# 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?

  - ## performance

  - ## power efficiency

  - ## fault tolerance

# 6.2 The Difficulty of Creating Parallel Programs

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

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

  - ___scheduling___
  - ___load___ ___balancing___
  - ___synchronization___ ___time___
  - ___communication___ ___overhead___
  - ___Amdahl's___ ___law___

# 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?

$$\frac{ET\ affected}{Amount\ of\ improvement}+ET\ unaffected = ET\ after \qquad Speedup = \frac{ET\ before}{ET\ after}$$

$$90 = \frac{ET\ b}{(ET\ b - ET\ aff) + ET\ aff/100} \cdot \frac{\frac{1}{ETb}}{\frac{1}{ETb}} \qquad\qquad fract = \frac{ET\ aff}{ET\ b}$$

$$90 = \frac{1}{1 - fract + fract/100}$$

$$90 - 90\ fract + 0.9\ fract = 1 \qquad\qquad Sequential = 1 - fract = 0.001 = 0.1\%$$
$$89 = 89.1\ fract$$
$$Fract = 0.999$$

# 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?

Scalar 10t
Matrix 100 r
ET before = 110 t

10 processors
ET after = 10t + 110t/10 = 20t
Speedup = 110t/20t = 5.5
Efficiency = 5.5/10 = 55%

40 processors
ET aff = 10t + 100t/40 = 12.5t
Speedup = 110t/12.5t = 8.8
Efficiency = 8.8/40 = 22%

# 6.2 Speedup Challenge – Bigger Problem

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

  How many additions? 400t  Still doing scalar sum of 10

  10 processors

  Original Time = 410t

  Parallel Time = 10t + 400t/10 = 50t

  Speedup = 410t/50t = 8.2

  40 processors

  Parallel Time = 10t + 400t/40 = 20t

  Speedup = 410t/20t = 20.5

- Strong scaling – measure speedup with fixed problem

- Weak scaling – program size grows proportionally to the number of processors

# 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%.

5% load
P5% = 20t
39 processors share 380t
Max (20t, 380t/39) + 10t = 30t
Speedup = 410t/30t = 13.67

12.5% load
P12.5% = 50t
39 processors share 350t
Max (50t, 350t/39) + 10t = 60t
Speedup = 410t/60t = 6.83

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

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

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

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

- The _array_ _processors_ that inspired the SIMD category faded into history but two _current_ interpretations of SIMD remain _active_ today.
  - _SIMD in x86 Multimedia Extensions_
  - _Vector_

# 6.3 x86 Multimedia Extensions

- The most __widely__ 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 __multimedia__ programs.

