

CS/ECE 5381/7381  
Computer Architecture  
Spring 2023

Dr. Manikas

Computer Science

Lecture 17: Apr. 4, 2023

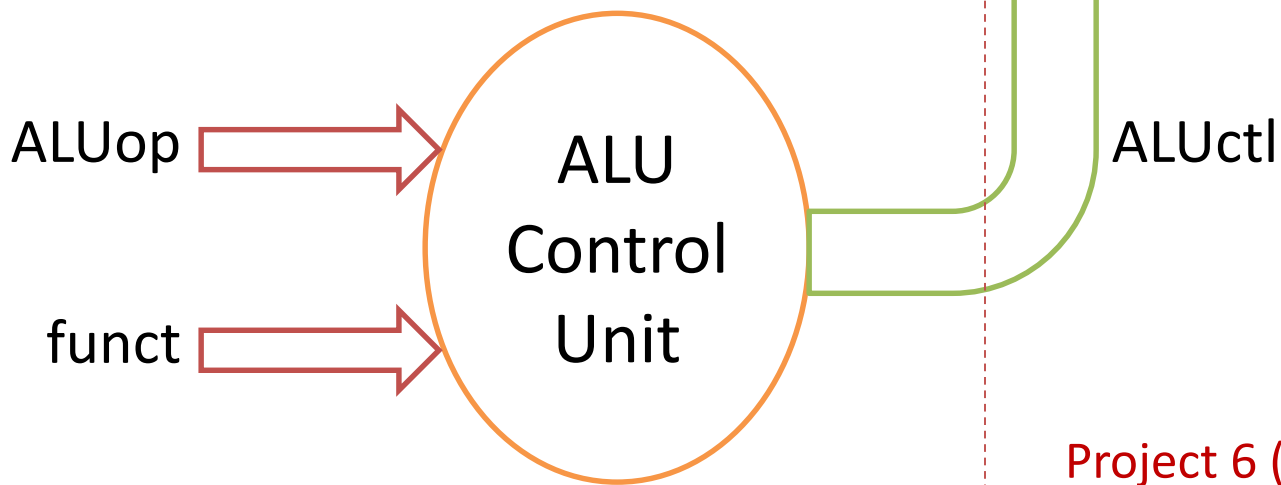
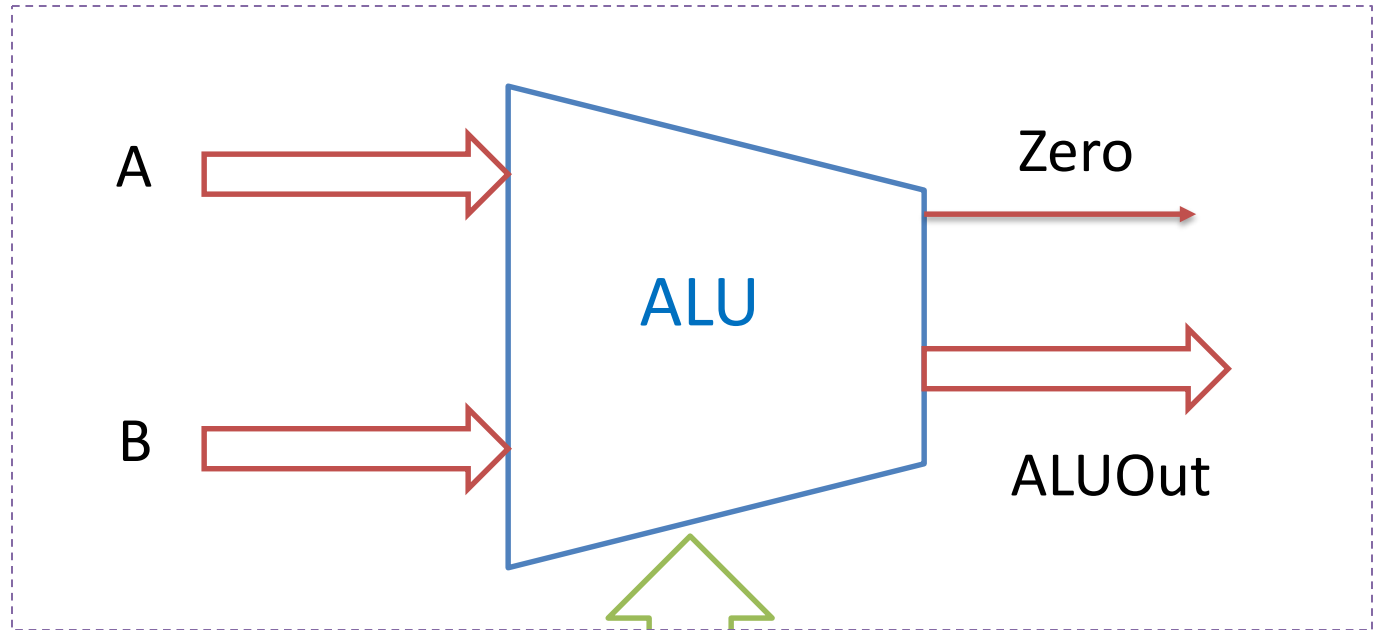
# Project 5

