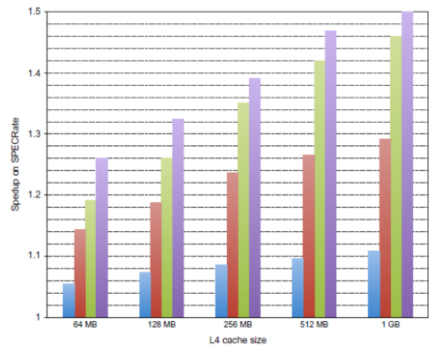
Advanced Optimizations

### Compiler Prefetching

- Insert prefetch instructions before data is needed
- Non-faulting: prefetch doesn't cause exceptions
- Register prefetch
  - Loads data into register
- Cache prefetch
  - Loads data into cache
- Combine with loop unrolling and software pipelining

Advanced Optimizations

### Use HBM to Extend Hierarchy



Advanced Optimizations

### Use HBM to Extend Hierarchy

- 128 MiB to 1 GiB
- Smaller blocks require substantial tag storage
- Larger blocks are potentially inefficient
- One approach (L-H):
  - Each SDRAM row is a block index
  - Each row contains set of tags and 29 data segments
  - 29-set associative
  - Hit requires a CAS

Advanced Optimizations

### Summary

Technique	Hit time	Bandwidth	Miss penalty	Miss rate	Power consumption	Hardware cost/complexity	Comment
Small and simple caches	+			-	+	0	Trivial; widely used
Way-predicting caches	+				+	1	Used in Pentium 4
Pipelined & banked caches	-	+				1	Widely used
Nonblocking caches		+	+			3	Widely used
Critical word first and early restart			+			2	Widely used
Merging write buffer			+			1	Widely used with write through
Compiler techniques to reduce cache misses				+		0	Software is a challenge, but many compilers handle common linear algebra calculations
Hardware prefetching of instructions and data		+	+		-	2 instr., 3 data	Most provide prefetch instructions; modern high-end processors also automatically prefetch in hardware
Compiler-controlled prefetching			+	+		3	Needs nonblocking cache; possible instruction overhead in many CPUs
HBM as additional level of cache	+/-		-	+	+	3	Depends on new packaging technology. Effects depend heavily on hit rate improvements

Advanced Optimizations

Virtual Memory and Virtual Machines

- Protection via virtual memory
  - Keeps processes in their own memory space
- Role of architecture
  - Provide user mode and supervisor mode
  - Protect certain aspects of CPU state
  - Provide mechanisms for switching between user mode and supervisor mode
  - Provide mechanisms to limit memory accesses
  - Provide TLB to translate addresses

Virtual Memory and Virtual Machines

MK

Copyright © 2019, Elsevier Inc. All rights Reserved

43

Impact of VMs on Virtual Memory

- Each guest OS maintains its own set of page tables
  - VMM adds a level of memory between physical and virtual memory called “real memory”
  - VMM maintains shadow page table that maps guest virtual addresses to physical addresses
    - Requires VMM to detect guest’s changes to its own page table
    - Occurs naturally if accessing the page table pointer is a privileged operation

Virtual Memory and Virtual Machines

MK

Copyright © 2019, Elsevier Inc. All rights Reserved

46

Virtual Machines

- Supports isolation and security
- Sharing a computer among many unrelated users
- Enabled by raw speed of processors, making the overhead more acceptable
- Allows different ISAs and operating systems to be presented to user programs
  - “System Virtual Machines”
  - SVM software is called “virtual machine monitor” or “hypervisor”
  - Individual virtual machines run under the monitor are called “guest VMs”

Virtual Memory and Virtual Machines

MK

Copyright © 2019, Elsevier Inc. All rights Reserved

44

Extending the ISA for Virtualization

- Objectives:
  - Avoid flushing TLB
  - Use nested page tables instead of shadow page tables
  - Allow devices to use DMA to move data
  - Allow guest OS’s to handle device interrupts
  - For security: allow programs to manage encrypted portions of code and data

Virtual Memory and Virtual Machines

MK

Copyright © 2019, Elsevier Inc. All rights Reserved

47

Requirements of VMM

- Guest software should:
  - Behave on as if running on native hardware
  - Not be able to change allocation of real system resources
- VMM should be able to “context switch” guests
- Hardware must allow:
  - System and use processor modes
  - Privileged subset of instructions for allocating system resources

Virtual Memory and Virtual Machines

MK

Copyright © 2019, Elsevier Inc. All rights Reserved

45

Fallacies and Pitfalls

- Predicting cache performance of one program from another
- Simulating enough instructions to get accurate performance measures of the memory hierarchy
- Not delivering high memory bandwidth in a cache-based system

Virtual Memory and Virtual Machines

MK

Copyright © 2019, Elsevier Inc. All rights Reserved

48