# CPE 431/531

# Chapter 6 – Parallel Processors from Client to Cloud

# Dr. Rhonda Kay Gaede



THE UNIVERSITY OF
ALABAMA IN HUNTSVILLE

# 6.1 Motivation

- ## Why multiprocessors?

    - _____

    - _____

    - _____

# 6.2 The Difficulty of Creating Parallel Programs

- The difficulty with parallelism is not the _____, it's the _____.

- Why is it difficult to write parallel processing programs that are fast?

  - _____

  - _____ _____

  - _____ _____

  - _____ _____

  - _____ _____

# 6.2 Speedup Challenge

- Suppose you want to achieve a speedup of 90 times faster with 100 processors. What percentage of the original computation can be sequential?

# 6.2 Speedup Challenge – Bigger Problem

- Suppose you want to perform two sums: one is a sum of 10 scalar variables and one is a matrix sum of a pair of two-dimensional arrays, size 10 by 10. What speedup do you get with 10 versus 40 processors?

# 6.2 Speedup Challenge – Bigger Problem

- Next, calculate the speed-ups assuming the matrices grow to 20 by 20

- Strong scaling

- Weak scaling

# 6.2 Speedup Challenge: Balancing Load

- To achieve the speed-up of 20.5 on the previous larger problem with 40 processors, we assumed the load was perfectly balanced (each processor did 2.5 % of the work). Instead, show the impact on speed-up if one processor's load is higher than all the rest. Calculate at 5% and 12.5%.

# 6.3 SISD, MIMD, SIMD, SPMD, and Vector

- SISD is the normal case – single instruction, single data.

- MIMD – multiple instruction, multiple data - is _____ possible but programmers normally write a _____ _____ that runs on all processors relying on _____ statements when _____ processors should execute _____ sections of code. This style is single _____, multiple data.

- SIMD – single instruction, multiple data – operate on _____ of data. SIMD needs only ____ _____ of the code that is being simultaneously executed. SIMD works best when dealing with _____ in ____ loops. _____ _____ _____

- The _____ _____ that inspired the SIMD category faded into history but two _____ interpretations of SIMD remain _____ today.
  - _____
  - _____

# 6.3 x86 Multimedia Extensions

- The most _____ used variation of SIMD is the basis of the hundreds of MMX and SSE instructions of the x86 processor. These instructions were added to improve performance of _____ programs.

  - The hardware allows flexible ALU operations – one 64-bit or two 32-bit or for 16-bit or eight 8-bit

  - Loads and stores are simply as wide as the widest ALU.

  - SSE now supports simultaneous execution of a pair of 64-bit floating-point numbers

# 6.3 Vector

- An older and more elegant interpretation of SIMD is called a vector architecture, which has been closely identified with Cray Computers.