- Due TODAY, Apr. 4 (11:59 pm)
- Xcelium tool
- Assignment:
  - Modify Verilog code for basic MIPS ALU
    - Add arithmetic and logic functions

# Project 6 (7381 only)

- Due Thur., Apr. 6 (11:59 pm)
- For CS/ECE 7381 students ONLY
- Additional Verilog programming assignment using Xcelium tool
  - Develop control unit for ALU

## Project 5 (5/7381)



## Project 6 (7381)

# Memory Design Hierarchy

(Chapter 2, Hennessy and Patterson)

Note: some course slides adopted  
from publisher-provided material

# Outline

- 2.1 Introduction
- 2.2 Memory Technology and Optimizations
- 2.3 Ten Advanced Optimizations of Cache Performance
- 2.4 Virtual Memory and Machines

# 10 Advanced Cache Optimizations

- Reducing hit time
  1. Small and simple caches
  2. Way prediction
- Increasing cache bandwidth
  3. Pipelined caches
  4. Multibanked caches
  5. Nonblocking caches
- Reducing Miss Penalty
  6. Critical word first
  7. Merging write buffers

# 10 Advanced Cache Optimizations

- Reducing Miss Rate

## 8. Compiler optimizations



# 8. Reducing Misses by Compiler Optimizations

- *Merging Arrays*: improve spatial locality by single array of compound elements vs. 2 arrays
- *Loop Interchange*: change nesting of loops to access data in order stored in memory
- *Loop Fusion*: Combine 2 independent loops that have same looping and some variables overlap
- *Blocking*: Improve temporal locality by accessing “blocks” of data repeatedly vs. going down whole columns or rows

# Merging Arrays Example

```
/* Before: 2 sequential arrays */
```

```
int val[SIZE];
```

```
int key[SIZE];
```

```
/* After: 1 array of structures */
```

```
struct merge {
```

```
    int val;
```

```
    int key;
```

```
};
```

```
struct merge merged_array[SIZE];
```