    - The hardware allows flexible ALU operations – one 64-bit or two 32-bit or four 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
```
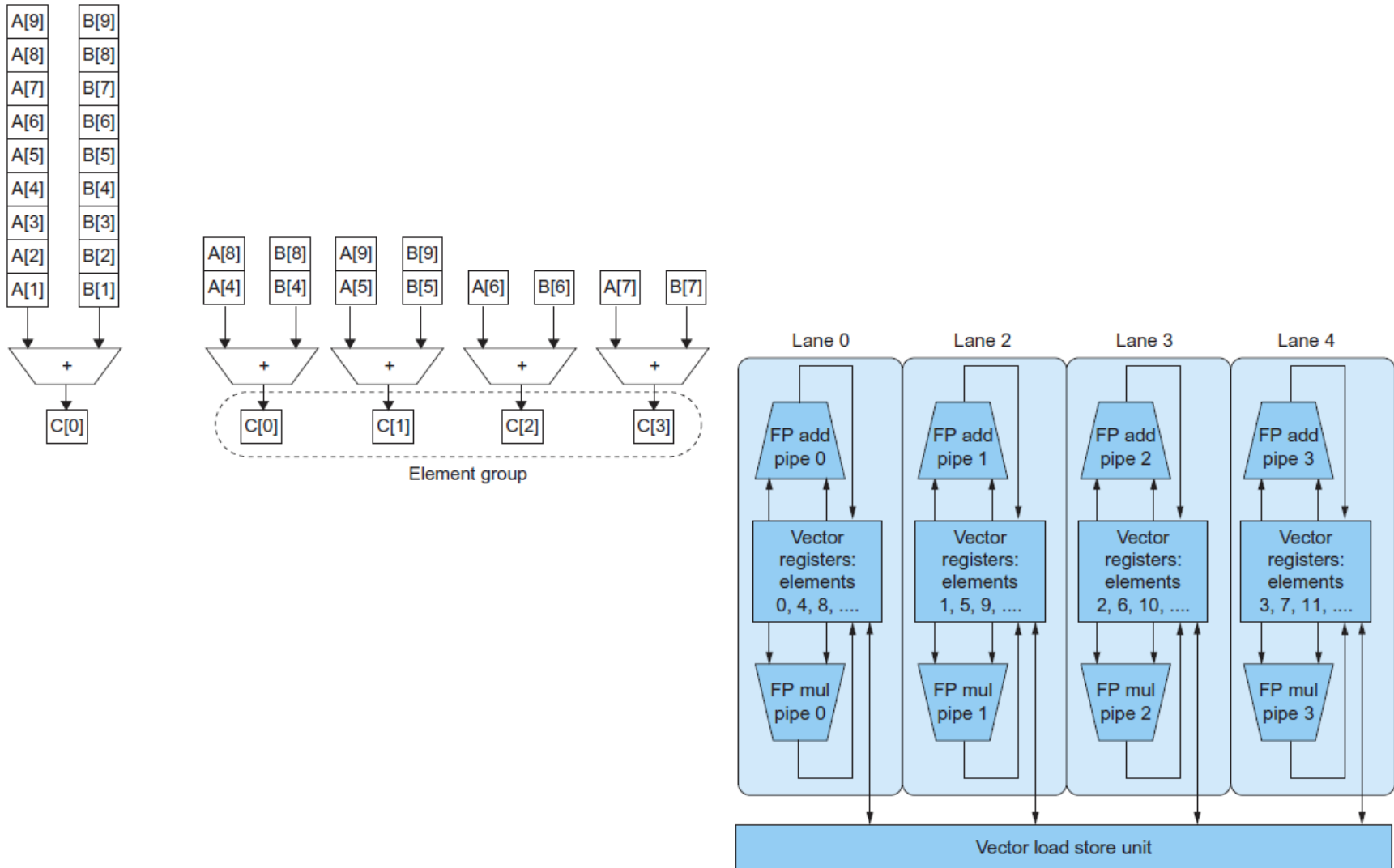
Vector
```
l.d      $f0, a($sp)
lv       $v1,0($s0)
mulvs.d  $v2,$v1,$f0
lv       $v3,0($s1)
addv.d   $v4,$v4,$v3
sv       $v4,0($s1)
```

# 6.3 Comparisons

- Vector versus Scalar
  - A _single_ _vector_ instruction specifies a great deal of work, the instruction _fetch_ and _decode_ bandwidth is greatly reduced.
  - Hardware does not have to check for _data_ _hazards_ _within_ a vector instruction.
  - Vector architectures and compilers have worked well for _data-level_ _parallelism_.
  - Hardware need only check for _data_ hazards _once_ between two vector _instructions_, not _once_ for every _vector_ _element_.
  - Vector instructions that access memory have a _known_ _access_ _pattern_, memory system can be adjusted accordingly.
  - Replacing a _loop_ with a _vector_ _instruction_ reduces _control_ hazards.

- Vector versus Multimedia Extensions
  - Vector specifies the number of operands in _registers_, not in _opcodes_.
  - Vector data transfers need _not_ be _contiguous_.
  - Vector _evolves_ over time more _gracefully_.

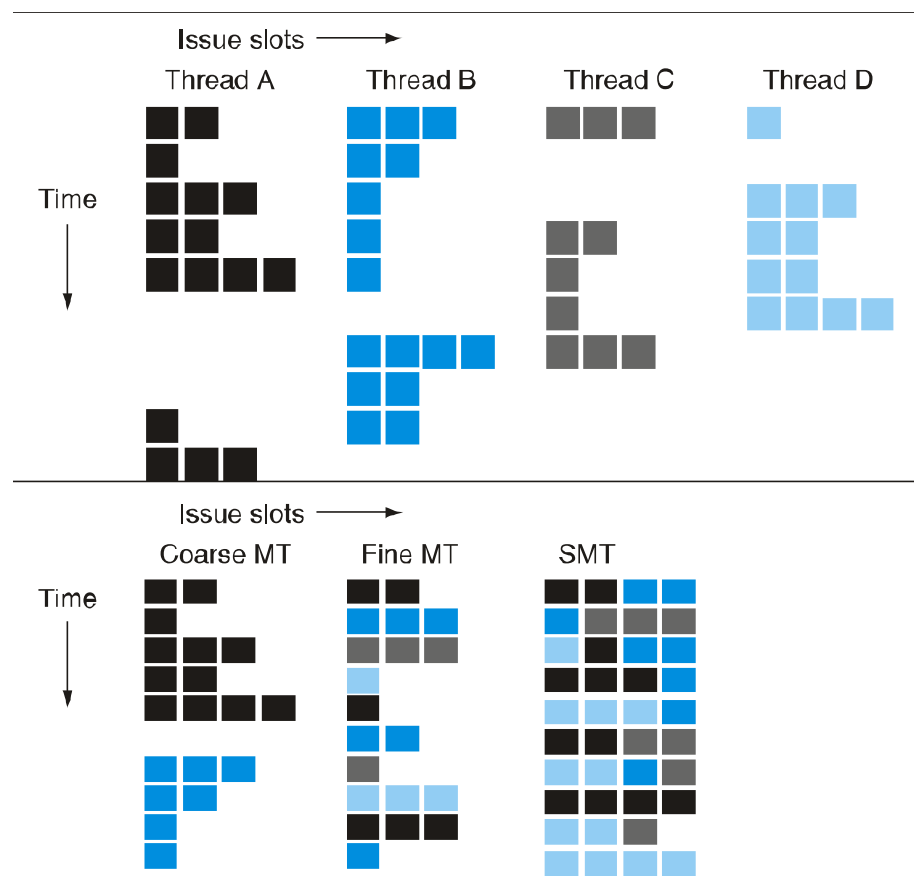# 6.3 Improving the Performance of Vector

# 6.4 Hardware Multithreading

- Hardware multithreading allows __multiple__ __threads__ to __share__ the __functional__ units of a single processor in __overlapping__ fashion.

- The processor must __duplicate__ the independent __state__ of each __thread__.

- The hardware must support the __ability__ to __change__ to different __thread__ relatively quickly.