- Consider $Y = a \times X + Y$

Original

```
        l.d    $f0,a($sp)
        addiu $t1,$s0,#512
loop:  l.d    $f2,0($s0)
        mul.d $f2,$f2,$f0
        l.d    $f4,0($s1)
        add.d $f4,$f4,$f2
        s.d    $f4,0($s1)
        addiu $s0,$s1,#8
        addiu $s1,$s1,#8
        subu  $t0,$t1,$s0
        bne    $t0,$zero,loop
```
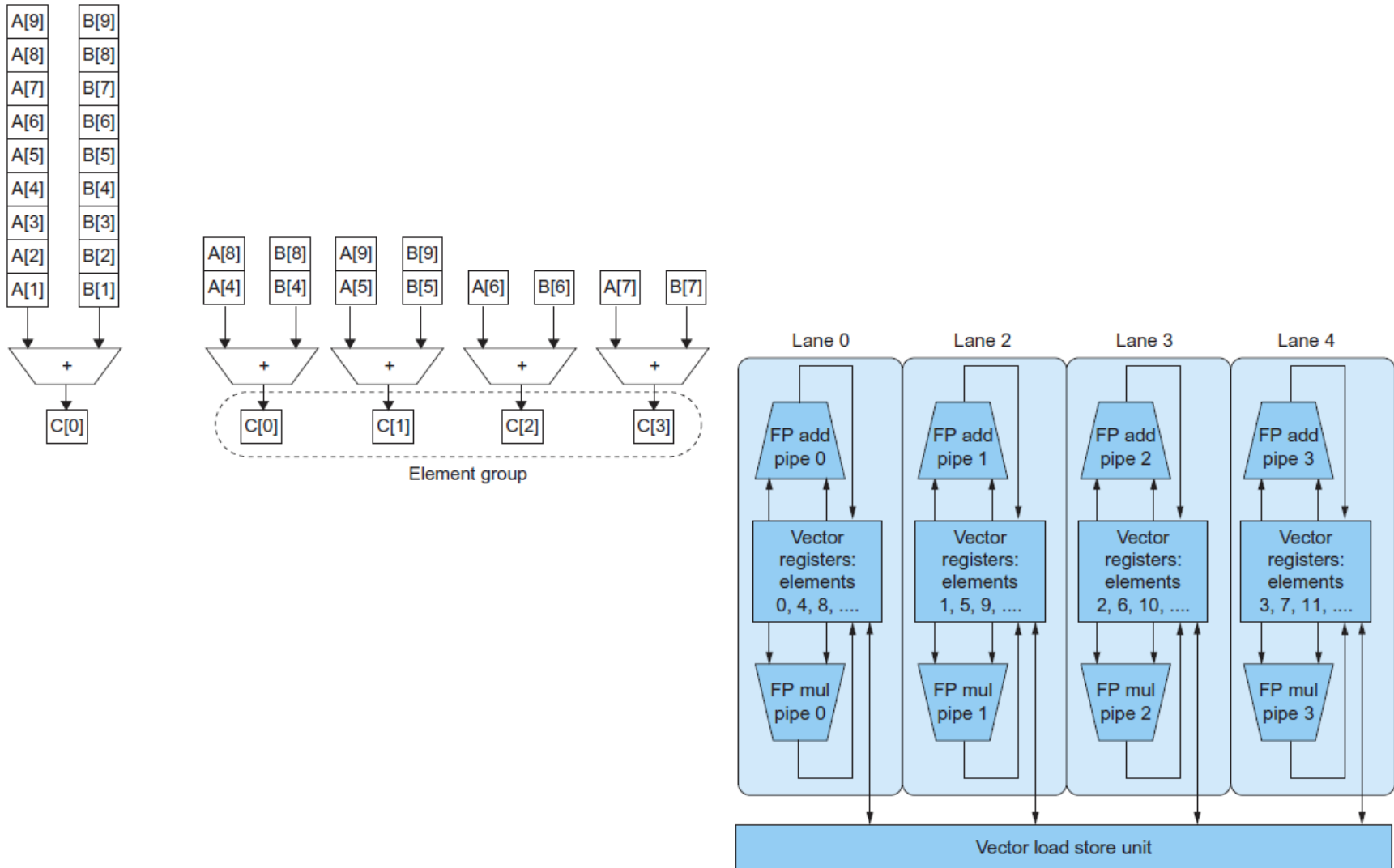
Vector

```
l.d      $f0, a($sp)
lv       $v1,0($s0)
mulvs.d $v2,$v1,$f0
lv       $v3,0($s1)
addv.d  $v4,$v4,$v3
sv       $v4,0($s1)
```

# 6.3 Comparisons

- Vector versus Scalar
  - A _____ _____ instruction specifies a great deal of work, the instruction _____ and _____ bandwidth is greatly reduced.
  - Hardware does not have to check for _____ _____ _____ a vector instruction.
  - Vector architectures and compilers have worked well for _____