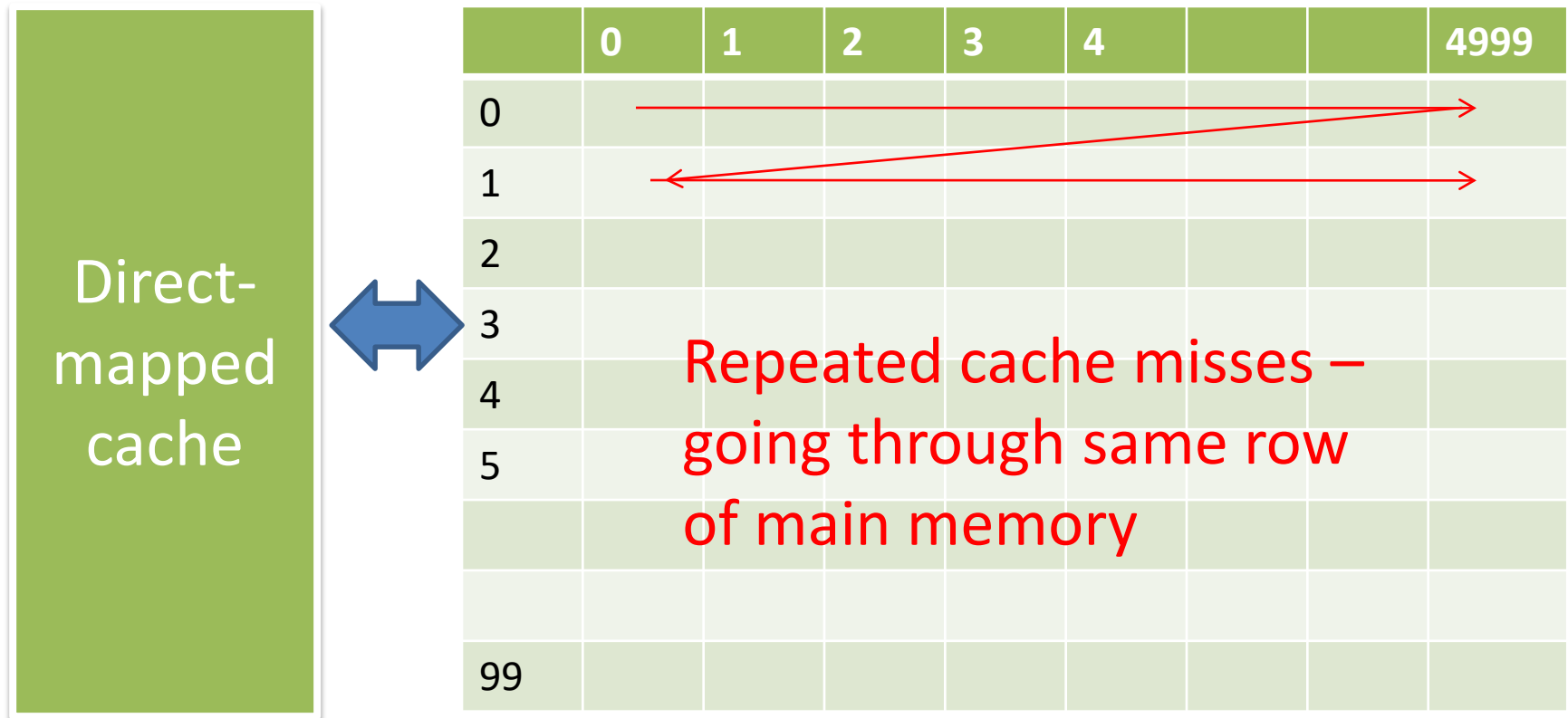
Reducing conflicts between **val** and **key**; improve spatial locality