- __Fine-grained__ multithreading switches between __threads__ on each __instruction__, often done round robin.
    - Hides __throughput__ losses by doing useful work during __thread__ __stalls__.
    - Inserts __latency__ for threads with __no__ __stalls__.

- __Coarse-grained__ multithreading __switches__ threads only on __costly__ stalls, such as __second-level__ __cache__ misses.
    - It is limited in its ability to overcome __throughput__ losses, especially from __shorter__ stalls. The major problem is pipeline __fill__ __time__.

# 6.4 Simultaneous Multithreading

- Simultaneous multithreading uses the resources of a __multiple-issue__, __dynamically__ scheduled processor to exploit __thread-level__ parallelism at the same time it exploits __instruction-level__ parallelism.

- The key insight is that multiple-issue processors often have more __functional__ __unit__ __parallelism__ than a __single__ thread can effectively use.

- __Resolution__ of dependences can be handled by the __dynamic__ __scheduling__ capability.



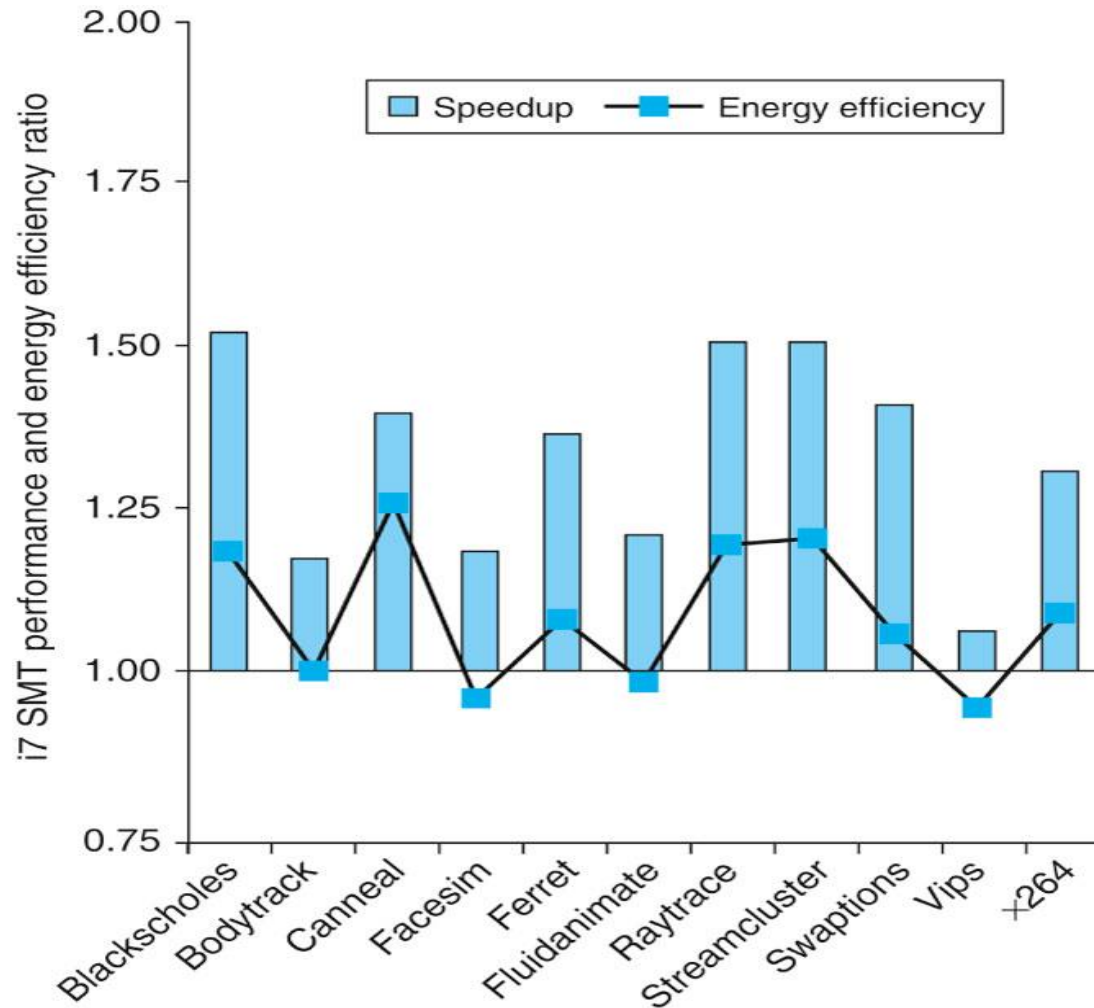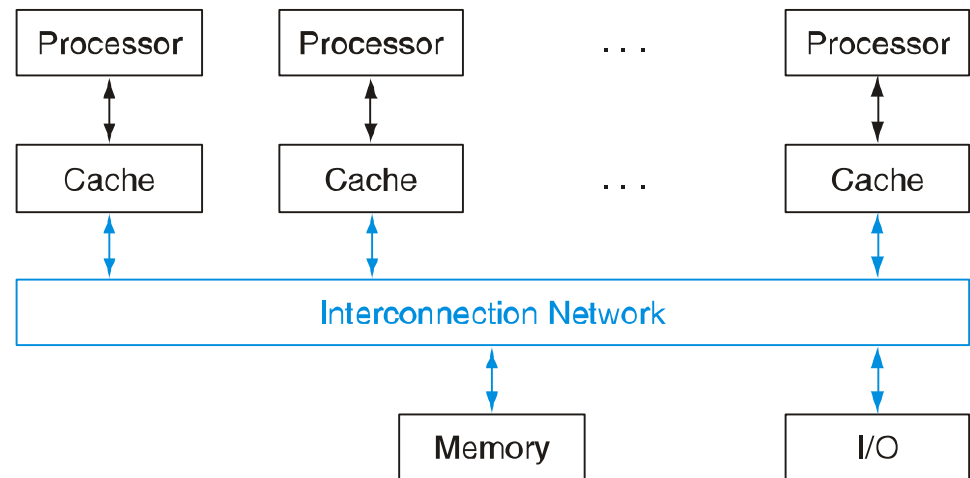Utilization 19/40          20/40          39/40

# 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].

THE UNIVERSITY OF
ALABAMA IN HUNTSVILLE

# 6.5 Multicore and Other Shared Memory Multiprocessors

- A shared memory multiprocessor (SMP) is one that offers the programmer a _single_ _physical_ _address_ _space_ across all processors
- Processor communicate through _shared_ _variables_ in memory.
- SMPs come in two flavors
  - _UMA_
  - _NUMA_

- Processors need to coordinate when sharing data, this process is called _synchronization_, processors must acquire a _lock_

