Computer Architecture
A Quantitative Approach, Sixth Edition

# Chapter 2

## Memory Hierarchy Design
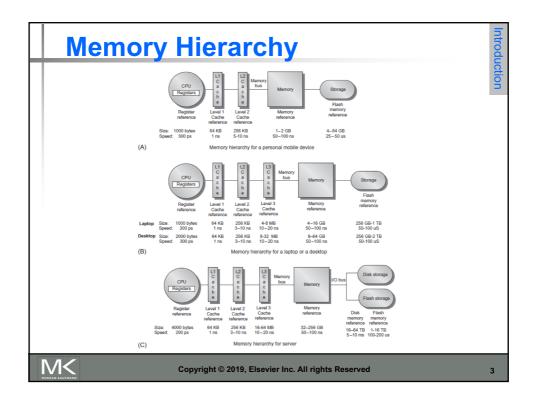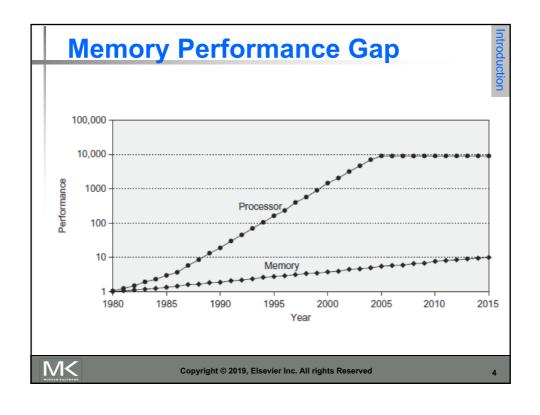
1

# Introduction

Introduction

- Programmers want unlimited amounts of memory with low latency
- Fast memory technology is more expensive per bit than slower memory
- Solution: organize memory system into a hierarchy
  - Entire addressable memory space available in largest, slowest memory
  - Incrementally smaller and faster memories, each containing a subset of the memory below it, proceed in steps up toward the processor
- Temporal and spatial locality insures that nearly all references can be found in smaller memories
  - Gives the allusion of a large, fast memory being presented to the processor

2

# Memory Hierarchy

Introduction



(A) Memory hierarchy for a personal mobile device

(B) Memory hierarchy for a laptop or a desktop

(C) Memory hierarchy for server

3

# Memory Performance Gap

Introduction

4

# 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

5

# Performance and Power

Introduction

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

6

# Memory Hierarchy Basics

- 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

# Memory Hierarchy Basics

- *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

# Memory Hierarchy Basics

- 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

9

# Memory Hierarchy Basics

$$\frac{\text{Misses}}{\text{Instruction}} = \frac{\text{Miss rate} \times \text{Memory accesses}}{\text{Instruction count}} = \text{Miss rate} \times \frac{\text{Memory accesses}}{\text{Instruction}}$$

$$\text{Average memory access time} = \text{Hit time} + \text{Miss rate} \times \text{Miss penalty}$$

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

10

# Memory Hierarchy Basics

- 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

11

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
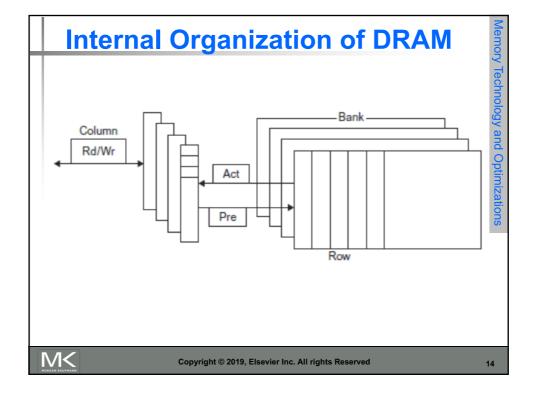
12