# Loop Interchange Example

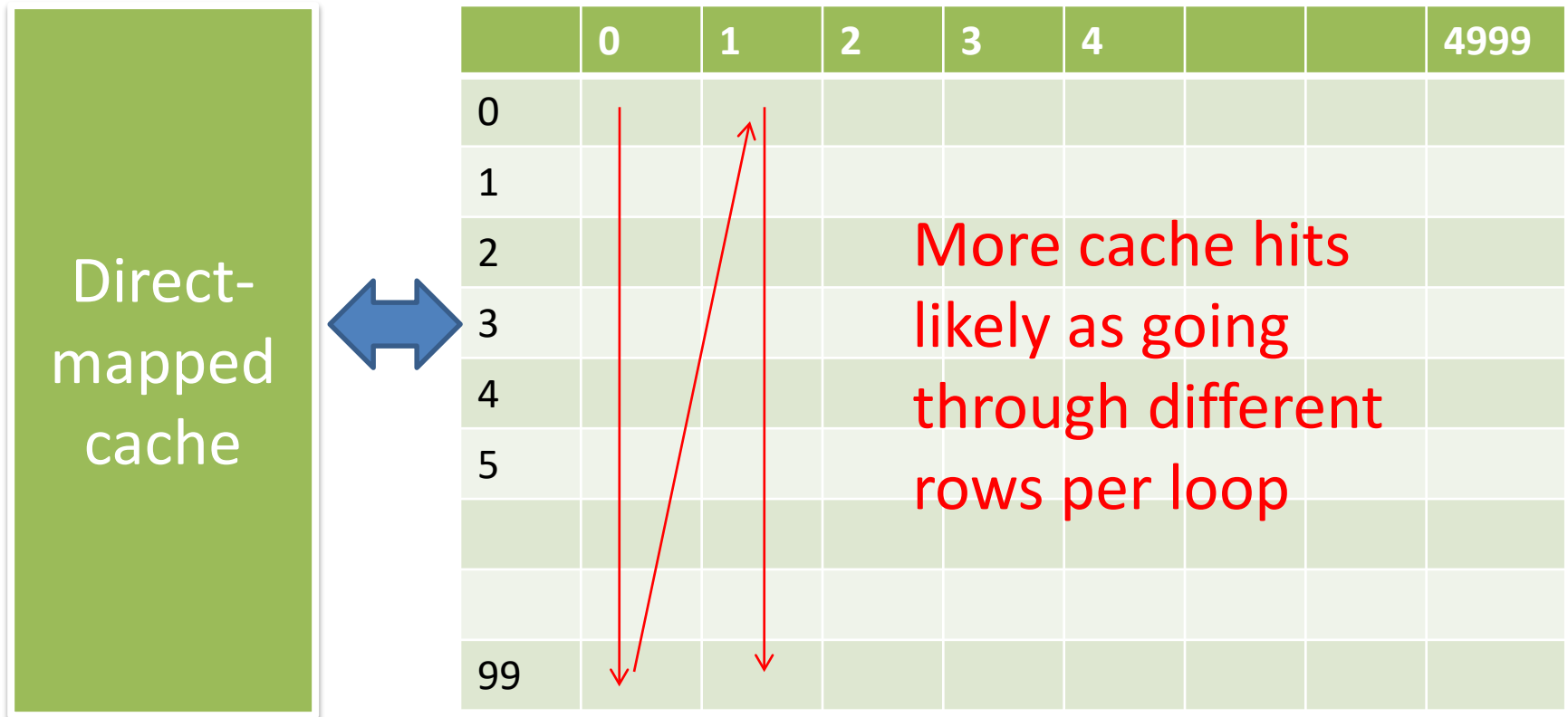
```
/* Before */  
for (k = 0; k < 100; k = k+1)  
    for (j = 0; j < 100; j = j+1)  
        for (i = 0; i < 5000; i = i+1)  
            x[i][j] = 2 * x[i][j];  
  
/* After */  
for (k = 0; k < 100; k = k+1)  
    for (i = 0; i < 5000; i = i+1)  
        for (j = 0; j < 100; j = j+1)  
            x[i][j] = 2 * x[i][j];
```

Sequential accesses instead of striding through memory every 100 words; improved spatial locality

# Before Loop Interchange



# After Loop Interchange



# Loop Fusion Example - Before

```
/* Before */
```

```
for (i = 0; i < N; i = i+1)
    for (j = 0; j < N; j = j+1)
        a[i][j] = 1/b[i][j] * c[i][j];
```

```
for (i = 0; i < N; i = i+1)
    for (j = 0; j < N; j = j+1)
        d[i][j] = a[i][j] + c[i][j];
```

# Loop Fusion Example - After

```
/* After */  
for (i = 0; i < N; i = i+1)  
    for (j = 0; j < N; j = j+1)  
        { a[i][j] = 1/b[i][j] * c[i][j];  
          d[i][j] = a[i][j] + c[i][j]; }
```

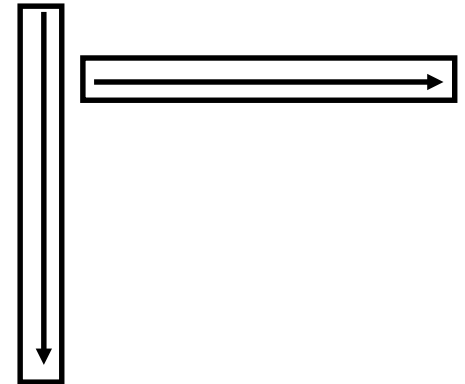
2 misses per access to a & c vs. one miss per access;  
improve spatial locality

