

Parallel Processing



Parallel Processing

0907536

Instructor: Dr. Fahed Jubair
Computer Engineering Department
University of Jordan



Parallel Processing



- Parallel processing is using *multiple processors in parallel* to solve a computation problem **more quickly** than a single processor
- Parallel computing requires parallel machines + parallel programs
 - Parallel machines (aka multiprocessors) have hardware organizations such that **multiple processors can perform multiple jobs in parallel**
 - Parallel programs are programs that **explicitly** specify how computation and data are divided among the multiple processors of a parallel machine

Class Objective



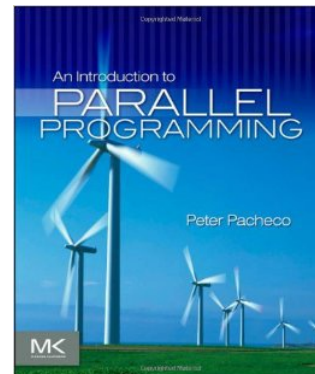
1. Study the most common parallel machine architectures
2. Shared-memory programming with OpenMP
3. Distributed-Memory programming with MPI
4. GPU programming with CUDA
5. Study some commonly-used parallel algorithms
 - Examples: parallel sorting and parallel matrix multiplication

© All Right Reserved

Textbook



Peter Pacheco, *Introduction to Parallel Programming*, 1st edition, 2011



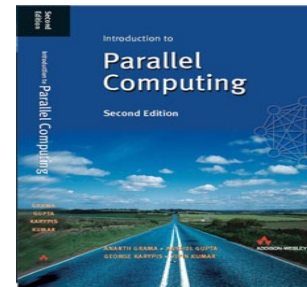
<http://www-users.cs.umn.edu/~karypis/parbook/>

© All Right Reserved

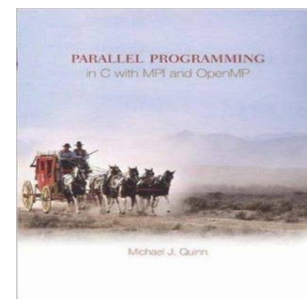
Additional References



Ananth Grama, Anshul Gupta, George Karypis, Vipin Kumar., Introduction to Parallel Computing, 2nd edition, 2010



Michael J. Quinn, Parallel programming in C with MPI and OpenMP, 2003



© All Right Reserved

Prerequisites



All students are assumed to be familiar with

- Basic data structures
- Execution analysis of algorithms
- Basic design principles of uniprocessor architecture
- Writing C and C++ programs

© All Right Reserved

Instructor Information



- Fahed Jubair
- B.Sc., University of Jordan
- Ph.D., Purdue University



© All Right Reserved

Grading Policy



- | | |
|---------------------------|-----|
| • Programming Assignments | 20% |
| • Midterm Exam | 30% |
| • Final Exam | 50% |

© All Right Reserved

Class Policy



- Attendance is important
- Do not come late
- Cheating will **NOT** be tolerated
- No makeup exams
- Course website: Microsoft teams
- Contact: fjubair@ju.edu.jo

© All Right Reserved

Questions?



© All Right Reserved

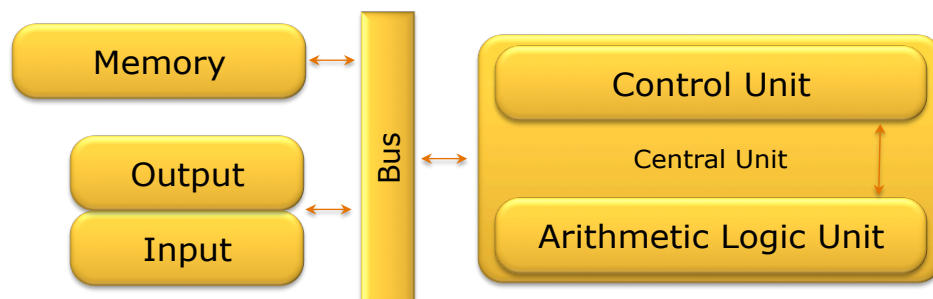
From Sequential to Parallel Computing

Instructor: Dr. Fahed Jubair
Computer Engineering Department
University of Jordan

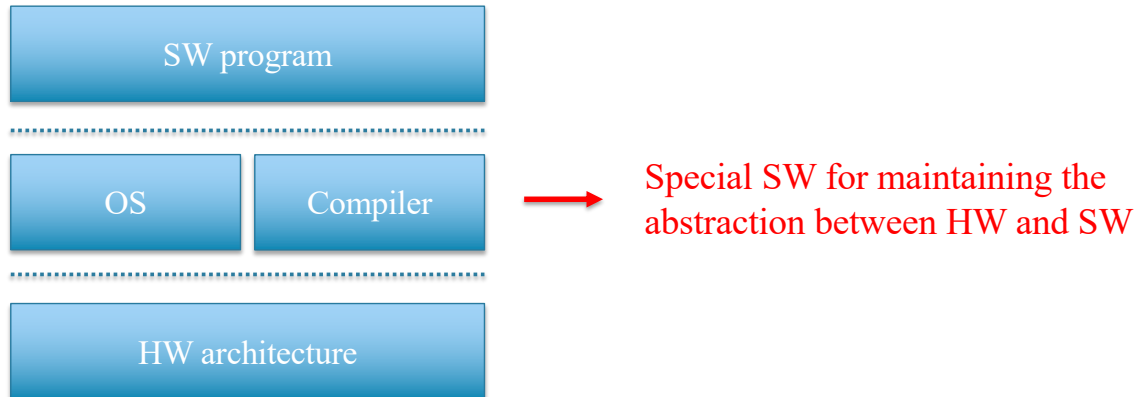
Serial Computers



- Defined by the **von Nuemann** architecture
- Consists of a central processing unit (CPU), main memory and an interconnection between the CPU and memory
- Executes (more or less) **a single job at a time**



HW-SW Abstraction in Serial Computers



© All Right Reserved.