```
┌───────────┐   ┌───────────┐           ┌───────────┐
│ Processor │   │ Processor │   . . .    │ Processor │
└───────────┘   └───────────┘           └───────────┘
      ↕               ↕                        ↕
┌───────────┐   ┌───────────┐           ┌───────────┐
│   Cache   │   │   Cache   │   . . .    │   Cache   │
└───────────┘   └───────────┘           └───────────┘
      ↕               ↕                        ↕
┌──────────────────────────────────────────────────┐
│              Interconnection Network               │
└──────────────────────────────────────────────────┘
              ↕                        ↕
        ┌──────────┐            ┌──────────┐
        │  Memory  │            │   I/O    │
        └──────────┘            └──────────┘
```
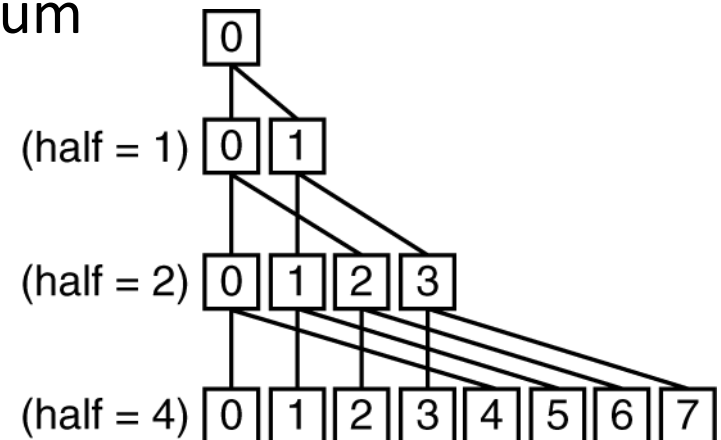
# 6.5 Shared Address Space Parallel Program (1)

- 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 partial sums
- Need to combine them into single sum
- Do so using a reduction

# 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 limited but popular example is OpenMP

  – OpenMP is an Application Programmer Interface along with a set of compiler directives, environment variables, and runtime library routines.

  – It offers a portable, scalable,and simple programming model for shared memory multiprocessors.

  – Its primary goal is to parallelize loops and to perform reductions.

  – OpenMP extends C using pragmas, 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 a single number */
```