# Blocking Example - before

```
/* Before */  
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;  
        };
```

- Two Inner Loops:

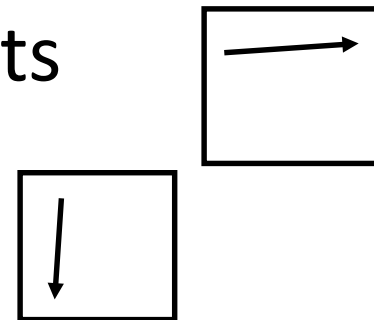
- Read all NxN elements of z[]
- Read N elements of 1 row of y[] repeatedly
- Write N elements of 1 row of x[]





# Blocking Example – before (cont)

- Capacity Misses a function of N & Cache Size:
  - $2N^3 + N^2 \Rightarrow$  (assuming no conflict; otherwise ...)
- Idea: compute on BxB submatrix that fits

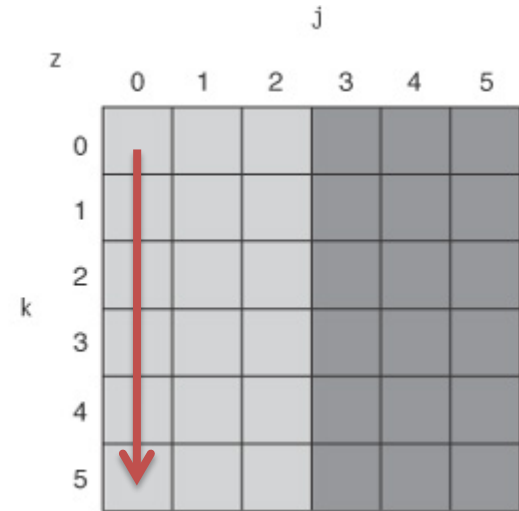
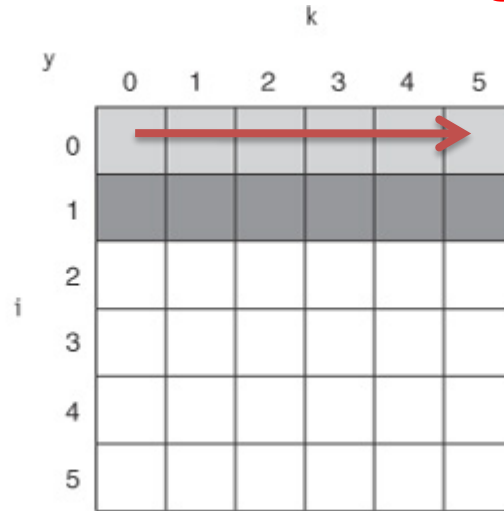
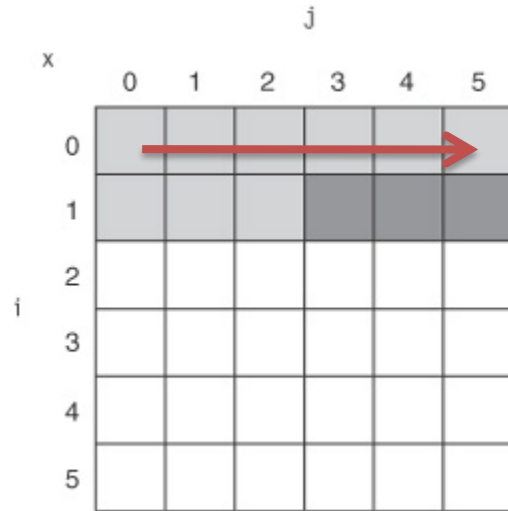