Programmers Want Simplicity

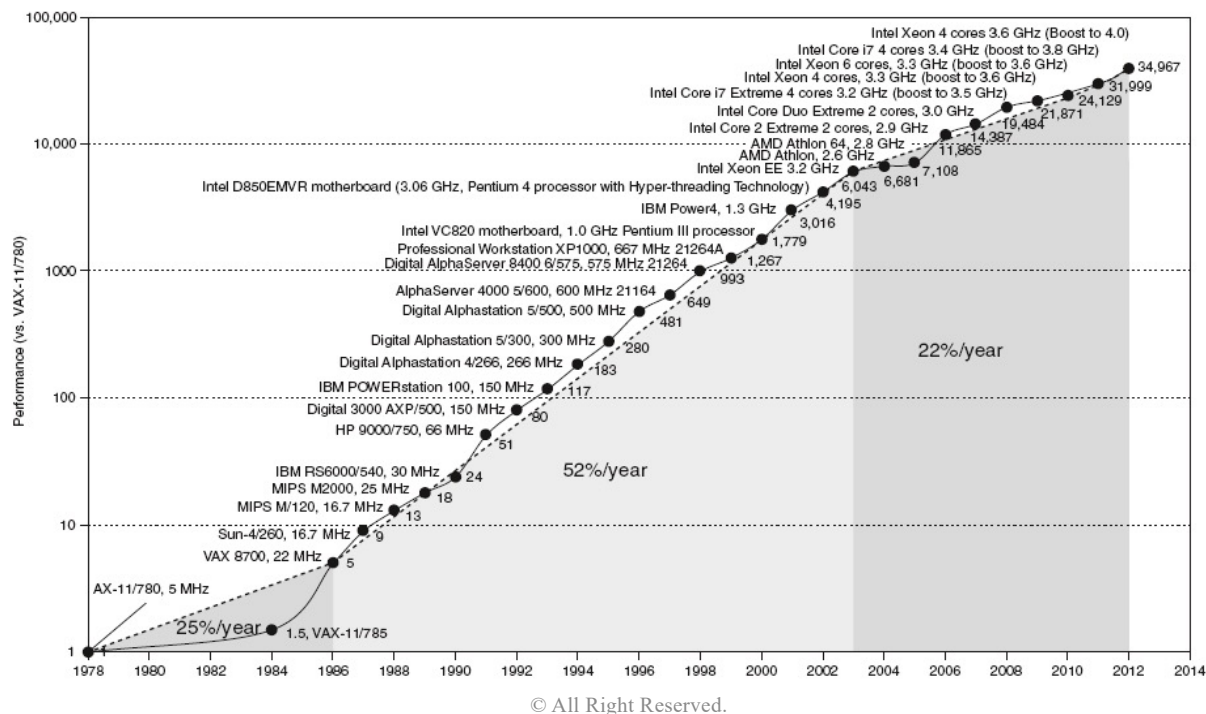


- Programmers want perfect HW-SW abstraction by writing programs that are
 - Unaware of any HW details
 - Serial, i.e., assume in-order, sequential execution
- Programmers also want their already-written programs to **automatically** run faster on **new** HW architectures without making any source-code changes

But how often do new HW architectures come out?

© All Right Reserved.

HW Technology Trend in Serial Computers



Serial Processors Are Secretly Not Serial



Due to Moore's law, HW architects were able to use more gate logic for exploiting **instruction level parallelism (ILP)** in CPUs

- Superscalar execution
 - Multiple functional units that execute instructions in parallel
- Instruction Pipelining
 - Overlapped execution of sequential instructions
- Speculative, out-of-order execution
 - Reorder instructions that have no dependency
 - Branch prediction and memory load value prediction
- Hardware multithreading
 - Interleaved execution of multiple threads on multiple functional units

Why Superscalars Are Still Serial Processors?



- Because they maintain the illusion of in-order, sequential execution to the programmer
 - “if no one saw it, then it didn't happen”
- In other words, parallelism in HW is invisible to SW
 - The HW-SW abstraction remains simple
 - Serial programs are enough 🤔

© All Right Reserved.

The Free Lunch is Over



As of 2003, the trend in doubling the clock speed per year has stopped due to the following three problems:

1. The power wall (**the main problem**)
2. The memory wall
3. The ILP wall

© All Right Reserved.

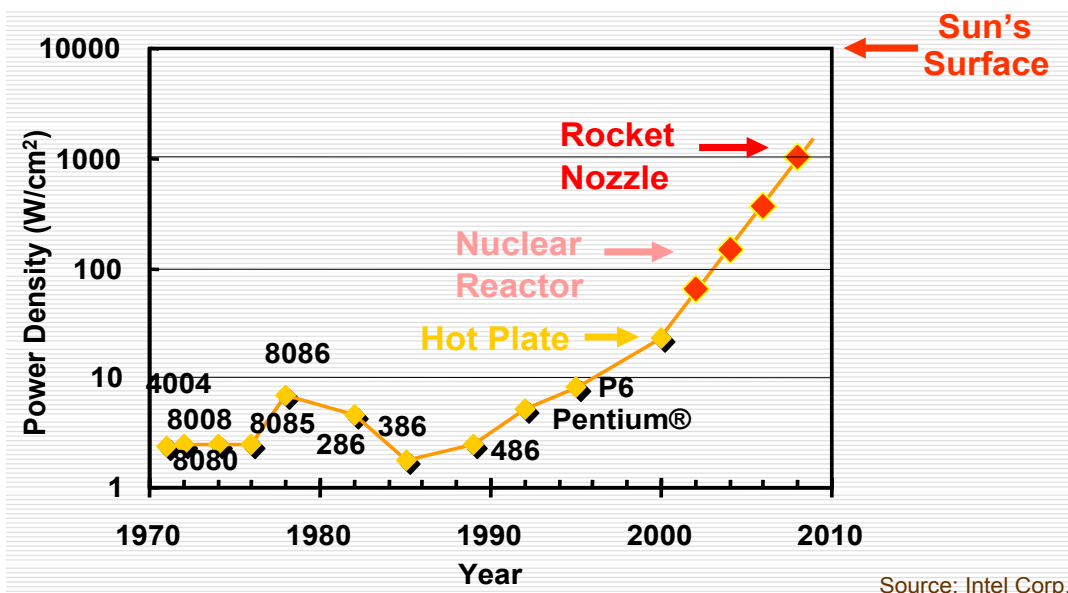


The Power Wall

- Physics problem: as clock frequency increases, **leakage power dissipation** gets exponentially worse
- Clock frequency increased by a factor of 4000 in less than two decades
- Any more significant increase in clock frequency requires heroic cooling that is not possible (integrated chips would simply melt!)

© All Right Reserved.