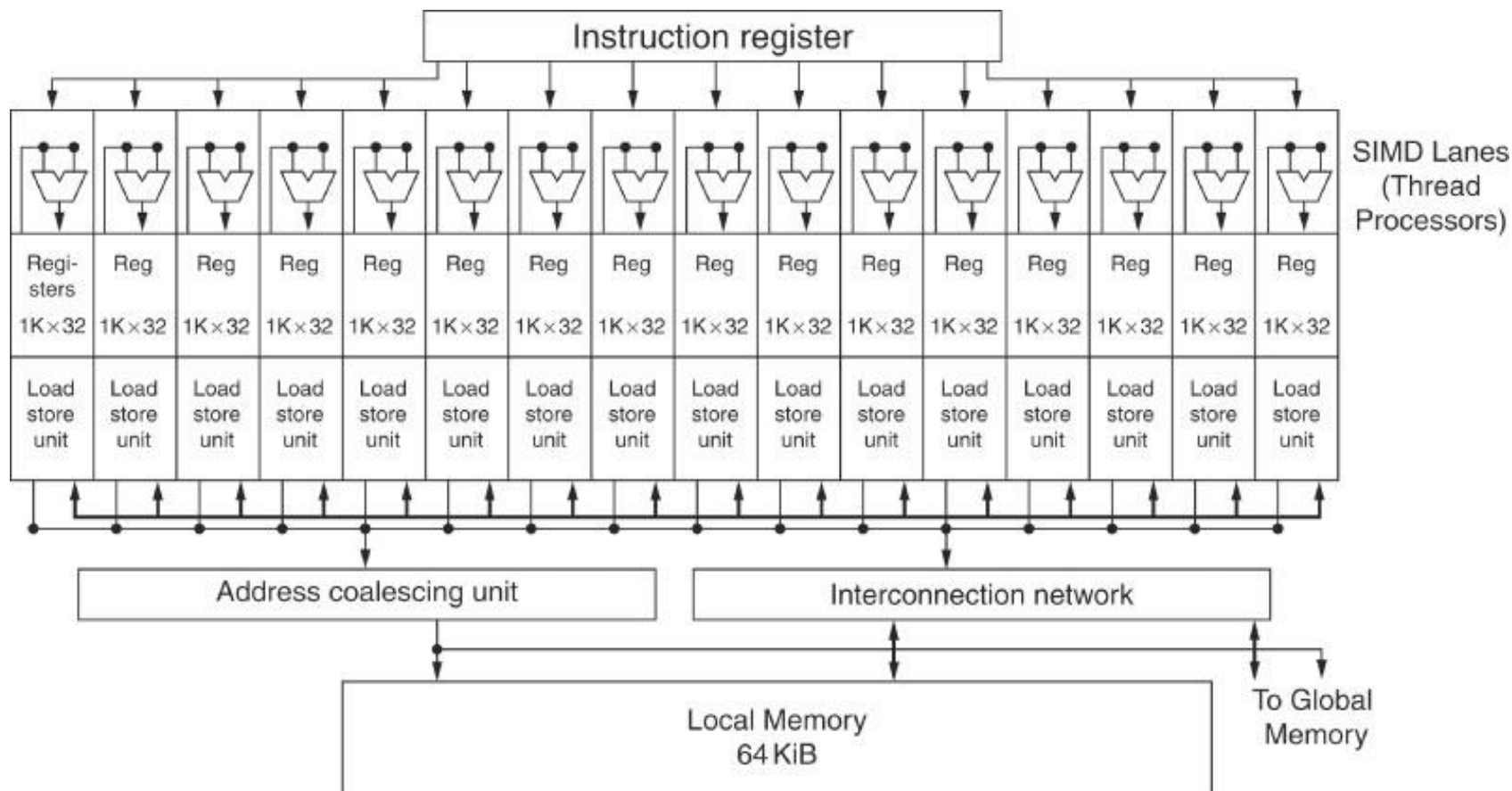# 6.6 Introduction to Graphics Processing Units