# 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)

13

# Internal Organization of DRAM
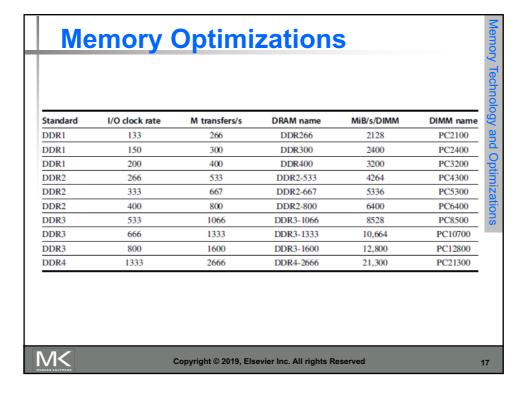
14

# 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

15

# 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 |

16

# 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 |

17

# 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

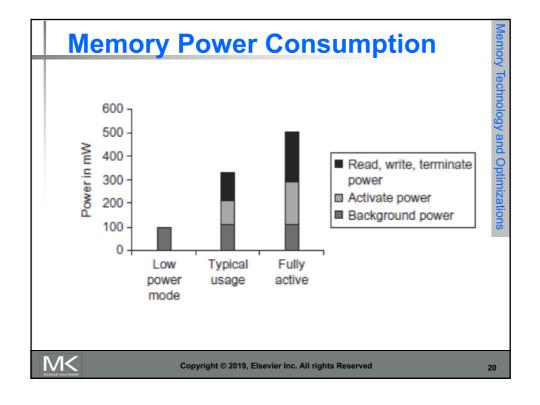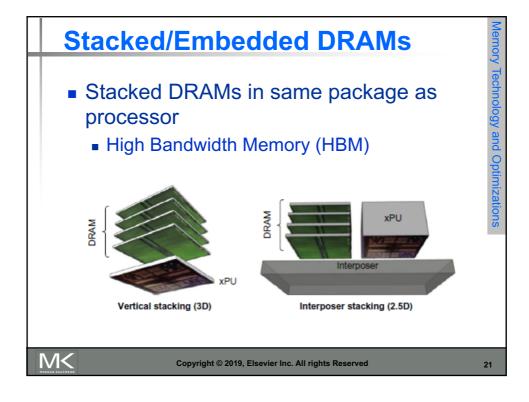- GDDR5 is graphics memory based on DDR3

18

# 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 socketted DIMM modules

Memory Technology and Optimizations

19

# Memory Power Consumption



Memory Technology and Optimizations

20

# Stacked/Embedded DRAMs

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



Vertical stacking (3D)          Interposer stacking (2.5D)

21

# 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

22

# 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/Memrister Memory
  - Possibly 10X improvement in write performance and 2X improvement in read performance

23

# 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

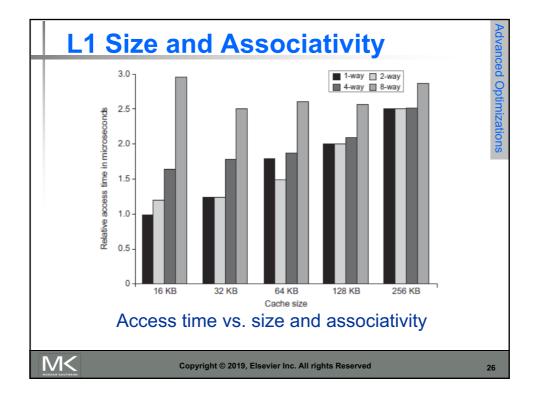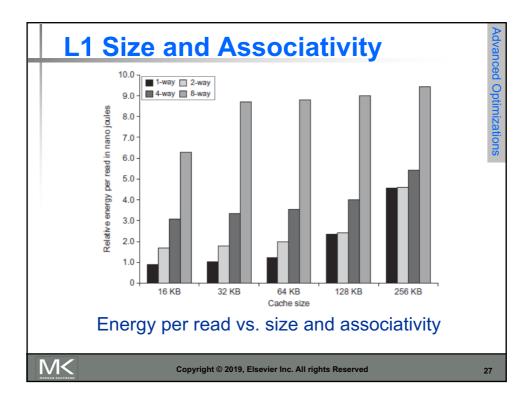- Chipkill:  a RAID-like error recovery technique

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

25

# L1 Size and Associativity

Access time vs. size and associativity

26

# L1 Size and Associativity



Energy per read vs. size and associativity

27

# 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

28

Advanced Optimizations

# 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
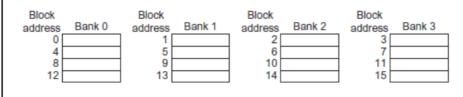- Makes it easier to increase associativity

29

# 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

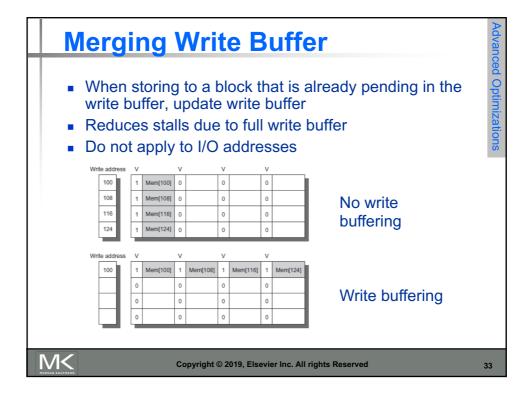| Block address | Bank 0 | Block address | Bank 1 | Block address | Bank 2 | Block address | Bank 3 |
|---|---|---|---|---|---|---|---|
| 0 | | 1 | | 2 | | 3 | |
| 4 | | 5 | | 6 | | 7 | |
| 8 | | 9 | | 10 | | 11 | |
| 12 | | 13 | | 14 | | 15 | |

30

# Nonblocking Caches

- Allow hits before previous misses complete
  - "Hit under miss"
  - "Hit under multiple miss"
- L2 must support this
- In general, processors can hide L1 miss penalty but not L2 miss penalty