The Kind of Power Dissipation needed To Keep Doubling The Clock Speed

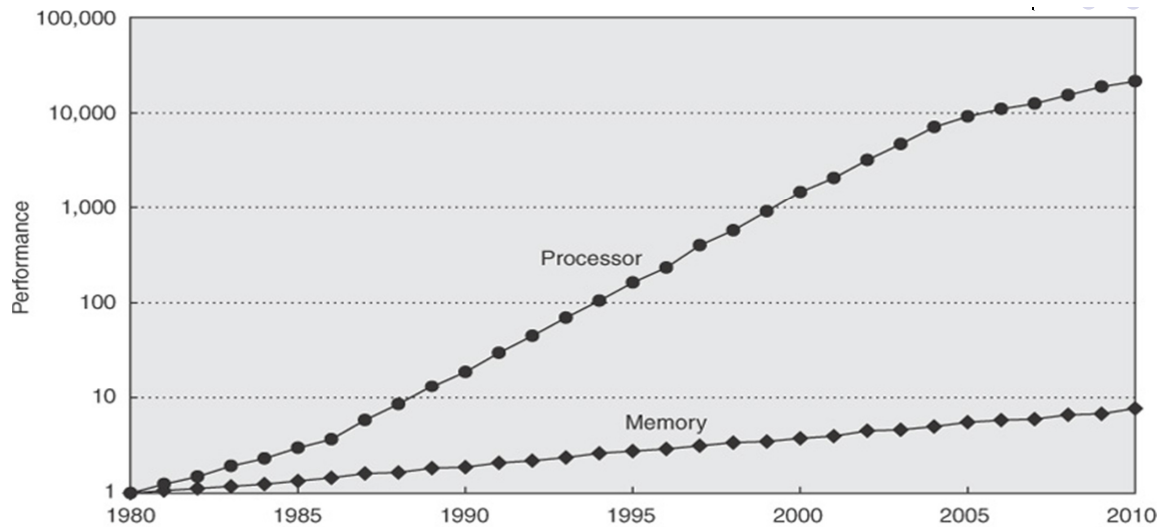


© All Right Reserved.



The Memory Wall

- Refers to the growing disparity of speed between CPU and memory outside the CPU chip



© All Right Reserved.



The Processor-DRAM Performance Gap

- From 1986 to 2000, CPU speed improved at an annual rate of 55% while memory access speed only improved at 10%
- Memory latency is a barrier to computer performance
- Current architectures have ever growing caches to bridge this gap between the CPU and memory
- Any more significant increase in CPU speed will impose a vary challenging task to bridge the gap between the CPU and memory

© All Right Reserved.

The ILP Wall



- To keep doubling the clock frequency, more aggressive ILP mechanisms are needed
- However, more aggressive ILP means deeper pipelines and more complex functional units, which makes the power leakage problem worse!

© All Right Reserved.

New Trend: Multicores



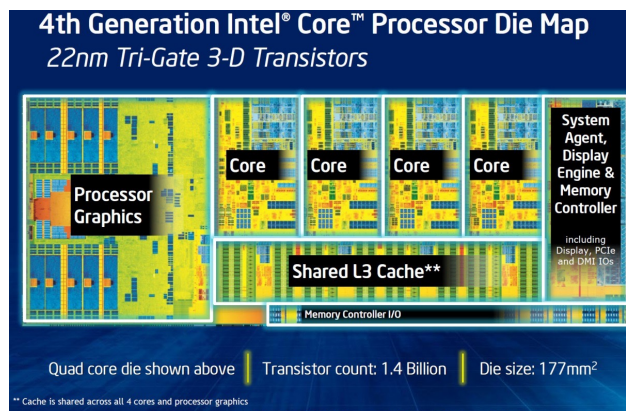
- Moore's law is still working and the number of transistors keeps increasing
- Give up the old trend: stop using these additional transistors to build a faster but more complex **uniprocessor**
- New trend: use these transistors to build **multi** simpler processors that cooperate together to perform faster
- We refer to these chips as **multicores**
- With multicores, leakage power increases linearly (not exponentially)

© All Right Reserved.

Multicore Processor Example: Intel Haswell



- June 2013
- 22 nm technology
- 1.4 billion transistors
- 4 cores
- 3.1 GHz
- 6 MB L3 Cache
- Hyperthreading



© All Right Reserved.

Sequential Programming No longer works



- All modern processors are multicores
- With multicores, we can simultaneously execute more than one job at a time (thus, improving throughput)
- However, multicores cannot automatically make the execution time (latency) of a serial program faster
- To make programs run faster, we need to **rewrite** them into an *explicitly parallel programs*
 - **The era of sequential programming is over!**

© All Right Reserved.

The Rise of Parallel Programming



- Explicit parallel programs **explicitly** specify how concurrent computation is performed by the multiple processors of a parallel machine
- **The merit:** parallel programs achieve faster computation
- **The problem:** parallel programs are *significantly* more difficult to write than sequential programs

© All Right Reserved.

Example: Parallel Sum



- Let us develop a parallel algorithm (written in pseudo-code) to compute the sum of 1000 integers stored in an array
 - Assume the algorithm will run on a multicore machine with 8 cores: core 0, core 1, ..., core 7
 - Let us use the term “thread” to refer to each sequence of instructions executed on a core
 - Therefore, we have 8 threads: thread 0, thread 1, ..., thread 7 such that thread i executes on core i
 - Let the variable *thread_id* refer to the thread number

Note that the sequential version is trivial

```
sum = 0
for ( i = 0; i < 1000; i++)
    sum = sum + A[ i ]
```

© All Right Reserved.

Example (cont.)

Parallel Sum Algorithm



- We have 1000 integers and 8 threads \Rightarrow make each thread computes the sum of 125 integers
- Our parallel sum algorithm works as follows:
 - i. First, each thread computes its partial sum independently
 - ii. Then, combine partial sums from each thread into a global sum

© All Right Reserved.

Example (cont.)

First Step: Partition The Work



- Let variables *lb* and *ub* determine the lower and upper array bounds for each thread
- Let variable *my_sum* determines the partial sum computed by each thread
- Note that each thread has its own **private** copy of *lb*, *ub* and *my_sum*

```
lb = 125 * thread_id
ub = lb + 124
my_sum = 0
for ( i = lb ; i < ub ; i++)
    my_sum = my_sum + A[ i ]
```

Wait! how does this code get executed by the 8-core parallel machine?

© All Right Reserved.

Example (cont.) Execution Model



- The same program is **duplicated** among threads
- This paradigm is referred to as the *single-program-multiple-data (SPMD)*

```
lb = 125 * thread_id
ub = lb + 124
my_sum = 0
for ( i = lb ; i < ub ; i++)
    my_sum = my_sum + A[ i ]
```

At runtime

Thread 0	lb = 0 , ub = 124
Thread 1	lb = 125 , ub = 249
Thread 2	lb = 250 , ub = 374
Thread 3	lb = 375 , ub = 499
Thread 4	lb = 500 , ub = 624
Thread 5	lb = 625 , ub = 749
Thread 6	lb = 750 , ub = 874
Thread 7	lb = 875 , ub = 999

© All Right Reserved.