# Blocking Example

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

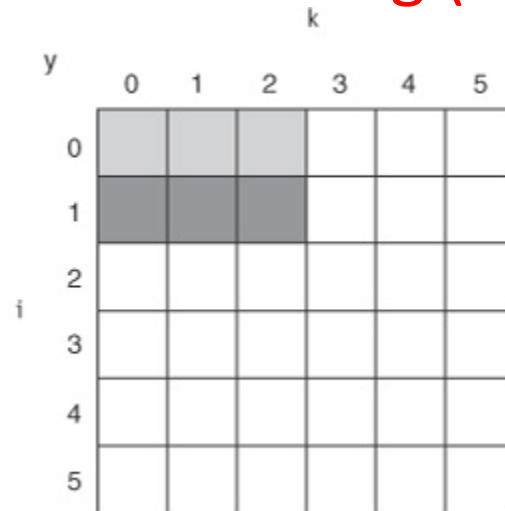
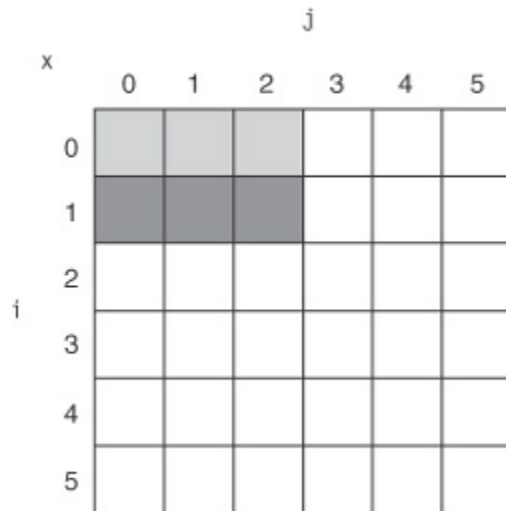
- B called *Blocking Factor*
- Capacity Misses from  $2N^3 + N^2$  to  $2N^3/B + N^2$

## Before blocking

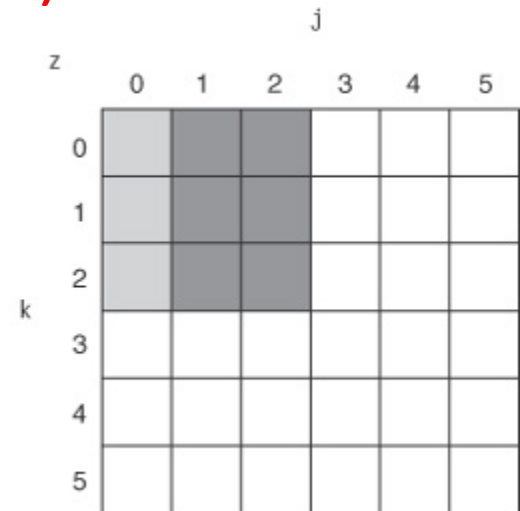


© 2007 Elsevier, Inc. All rights reserved.

## After blocking (B=3)



© 2007 Elsevier, Inc. All rights reserved.



# 10 Advanced Cache Optimizations

- Reducing miss penalty or miss rate via parallelism

9. Hardware prefetching

10. Software (compiler) prefetching

## 9. Reducing Misses by Hardware Prefetching of Instructions & Data

- Prefetching relies on having extra memory bandwidth that can be used without penalty
- **Instruction** Prefetching
  - Typically, CPU fetches **2** blocks on a miss: the requested block and the next consecutive block.
  - Requested block is placed in instruction cache when it returns, and prefetched block is placed into instruction stream buffer

# Array Storage in Memory

- We may view items as arrays (2D)
- But these are stored in main memory as 1D items
- How to translate? Two possibilities:
  - Row-Major Mapping
    - Used in C/C++, Python
  - Column-Major Mapping
    - Used in Fortran, MATLAB

# Row-Major Mapping

- Example 3 x 4 array:

a	b	c	d
e	f	g	h
i	j	k	l

- Convert into 1D array **y** by collecting elements by *rows*
- Within a row, elements are collected from left to right
- Rows are collected from top to bottom
- We get {a, b, c, d, e, f, g, h, i, j, k, l}

# Column-Major Mapping

a	b	c	d
e	f	g	h
i	j	k	l

- Convert into 1D array **y** by collecting elements by *columns*
- Within a column, elements are collected from top to bottom
- Columns are collected from left to right
- We get {a, e, i, b, f, j, c, g, k, d, h, l}



## Example 2.3-2

Assume that we have the following matrix stored in main memory in **column major order**:

$$\begin{bmatrix} A_{11} & A_{12} & A_{13} & A_{14} \\ A_{21} & A_{22} & A_{23} & A_{24} \\ A_{31} & A_{32} & A_{33} & A_{34} \\ A_{41} & A_{42} & A_{43} & A_{44} \end{bmatrix}$$

How are the matrix elements arranged in main memory?