- A major driving force for improving graphics processing was the gaming industry, a different development community than the one for CPUs.

- Key differences between GPUs and CPUs

  - GPUs are accelerators that supplement a CPU, they don't have to do everything.

  - GPU problems sizes are typically hundreds of megabytes to gigabytes, but not hundreds of gigabytes to terabytes.

- Different Architecture Features

  - GPUs do not rely on multilevel caches, they rely on having enough threads to hide memory latency.

  - The GPU main memory is oriented towards bandwidth rather than latency.

  - Each GPU processor is more highly multithreaded than a typical CPU, plus they have more processors.
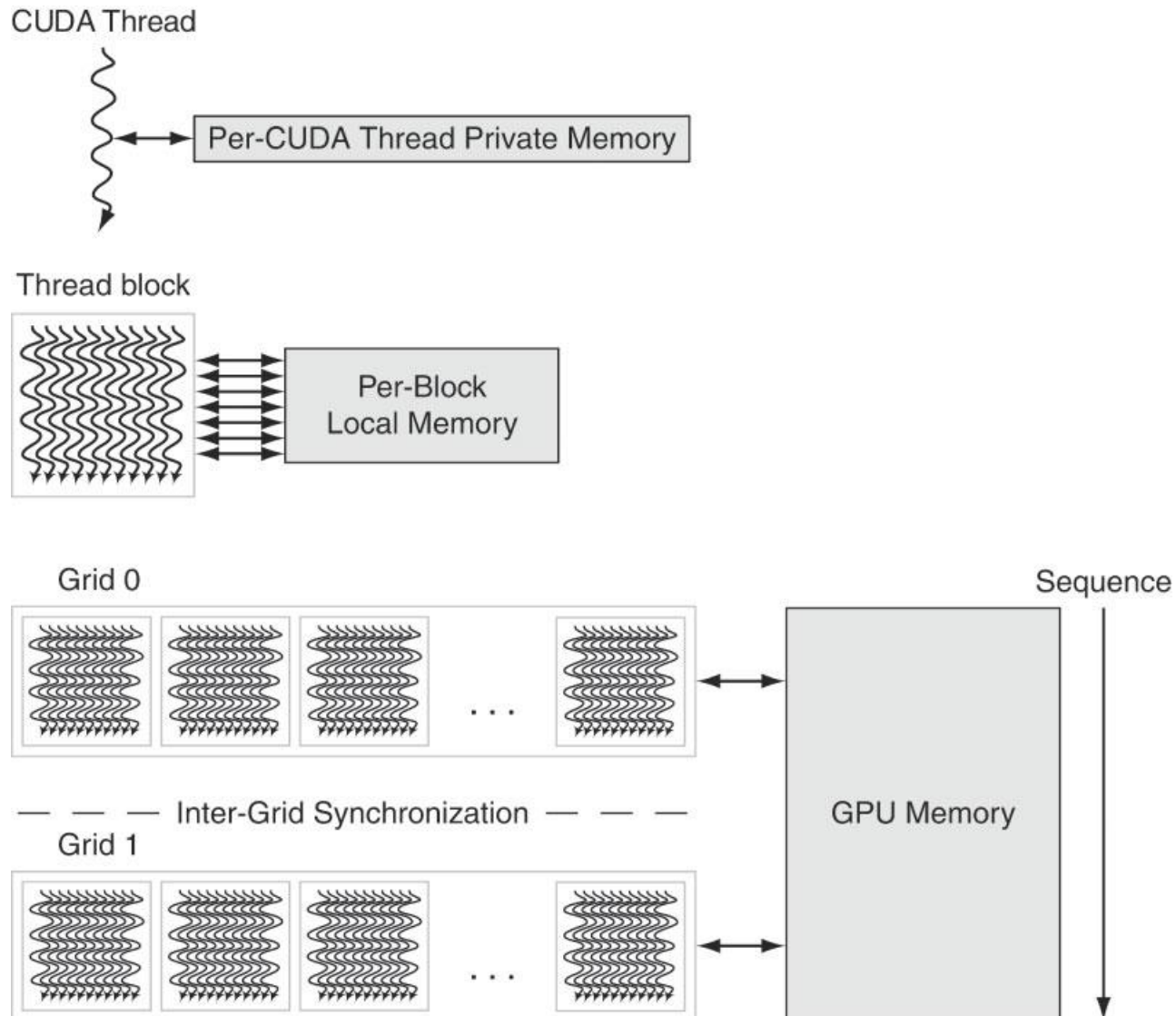
# 6.6 Programming GPUs

- Initially programmers only had graphics APIs and languages.

    - They developed C-inspired programming languages

    - NVIDIA Compute Unified Device Architecture (CUDA)

    - OpenCL is a multi-company initiative to develop a portable programming language

    - Unifying theme is CUDA thread

    - Compiler and hardware can gang thousands of CUDA threads together to utilize multithreading, MIMD, SIMD, and instruction-level parallelism

    - Threads are blocked together and executed in groups of 32 at a time

    - A multithreaded processor inside a GPU executes these blocks of threads, and a GPU consists of 8 to 32 of these multithreaded processors.