Example (cont.) Second Step: Combine Partial Sums



- Let variable *global_sum* be a global variable that each threads updates with its partial sum
- We will let thread 0 be responsible for initializing *global_sum*
- Note that there is a single copy of *global_sum* that is **shared** between all threads

```
if ( thread_id = 0 )
    global_sum = 0
lb = 125 * thread_id
ub = lb + 124
my_sum = 0
for ( i = lb ; i <= ub ; i++)
    my_sum = my_sum + A[ i ]
global_sum = global_sum + my_sum
```

- Note that the order at which threads update *global_sum* is **random**
- Here, the order is not important because **addition is associative**

© All Right Reserved.

Example (cont.)

Race Condition Scenario



```
if ( thread_id = 0 )
    global_sum = 0
lb = 125 * thread_id
ub = lb + 124
my_sum = 0
for ( i = lb ; i <= ub ; i++ )
    my_sum = my_sum + A[ i ]
global_sum = global_sum + my_sum
```

Assume threads 1 & 3 arrived here simultaneously

Also assume the current status of global sum and partial sums is:

my_sum (thread1) = 20

my_sum (thread3) = 33

$global_sum$ = 60

© All Right Reserved.

Example (cont.)

Race Condition Scenario



1	Thread 1 loads $global_sum$ from memory into register (i.e., R) $R = 60$	Thread 3 loads $global_sum$ from memory into register (i.e., R) $R = 60$
2	Thread 1 adds my_sum to R $R = 60 + 20 = 80$	Thread 3 adds my_sum to R $R = 60 + 33 = 93$
3	Thread 1 stores R into $global_sum$	Thread 3 stores R into $global_sum$

- Depending on which threads stores its value last, $global_sum$ could be either 80 or 93
- This is called a **race condition**, which occurs when at least two threads simultaneously access a shared variable and at least one of these threads is performing a write operation

© All Right Reserved.

Example (cont.) Synchronization



- To eliminate race conditions, programmers must **synchronize** between threads by forcing them to **sequentially** update the global sum by their partial sums
- Therefore, we will use a *critical section* as our form of synchronization
- A critical section is executed by a thread **atomically**, i.e., uninterrupted by any other thread

© All Right Reserved.

Example (cont.) Critical Section



```
if ( thread_id = 0 )
    global_sum = 0
lb = 125 * thread_id
ub = lb + 124
my_sum = 0
for ( i = lb ; i <= ub ; i++ )
    my_sum = my_sum + A[ i ]

critical section {
    global_sum = global_sum + my_sum
}
```

© All Right Reserved.

Example (cont.)

Another Race Condition Scenario



```
if ( thread_id = 0 )
    global_sum = 0
lb = 125 * thread_id
ub = lb + 124
my_sum = 0
for ( i = lb ; i <= ub ; i++)
    my_sum = my_sum + A[ i ]
critical section {
    global_sum = global_sum + my_sum
}
```

Thread 0 has the responsibility of initializing *global_sum* to 0

A race condition could occur if thread 0 is too slow that some other threads reached the critical section before thread 0 initializes *global_sum* to 0

© All Right Reserved.

Example (cont.)

Barriers



- One way to eliminate the data race shown in the previous slide is to force all threads to wait for thread 0 to initialize *global_sum* to 0
- To do so, we will use a *barrier*, which is a form of synchronization where all threads must wait until all other threads reach the same point in execution
- The hardware architecture must support the thread-waiting mechanisms of barriers

© All Right Reserved.

Example (cont.) Final Parallel Code



```
if ( thread_id = 0 )
    global_sum = 0

threads_barrier

lb = 125 * thread_id
ub = lb + 124
my_sum = 0
for ( i = lb ; i <= ub ; i++)
    my_sum = my_sum + A[ i ]

critical section {
    global_sum = global_sum + my_sum
}
```

Very different from the
trivial sequential code!

© All Right Reserved.

Example (cont.) Performance Analysis



- Ideally, our parallel algorithm aims to achieve 8-times faster execution than sequential version (i.e., achieves *linear speedup*)
- But parallelism has overheads:
 - Work partitioning overhead: extra computation needed to partition computation among threads
 - Synchronization overhead:
 - Barriers force fast threads to wait for slow threads
 - Critical sections are executed serially
 - The creation and termination of threads also have an overhead
- To write efficient parallel programs, programmers need to think about how to minimize parallelism overheads in their codes

© All Right Reserved.

Example (cont.) Scalability



- Ideally, we want our algorithm to perform faster when executed on a parallel machine with **higher number of cores**
 - i.e., our algorithm needs to be *scalable*
- The problem is parallelism overhead, which increases as the number of cores (or threads) increases
- This imposes yet another challenge to programmers: how to write *scalable* parallel programs!

© All Right Reserved.

Challenges in Writing Parallel Programs



1. Programmers need to **think parallel**
 - Rewrite the sequential algorithm into a parallel algorithm
 - In some cases, the parallel algorithm is a totally new algorithm
2. Programmers need to understand the architecture of the underlying parallel machine (**HW-SW abstraction is gone!**)
 - HW details affect the programmer's decision on how to find efficient partitioning of computation between threads
 - HW details affect the programmer's decision on how to choose efficient synchronization methods

© All Right Reserved.

Challenges in Writing Parallel Programs



3. Testing and debugging parallel programs is hard
 - For example, some race conditions are very subtle that may not occur for the first 999999 runs, and occurs on the 1000000th run
4. Parallelism **has an overhead**
 - If not properly handled, this overhead can kill the scalability of a parallel program

© All Right Reserved.

How To Write A Parallel Program?



- In general, there are three approaches for writing an explicit parallel program:
 1. Write from scratch
 2. Parallelize an existing sequential program
 3. Parallelize an existing sequential program by a compiler
- } Done manually

© All Right Reserved.

Wait! Why Not Always Use A Parallelizing Compiler?



- A parallelizing compiler converts a sequential program into a parallel program
- The merit: it is automatic!
- The problem: so far, parallelizing compilers have shown limited success due to the significant complexity of creating parallel programs
- Examples on existing parallelizing compilers: SUIF, PLUTO, Rose, Cetus, Polaris, and Intel C++ compiler

© All Right Reserved.