    _____.
  - Hardware need only check for _____ hazards _____ between two vector _____, not _____ for every _____ _____.
  - Vector instructions that access memory have a _____ _____ _____, memory system can be adjusted accordingly.
  - Replacing a _____ with a _____ _____ reduces _____ hazards.
- Vector versus Multimedia Extensions
  - Vector specifies the number of operands in _____, not in _____.
  - Vector data transfers need ____ be _____.
  - Vector _____ over time more _____.
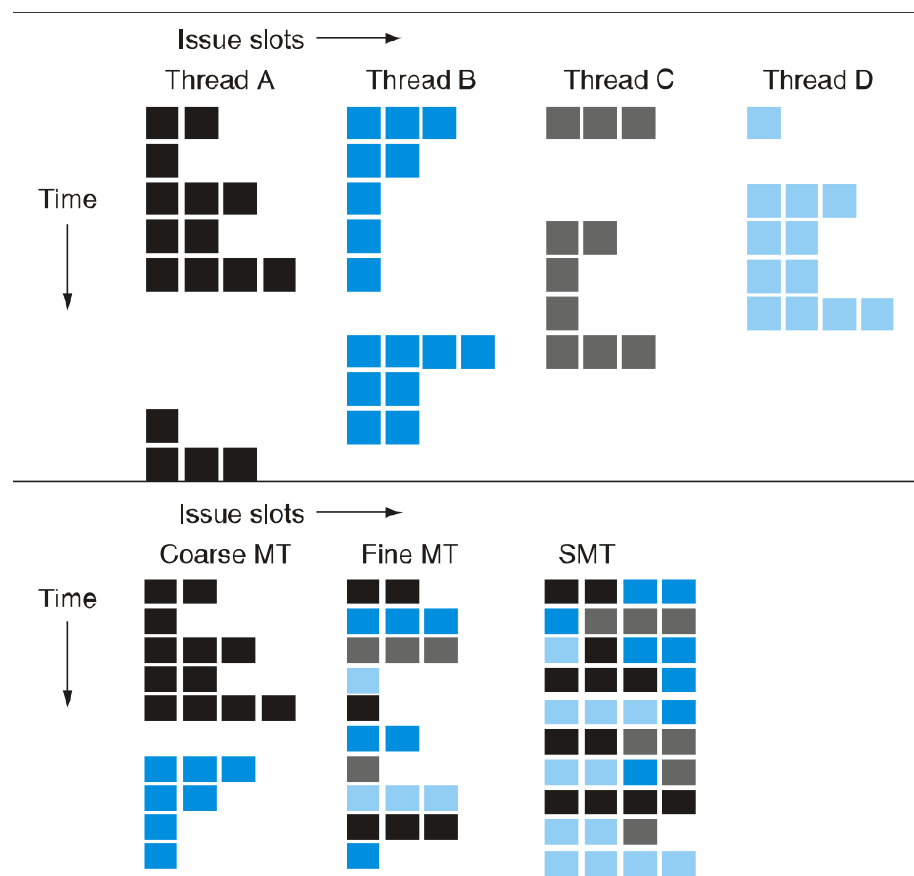
# 6.3 Improving the Performance of Vector

# 6.4 Hardware Multithreading

- Hardware multithreading allows _____ _____ to _____ the _____ units of a single processor in _____ fashion.

- The processor must _____ the independent _____ of each _____.

- The hardware must support the _____ to _____ to different _____ relatively quickly.

- _____ multithreading switches between _____ on each _____, often done round robin.

  - Hides _____ losses by doing useful work during _____ _____.

  - Inserts _____ for threads with __ _____.

- _____ multithreading _____ threads only on _____ stalls, such as _____ _____ misses.

  - It is limited in its ability to overcome _____ losses, especially from _____ stalls. The major problem is pipeline ___ _____.

# 6.4 Simultaneous Multithreading

- Simultaneous multithreading uses the resources of a _____, _____ scheduled processor to exploit _____ parallelism at the same time it exploits _____ parallelism.

- The key insight is that multiple-issue processors often have more _____ _____ _____ than a _____ thread can effectively use.