31

# Critical Word First, Early Restart

- Critical word first
  - Request missed word from memory first
  - Send it to the processor as soon as it arrives
- Early restart
  - Request words in normal order
  - Send missed work to the processor as soon as it arrives

- Effectiveness of these strategies depends on block size and likelihood of another access to the portion of the block that has not yet been fetched

32

# Merging Write Buffer

- When storing to a block that is already pending in the write buffer, update write buffer
- Reduces stalls due to full write buffer
- Do not apply to I/O addresses



No write buffering
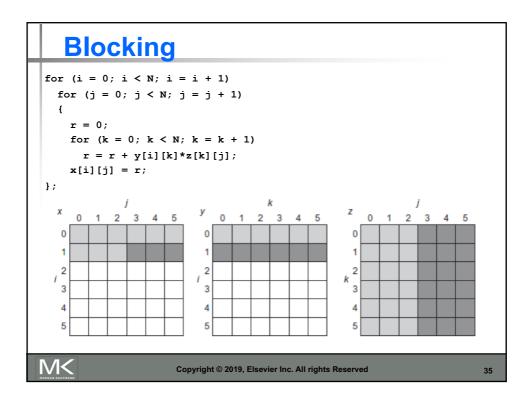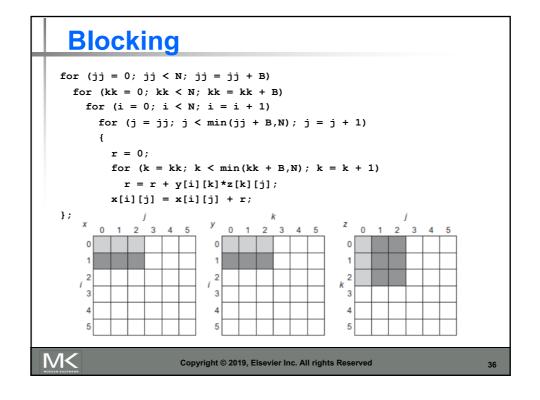
Write buffering

Advanced Optimizations

33

# Compiler Optimizations

- Loop Interchange
  - Swap nested loops to access memory in sequential order

- Blocking
  - Instead of accessing entire rows or columns, subdivide matrices into blocks
  - Requires more memory accesses but improves locality of accesses

Advanced Optimizations

34

# Blocking

```
for (i = 0; i < N; i = i + 1)
  for (j = 0; j < N; j = j + 1)
  {
    r = 0;
    for (k = 0; k < N; k = k + 1)
      r = r + y[i][k]*z[k][j];
    x[i][j] = r;
};
```

35

# Blocking

```
for (jj = 0; jj < N; jj = jj + B)
  for (kk = 0; kk < N; kk = kk + B)
    for (i = 0; i < N; i = i + 1)
      for (j = jj; j < min(jj + B,N); j = j + 1)
      {
        r = 0;
        for (k = kk; k < min(kk + B,N); k = k + 1)
          r = r + y[i][k]*z[k][j];
        x[i][j] = x[i][j] + r;
};
```

36

# Hardware Prefetching

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



Pentium 4 Pre-fetching

37

# 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

38

# 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
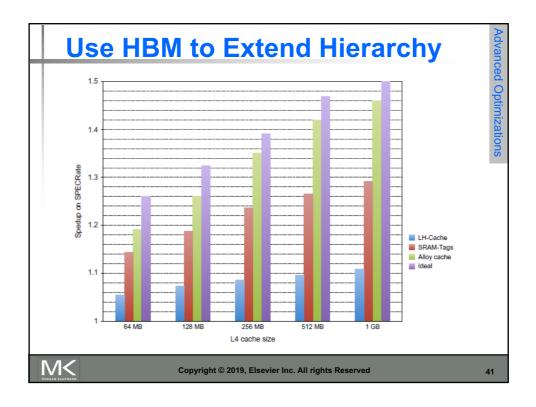  - Hit requires a CAS

39

# 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
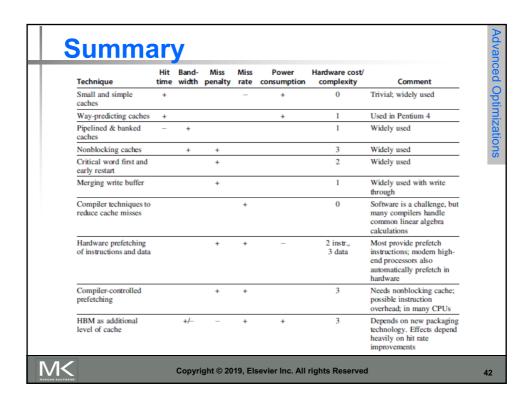    - Predict likely misses

40

# Use HBM to Extend Hierarchy

41

# Summary

| Technique | Hit time | Band-width | 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 |

42

## 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

43

Virtual Memory and Virtual Machines

## 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"

44

Virtual Memory and Virtual Machines

# 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

45

# 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

46

# 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

47

# 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

48