Parallel Program Types



Parallel applications can be classified depending on how problem partitioning is performed, as follows:

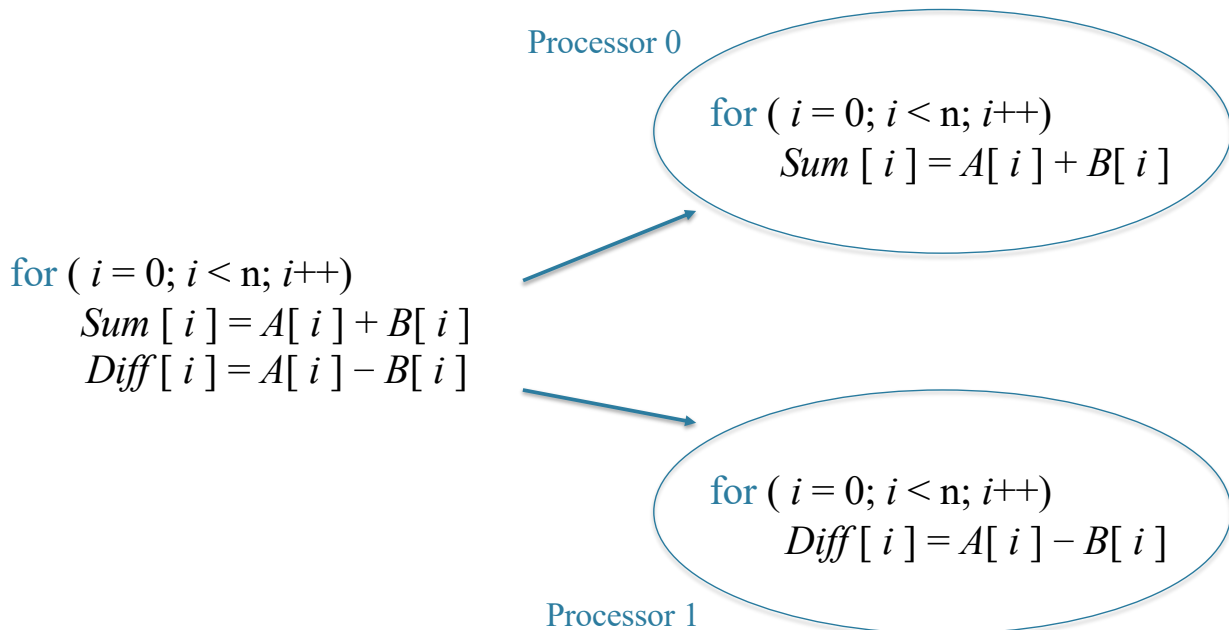
- Task-parallel applications
 - Partition various tasks carried out solving the problem among the processors
- Data-parallel applications
 - Partition the data used in solving the problem among the processors
 - Each processor carries out similar operations on it's part of the data

© All Right Reserved.



Example

Work Partitioning – Task Parallelism

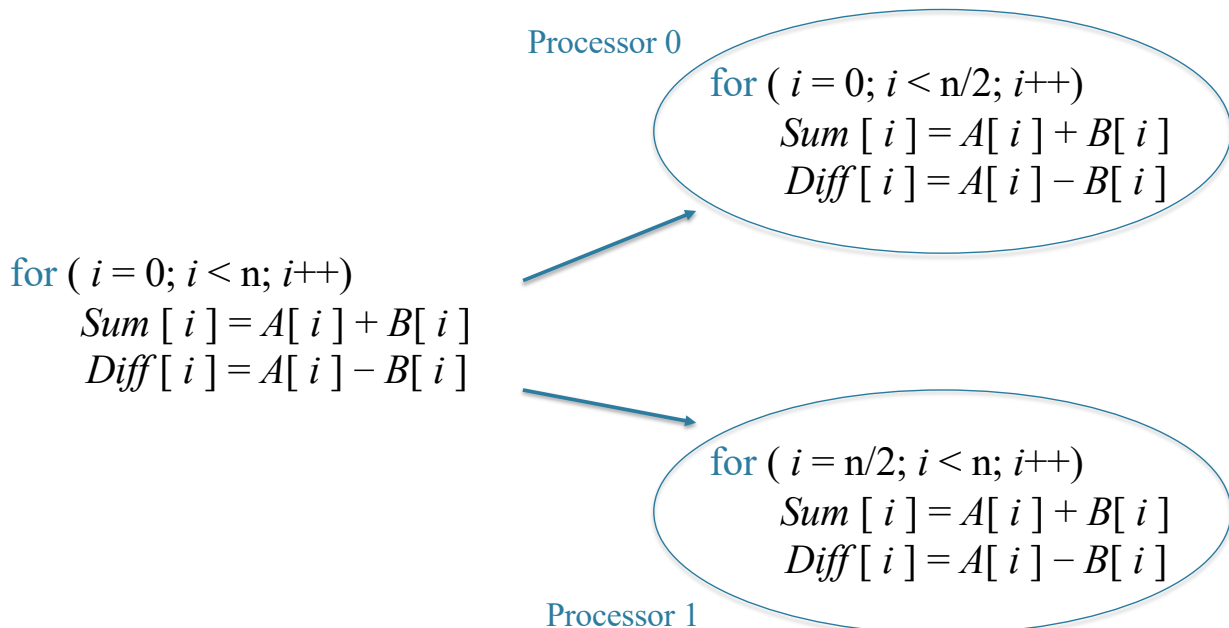


© All Right Reserved.



Example

Work Partitioning – Data Parallelism



© All Right Reserved.

Terminology

Concurrency vs Parallelism



- **Concurrency**: capability of a system to have two or more activities in progress at the same time
 - Activities can be independent
 - Do not necessarily demand parallel hardware but demands OS to support concurrency
- **Parallelism**: capability of a system to execute multiple activities simultaneously
 - Demands parallel hardware and concurrency support
 - Any parallel program is a concurrent program

© All Right Reserved.

Terminology

Programming Model



- In this course, we will learn *programming models* for parallel machines
- A programming model abstractly describes how a program for a parallel machine should be written
- A programming model can be implemented by
 - Using a new programming language
 - Extending an existing sequential programming language
 - A library that is invoked from a sequential code

© All Right Reserved.

Summary



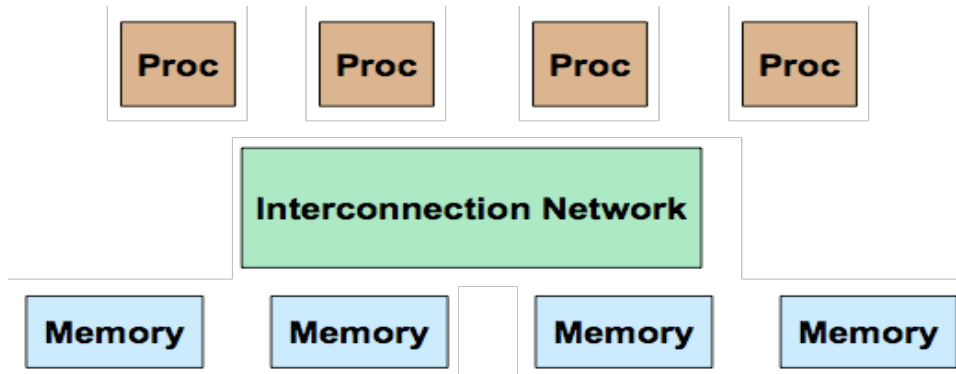
- The laws of physics have brought us to the doorstep of multicore technology
- Serial programs typically don't benefit from multiple cores
- Writing parallel programs is needed to achieve high-performance execution on parallel machines
- Parallel programs are usually very complex and therefore require advanced programming techniques and development

© All Right Reserved.