## Example 2.3-3

Next, assume that our 4x4 matrix contains **double-precision** numbers.

Assuming the IEEE floating-point standard (as used by MIPS), how many bytes per number?

## Example 2.3-4

Next, assume that we have a **fully-associative** cache that views our main memory as **64-byte** blocks.

How are the matrix elements stored in main memory blocks?

## Example 2.3-5

Next, assume that we now have a 256 x 256 matrix. How are the matrix elements arranged in the memory blocks?

(Assume double-precision numbers and 64-byte blocks as before)

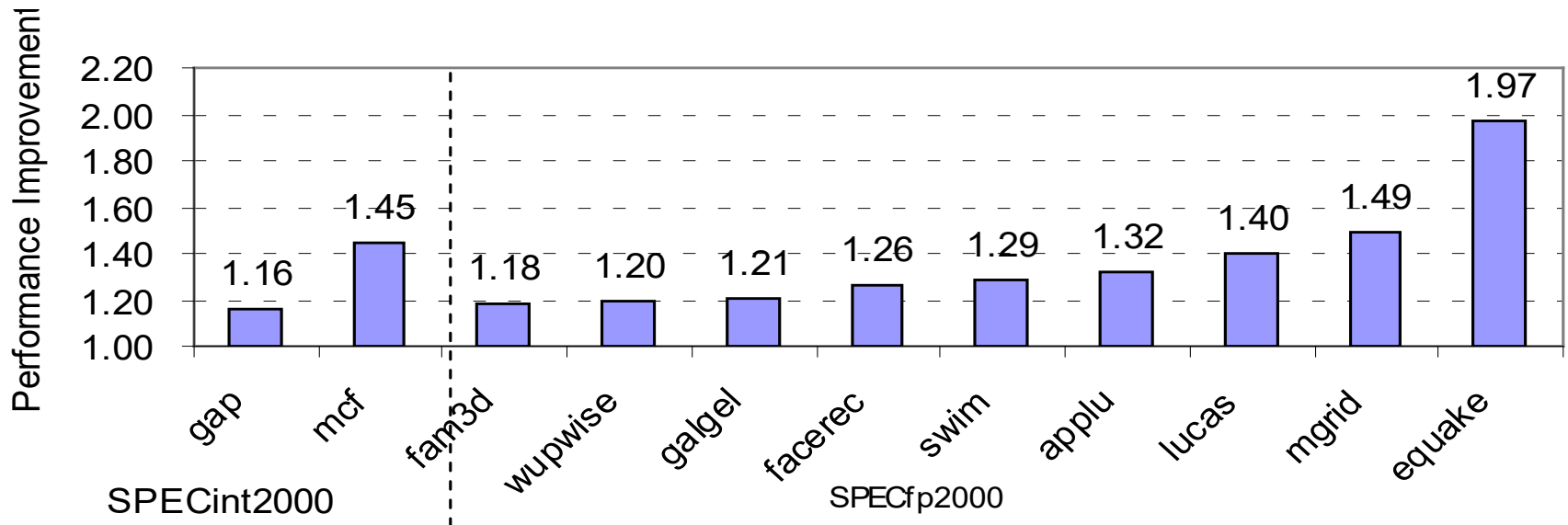
## Example 2.3-6

Given a matrix element  $A(\text{row}, \text{column})$  for our  $256 \times 256$  matrix, how do we identify the block in main memory where it is stored?