- _____ of dependences can be handled by the _____ _____ capability.
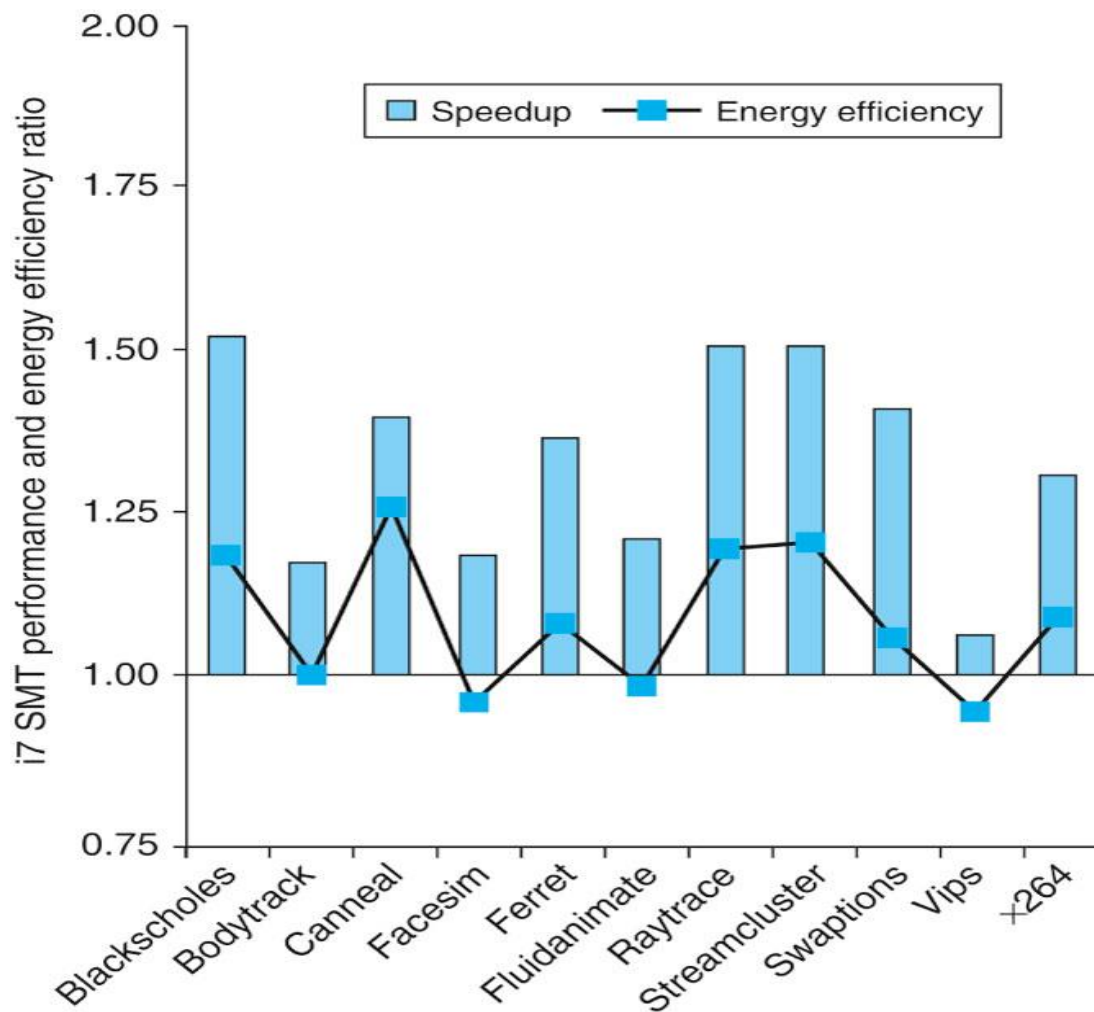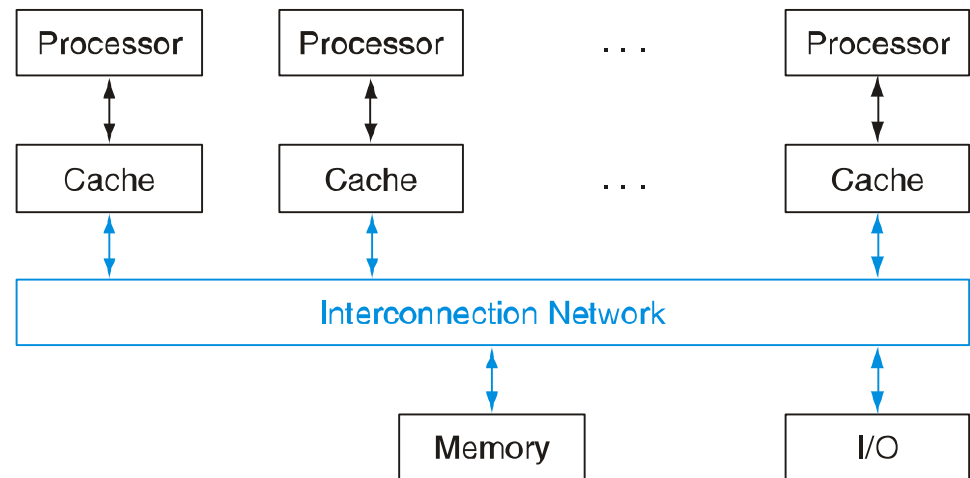
# 6.4 Multithreading Speedup



FIGURE 6.6   The speed-up from using multithreading on one core on an i7 processor averages 1.31 for the PARSEC benchmarks (see  Section 6.9) and the energy efficiency improvement is 1.07. This data was collected and analyzed by Esmaeilzadeh et. al. [2011].