Parallel Hardware

Instructor: Dr. Fahed Jubair
Computer Engineering Department
University of Jordan

An Abstract Parallel Machine



In general, parallel machines differ in their

- Concurrency model, which affects how instruction stream and data stream are partitioned among processors
- Communication model, which affects how processors and memory units are connected by the interconnection network

Flynn's Taxonomy



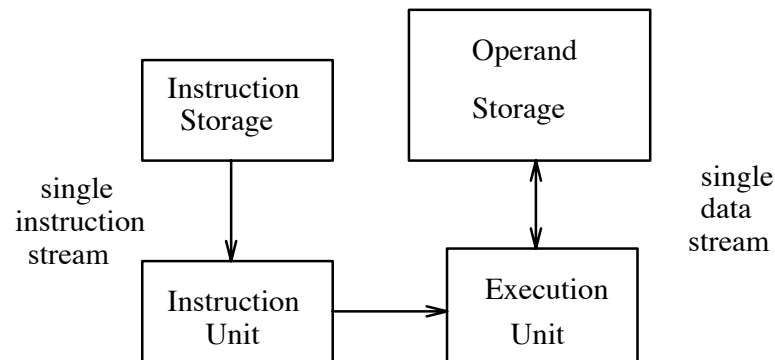
Processor architectures, in general, are classified based on instruction stream and data stream, as follows:

<p><i>serial</i></p> <p>SISD</p> <p>Single instruction stream Single data stream</p>	<p><i>parallel</i></p> <p>(SIMD)</p> <p>Single instruction stream Multiple data stream</p>
<p><i>mostly theoretical</i></p> <p>MISD</p> <p>Multiple instruction stream Single data stream</p>	<p><i>parallel</i></p> <p>(MIMD)</p> <p>Multiple instruction stream Multiple data stream</p>

Serial Computing Platforms



- **SISD processors**
 - Von Neumann uniprocessor architecture



© All Rights Reserved.

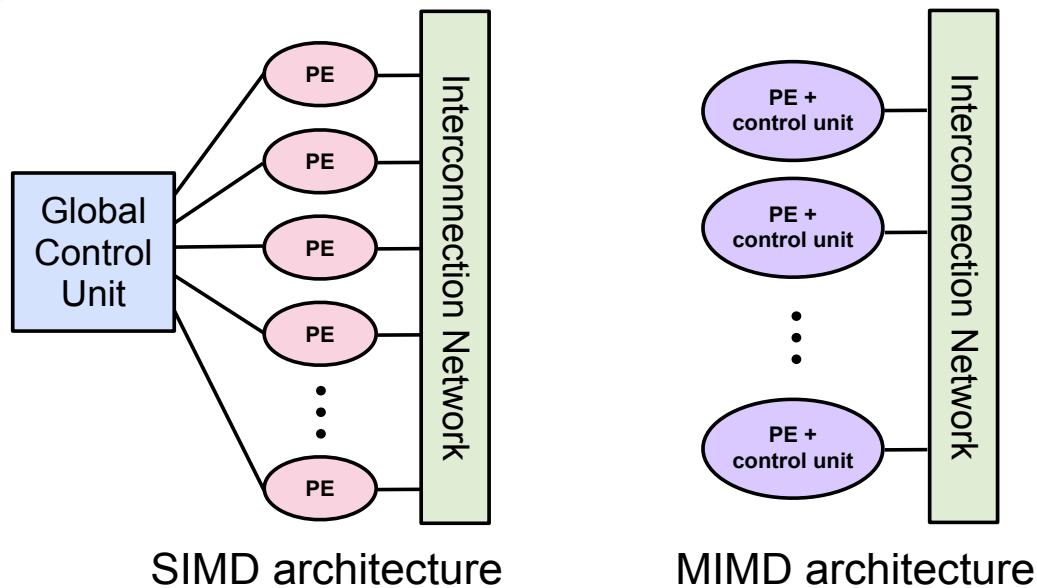
Parallel Computing Platforms



- **SIMD processors**
 - All processors execute the same instruction simultaneously, but while operating on different pieces of data
 - Work best with data-parallel applications
- **MIMD processors**
 - All processors execute in parallel but on different instructions and different pieces of data
 - Can handle both data-parallel and task-parallel applications

© All Rights Reserved.

SIMD and MIMD Processors



PE = Processing Element
© All Rights Reserved.

SIMD vs. MIMD



- SIMD platforms
 - Single control unit
 - Only one copy of program
 - Special purpose: not suited for all applications
- MIMD platforms
 - Multiple control units
 - Multiple programs
 - Suitable for broad range of applications



SIMD Example

Conventional SISD loop that sums two arrays

```
for ( i=1; i<=N; i++ )  
    sum[ i ] = A[ i ] + B[ i ]
```



SIMD code has a single instruction that executes simultaneously on N ALUs

```
sum[ 1 : N ] = A[ 1 : N ] + B[ 1 : N ]
```

© All Rights Reserved.



SIMD Processor Example: Vector Processor

- A vector processor executes *vector instructions*, which are instructions that operate on one-dimensional arrays of data called *vectors*
 - In contrast, conventional processors execute instructions that operate on individual data elements (called *scalars*)
- Vector processors have
 - Vectorized and pipelined functional unit: the same operation is applied to each element in the vector (or pairs of elements)
 - Vector registers: capable of storing a vector of data elements
- **The merit:** fast
- **The problem:** limited to data-parallel applications

© All Rights Reserved.