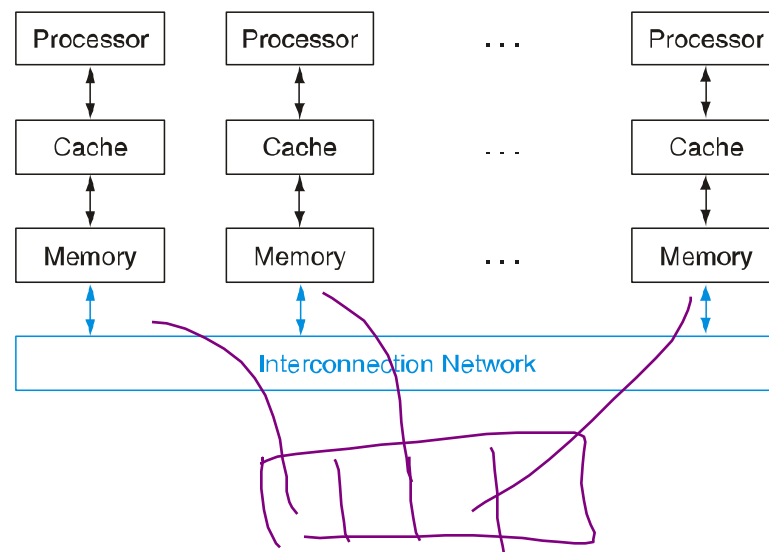# 6.6 Block Diagram of SIMD Processor

# 6.6 GPU Memory Structures



CUDA Thread

Per-CUDA Thread Private Memory

Thread block

Per-Block Local Memory

Grid 0

Sequence

Inter-Grid Synchronization

Grid 1

GPU Memory

# 6.6 Putting GPUs into Perspective

| Feature | Multicore with SIMD | GPU |
|---|---|---|
| SIMD processors | 4 to 8 | 8 to 16 |
| SIMD lanes/processor | 2 to 4 | 8 to 16 |
| Multithreading hardware support for SIMD threads | 2 to 4 | 16 to 32 |
| Typical ratio of single precision to double-precision performance | 2:1 | 2:1 |
| Largest cache size | 8 MB | 0.75 MB |
| Size of memory address | 64-bit | 64-bit |
| Size of main memory | 8 GB to 256 GB | 4 GB to 6 GB |
| Memory protection at level of page | Yes | Yes |
| Demand paging | Yes | No |
| Integrated scalar processor/SIMD processor | Yes | No |
| Cache coherent | Yes | No |

# 6.7 Message Passing Multiprocessors

- When an address space is not __shared__, communication occurs via explicit __message__ __passing__.
- __Communication__ occurs using __send__ and __receive__ messages.
- __Task-level__ parallelism and applications with little __communication__ do not require __shared__ addressing to run well. Examples include __web__ __search__, __mail__ __servers__, and __file__ __servers__.

- __Special__ __purpose__ inter-connection networks have been used, they provide __high__ performance at __much higher__ cost than LANs.

- More __common__ are __clusters__ that are collections of computers connected via __standard__ networks.

| Processor | Processor | ... | Processor |
| --- | --- | --- | --- |
| Cache | Cache | ... | Cache |
| Memory | Memory | ... | Memory |

Interconnection Network

# 6.7 Message Passing Program

```
sum = 0;
for (i = 0; i<1000; i = i + 1)
    sum = sum + AN[i];
limit = 100; half = 100;/* 100 processors */
repeat
    half = (half+1)/2; /* send vs. receive
                            dividing line */
    if (Pn >= half && Pn < limit)
      send(Pn - half, sum);
    if (Pn < (limit/2))
      sum = sum + receive();
    limit = half; /* upper limit of senders */
  until (half == 1); /* exit with final sum */
```

# 6.7 Hardware/Software Interface

- Message passing systems are easier for __hardware__ designers to __build__ .

- For programmers, there are fewer __side__ __effects__ , the communication is explicit, there is no guessing about the __cache__ __coherence__ performance

- However, it's harder to __port__ a __sequential__ program to a message-passing computer.

- Modern systems are a __hybrid__ ; __multicore__ multiprocessors use __shared__ physical memory and __nodes__ of a __cluster__ communicate with each other using __message__ __passing__ .

- The weakness of __separate__ memories for user memory from a parallel programming perspective turns into a strength in system __dependability__ .

- Computers can be __replaced__ in a cluster without __bringing__ the system __down__ .

- Work can also be more easily __reallocated__ from __failing__ servers.

- Systems can be more easily __expanded__ using __clusters__ .

# 6.7 Warehouse Scale Computers (WSCs)

- The most popular framework for batch processing in a WSC is MapReduce and the open source version, Hadoop, inspired by _Lisp_ functions of the same name.
- WSCs require innovations in _power_ _distribution_, _cooling_, _monitoring_, and _operations_, they are a modern descendant of the 1970s supercomputer.
- The 1970s supercomputer provided the _few_ _companies_ that could afford it, high performance computing for _scientists_ and _engineers_
- Warehouse Scale Computers make it possible for us to have _Internet_ _sensations_ on _YouTube_.

# 6.7 Warehouse Scale Computers (WSCs)

- WSCS have three major distinctions
  - _Ample_ , _easy_ parallelism - or _request-level_ parallelism
  - _Operations_ Costs Count – not just _purchase_ price, energy, power distribution, and cooling represent more than _30%_ of the costs of a WSC over _10_ _years_
  - _Scale_ and the Opportunities/Problems Associated with _scale_
    - 100,000 servers mean _volume_ _discounts_
    - _Time_ on servers can be _sold_
    - 100,000 servers mean lots of _failures_ - _disks_ and _servers_
    - _Fault_ _tolerance_ takes on more importance

# 6.7 Warehouse Scale Computers (WSCs)

- In __2012__, Amazon Web Services announced the amount of new server capacity it adds __every__ __day__ is sufficient to support all of Amazon in 2003 when it was a company with __6000__ employees and __$5.2 billion__ in annual revenue.

- The growth of cloud computing could be slowed by __security__ concerns, __privacy__ concerns, __standards__ and the __rate__ of __growth__ of Internet bandwidth.