# 6.5 Multicore and Other Shared Memory Multiprocessors

- A shared memory multiprocessor (SMP) is one that offers the programmer a _____ _____ _____ _____ across all processors
- Processor communicate through _____ _____ in memory.
- SMPs come in two flavors
    - _____
    - _____

- Processors need to coordinate when sharing data, this process is called _____, processors must acquire a _____
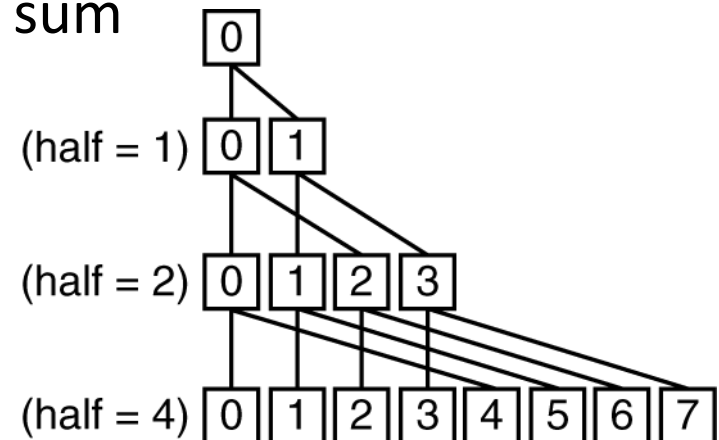
# 6.5 Shared Address Space Parallel Program (1)

- Suppose we want to sum 64,000 numbers on an SMP with UMA. Let's assume we have 64 processors.

```
sum[Pn] = 0;
for (i = 1000*Pn; i < 1000*(Pn+1); i = i + 1)
  sum[Pn] = sum[Pn] + A[i]; /* sum the assigned areas */
```

- After execution of this code, there are 64 _____ sums
- Need to _____ them into _____ sum
- Do so using a _____

# 6.5 Shared Address Space Parallel Program (2)

```
half = 64; /* 64 processors in multiprocessor */
repeat
  synch(); /* wait for partial sum completion */
  if (half%2 != 0 && Pn == 0)
    sum[0] = sum[0] + sum[half-1];
    /* Conditional sum needed when half is odd;
        Processor0 gets missing element */
  half = half/2; /* dividing line on who sums */
  if (Pn < half) sum[Pn] = sum[Pn] + sum[Pn+ half];
until (half == 1); /* exit with final sum in sum[0] */
```

# 6.5 A Parallel Programming System

- A _____ but _____ example is _____
  - _____ is an _____ _____ _____ along with a set of _____ _____, _____ variables, and _____ library routines.
  - It offers a _____, _____, and _____ programming model for shared memory multiprocessors.
  - Its primary goal is to _____ _____ and to perform _____.
  - _____ extends __using _____, commands to the C macro processor

# 6.5 OpenMP Example

```
Cc -fopenmp foo.c
#define  P 64 /* define a constant */
#pragma omp parallel num_threads(P)

#pragma omp parallel for
for (Pn = 0; Pn < P; Pn +=1)
  for (1000*Pn; i < 1000*(Pn +1); i +=1)
    sum[Pn] += A[i]; /* sum the assigned areas */

#pragma omp parallel for reduction(+ : FinalSum)
for (i = 0; i < P; i += 1)
  FinalSum += sum[i]; /* Reduce to asingle number */
```

# 6.6 Introduction to Graphics Processing Units

- A major driving force for improving graphics processing was the _____ _____, a different _____ _____ than the one for CPUs.

- Key differences between GPUs and CPUs

  - GPUs are _____ that _____ a CPU, they don't have to do _____.

  - GPU problems sizes are typically hundreds of _____ to _____, but not hundreds of _____ to _____.

- Different Architecture Features

  - GPUs do not rely on _____ caches, they rely on having enough _____ to _____ memory latency.

  - The GPU main memory is oriented towards _____ rather than _____.

  - Each GPU processor is more highly _____ than a typical CPU, plus they have _____ _____.

# 6.6 Programming GPUs

- Initially programmers had only _____ _____ and _____.
  - They developed _____ _____languages
  - NVIDIA _____ _____ _____ _____ (CUDA)
  - OpenCL is a _____ initiative to develop a _____ programming language
  - Unifying theme is _____ _____
  - _____ and _____ can gang thousands of CUDA threads together to utilize _____, _____, _____, and _____
  - Threads are _____ together and executed in _____ of ____at a time
  - A _____ processor inside a GPU executes these blocks of threads, and a GPU consists of ___ to _____ of these _____ processors.