Vector Processors in Modern Architectures



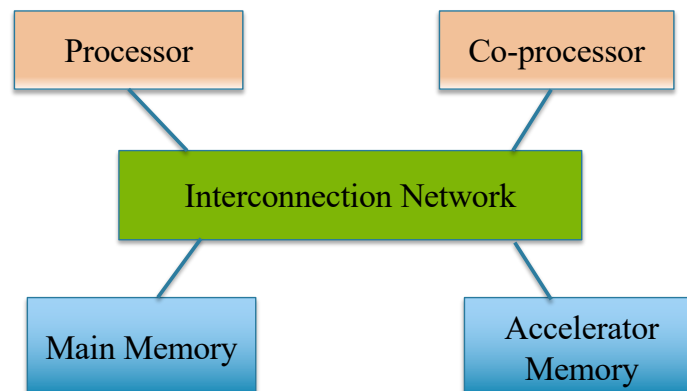
- Today's processors usually have SIMD units or vector processors as co-processors (aka *accelerators*)
- On-chip SIMD units:
 - For example, modern Intel processors have *on-chip graphical processing units (GPUs)* for fast image processing
 - Intel *instruction set architecture* has special instructions for multimedia support (e.g., referred to as multimedia extensions or MMX)
- Off-chip co-processors:
 - Typically sold as PCIe cards
 - Architectures that have both conventional processors and co-processors are referred to as *heterogeneous architectures*

© All Rights Reserved.

Heterogeneous Architecture



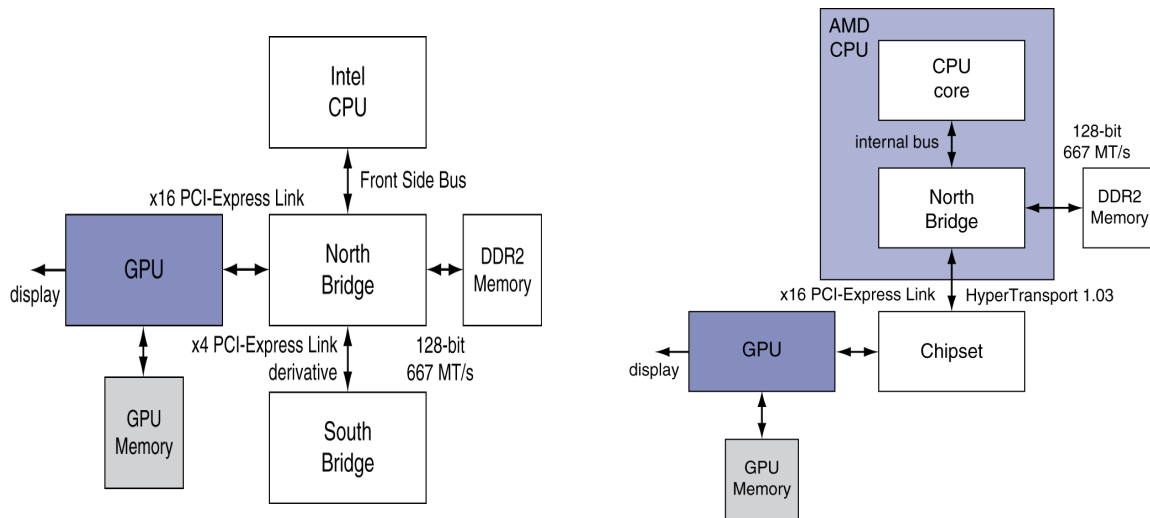
Abstract
heterogeneous
architecture



- Heterogeneous platform: Processor (on-chip) + Co-Processor (on-chip or off-chip)
- In contrast, a homogeneous platform has only on-chip processors

© All Rights Reserved.

Heterogeneous Architecture Example: Off-Chip GPUs



© All Rights Reserved.

Accelerator Example: NVIDIA Tesla K40 (July 2013)



- 2880 streaming cores
- 732 MHz
- 12GB memory
- Memory B/W 288 GB/s
- 4.29 TFLOP for single-precision
- 1.43 TFLOPS for double-precision
- 235 Watts
- Kepler GK110 architecture



© All Rights Reserved.

Achieving High-Performance Computing On Accelerators



- Programmers write an explicit parallel program that uses the conventional CPU for executing the sequential parts of the program and the accelerators for executing the parallel parts of the program (generally speaking)
 - The term *heterogeneous computing* refers to a parallel program being executed on a heterogeneous architecture
- *OpenCL*, *OpenACC*, and *CUDA* are widely-used programming frameworks for accelerators

© All Rights Reserved.

MIMD Processors



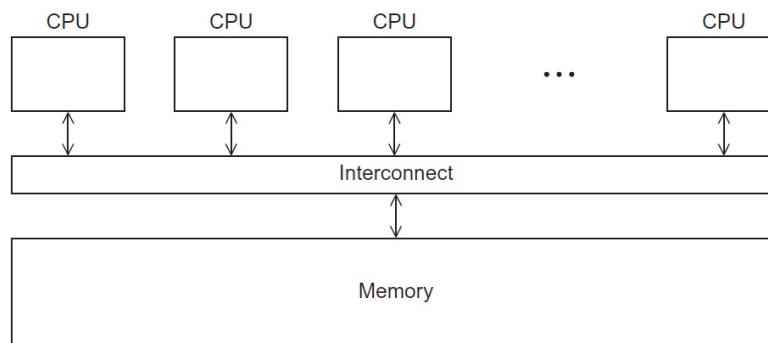
- **Execute different programs on different processors**
- MIMD processors are classified based on their communication model, as follows:
 - Shared-memory architectures
 - All processors can **directly** access a common memory
 - Processors **communicate implicitly** by reading or writing shared data structures in memory (for example, via load and store instructions)
 - Distributed-memory architectures:
 - Each processor can only access **local** memory
 - Processors **communicate explicitly** by exchanging messages (for example, through a LAN network)

© All Rights Reserved.

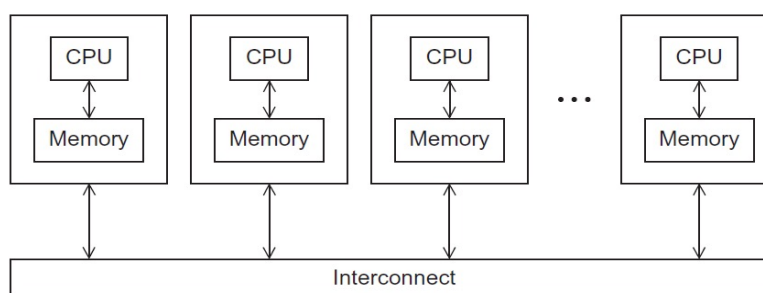
Abstract View of MIMD Architectures



Shared-memory architecture



Distributed-memory architecture



© All Rights Reserved.