- https://www.youtube.com/watch?v=zRwPSFpLX8I

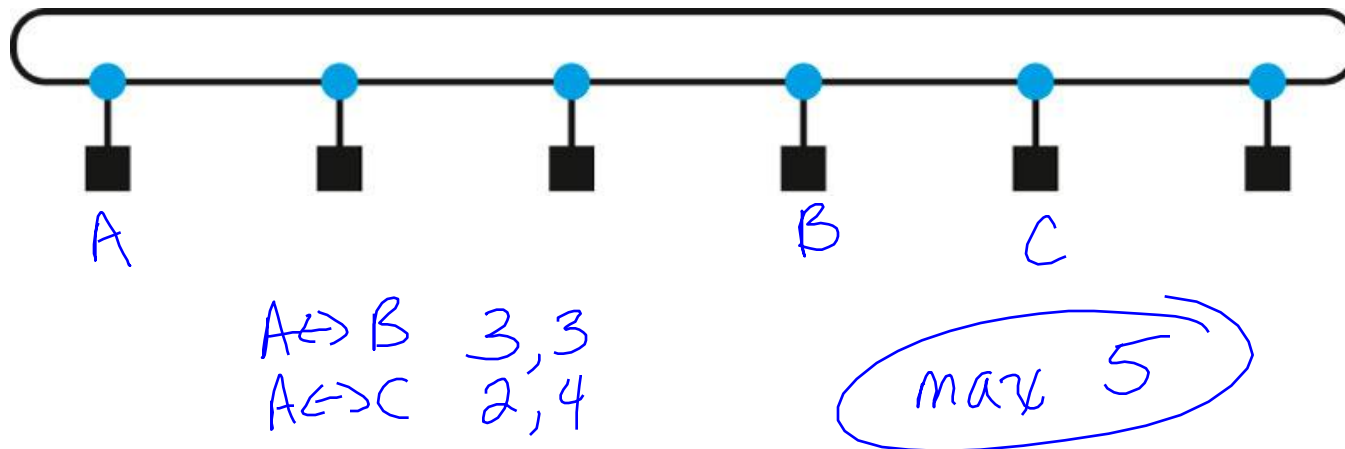- https://www.youtube.com/watch?v=XZmGGAbHqa0

# 6.8 Introduction to Multiprocessor Network Topologies

- The popularity of __cloud__ __computing__ leads to a need for __interconnection__ networks between __nodes__ in a warehouse scale computer.

- The increasing number of __cores__ per chips means we need networks __inside__ a chip as well.

- Network costs include the __number__ of __switches__, the __number__ of __links__ on a switch to connect to the network, the __width__ per __link__, and __length__ of the link when the network is mapped onto __silicon__.

- Network performance includes
  - The __latency__ on an __unloaded__ network to send and receive a message
  - The __throughput__ in terms of the __maximum__ number of messages that can be __transmitted__ in a given time period
  - __Delays__ caused by __contention__ for a portion of the network
  - Variable performance depending on the __pattern__ of communication.
  - __Fault__ __tolerance__
  - __Energy__ __efficiency__

# 6.8 Introduction to Multiprocessor Network Topologies

- Networks are normally drawn as graphs
  - *Edge* represents *link*
  - *Node* represents *computers*
  - Links are *bidirectional*
  - Networks consist of *switches*

- First example is a *ring*
  - Capable of *many* *simultaneous* transfers

A          B  C

A<=>B  3,3
A<=>C  2,4

*max 5*

# 6.8 Network Performance Metrics

- Network Bandwidth - the bandwidth of __each__ __link__ multiplied by the __number__ of __links__ (__best__ __case__)
  - Ring - *P times the bandwidth of the link*    *P - number of processors*
  - Bus – *the bandwidth of the bus*

- Bisection Bandwidth (__worst__ __case__) – the __sum__ of the bandwidth of the links that __cross__ the __divide__ between the two __halves__ of a machine
  - Ring – *Two times the link bandwidth*
  - Bus – *the bandwidth of the bus*
  - Some network topologies are not __symmetric__
    - Where do we draw the line?
    - Calculate __all__ __possible__ bisection bandwidths and pick the __smallest__
    - Parallel programs are often __limited__ by the __weakest__ link in the communication chain

- Diameter – __Maximum__ __distance__ between two processors
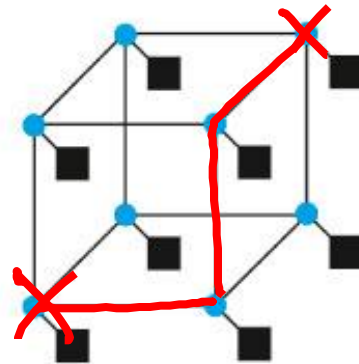
# 6.8 Network Topologies

- Each processor in a __ring__ network connects to __two__ other processors
  - Diameter of the network __$P-1$__

- In a __fully__ __connected__ network, every processor has a link to __every__ __other__ processor
  - Network bandwidth – __$P \times (P-1)/2$__
  - Bisection bandwidth – __$(P/2)^2$__
  - Diameter __1__
  - High __performance__ at high __cost__

- Many networks have been proposed between these two __extremes__ - their success largely depends on the __communication__ __workload__ of parallel programs

- Only a few have been used in __commercial__ parallel processors

# 6.8 Commercial Parallel Processor Topologies



a. 2-D grid or mesh of 16 nodes

b. *n*-cube tree of 8 nodes ($8 = 2^3$ so $n = 3$)