## 9. Reducing Misses by Hardware Prefetching of Instructions & Data (cont)

- Data Prefetching
  - Pentium 4 can prefetch data into L2 cache from up to 8 streams from 8 different 4 KB pages
  - Prefetching invoked if 2 successive L2 cache misses to a page, if distance between those cache blocks is  $< 256$  bytes

## 9. Reducing Misses by Hardware Prefetching of Instructions & Data (cont)



# 10. Reducing Misses by Software Prefetching Data

- Data Prefetch
  - Load data into register (HP PA-RISC loads)
  - Cache Prefetch: load into cache (MIPS IV, PowerPC, SPARC v. 9)
- Issuing Prefetch Instructions takes time
  - Is cost of prefetch issues < savings in reduced misses?
  - Higher superscalar reduces difficulty of issue bandwidth

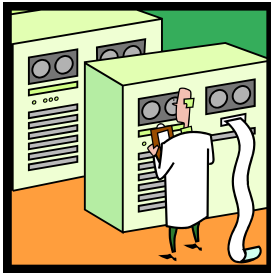


# Outline

- 2.1 Introduction
- 2.2 Memory Technology and Optimizations
- 2.3 Ten Advanced Optimizations of Cache Performance
- 2.4 Virtual Memory and Machines

# Introduction to Virtual Machines

- VMs developed in late 1960s
  - Remained important in mainframe computing over the years
  - Largely ignored in single user computers of 1980s and 1990s
  - But used in *networked* computers



# Introduction to Virtual Machines (cont)

- Recently regained popularity due to
  - increasing importance of isolation and security in modern systems,
  - failures in security and reliability of standard operating systems,
  - sharing of a single computer among many unrelated users,
  - and the dramatic increases in raw speed of processors, which makes the overhead of VMs more acceptable

# What is a Virtual Machine (VM)?

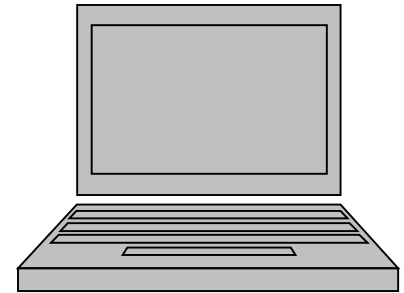
- Broadest definition includes all emulation methods that provide a standard software interface, such as the Java VM
- “(Operating) **System Virtual Machines**” provide a complete system level environment at binary ISA
  - Here assume ISAs always match the native hardware ISA

# What is a Virtual Machine (VM)?

- Present illusion that VM users have entire computer to themselves, including a copy of OS
- Single computer runs multiple VMs, and can support a multiple, different OSes
  - On conventional platform, single OS “owns” all HW resources
  - With a VM, multiple OSes all share HW resources
- Underlying HW platform is called the **host**, and its resources are shared among the **guest** VMs

# VM Examples?

- We have used virtual machines recently in this course
  - X-windows emulators
    - Xming (Windows)
    - Mac OS X
  - Your PC acts like a Unix machine
    - It is the “guest” VM
    - Lyle Linux Server is the “host” VM



# Data-Level Parallelism

(Chapter 4, Hennessy and Patterson)

Note: some course slides adopted  
from publisher-provided material

# Outline

- 4.1 Introduction
- 4.2 Vector Architecture
- 4.4 Graphics Processing Units (GPU's)



# Introduction

- SIMD = Single Instruction stream, Multiple Data streams
  - Single set of instructions accesses multiple sets of data
- MIMD = Multiple Instruction streams, Multiple Data streams

# Introduction

- SIMD architectures can exploit significant data-level parallelism for:
  - matrix-oriented scientific computing
  - media-oriented image and sound processors
- SIMD is more energy efficient than MIMD
  - Only needs to fetch one instruction per data operation
  - Makes SIMD attractive for personal mobile devices
- SIMD allows programmer to continue to think sequentially

# SIMD Parallelism

- Vector architectures
- SIMD extensions
- Graphics Processor Units (GPUs)

# Outline

- 4.1 Introduction
- 4.2 Vector Architecture
- 4.4 Graphics Processing Units (GPU's)