Shared-Memory vs Distributed-Memory



	Shared-memory	Distributed-memory
Execution units	Threads	Processes
Memory sharing	All threads share the same data space	Each process has its own data space
Operating system	Typically, all threads run within the same OS	Typically, each process runs within a separate OS
Communication method	Threads communicate implicitly by accessing shared variables in memory	Processes communicate explicitly by exchanging messages
Communication speed	Communication between processors in shared-memory is typically faster than distributed-memory	
Scalability	Distributed-memory platforms are typically more scalable (have higher number of processors) than shared-memory platforms	

© All Rights Reserved.

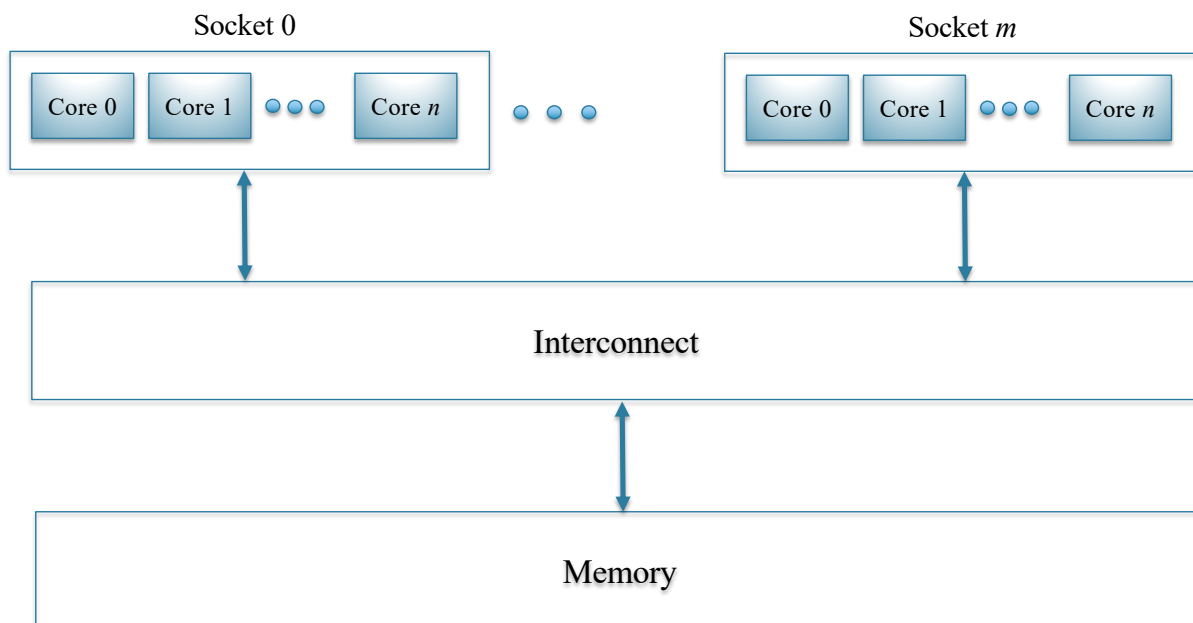
Shared-Memory Platforms



- Most widely available shared-memory architectures are *multicores* (aka *chip multiprocessor* or *CMP*)
- Multicores have multiple processing units (cores) that are packaged within the same chip and share the same physical memory
- Multicore architectures have two types:
 - **UMA (Uniform Memory Access)**
 - Time taken by a core to access any memory word is the same
 - **NUMA (Non-Uniform Memory Access)**
 - Time taken by a core to access memory words vary

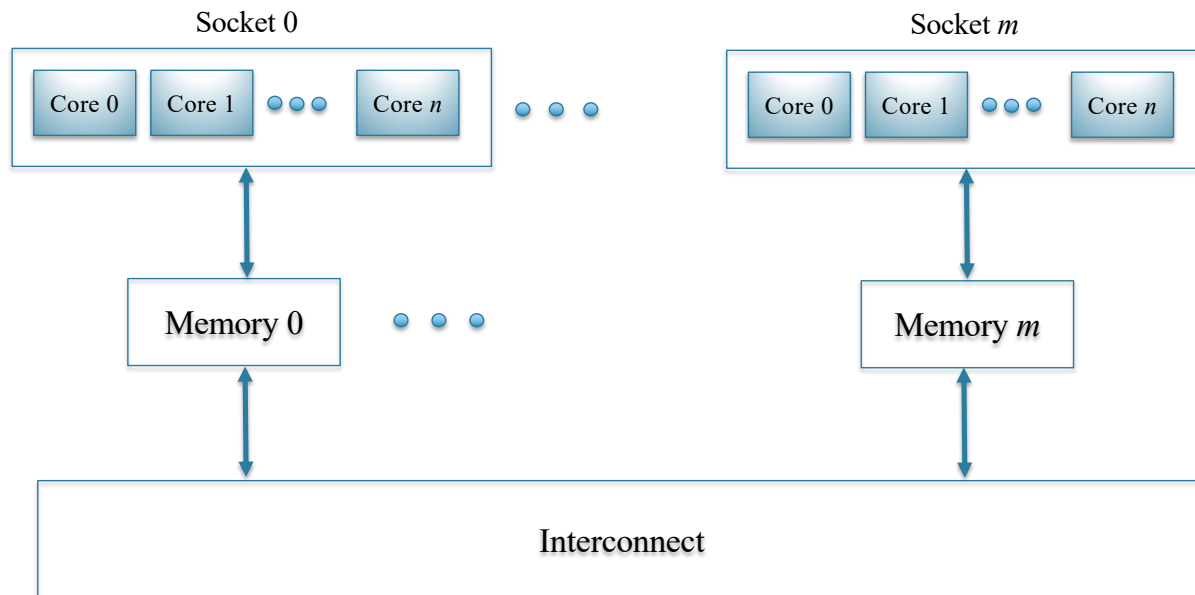
© All Rights Reserved.

Abstract View of UMA



© All Rights Reserved.

Abstract View of NUMA



© All Rights Reserved.

UMA vs NUMA



UMA

- Data placement is unimportant
⇒ simpler HW and SW
- Less scalable

NUMA

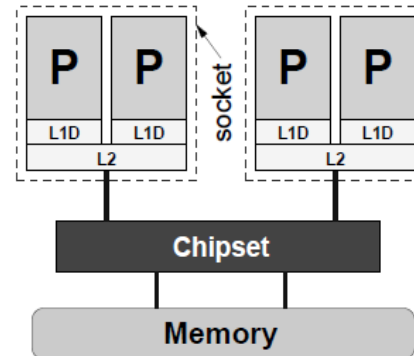
- Data placement affects performance
⇒ more complex HW and SW
- More scalable

© All Rights Reserved.



UMA Example

- Two dual-core sockets
- Each core has its own L1 cache
- The two cores of each socket share L2 cache

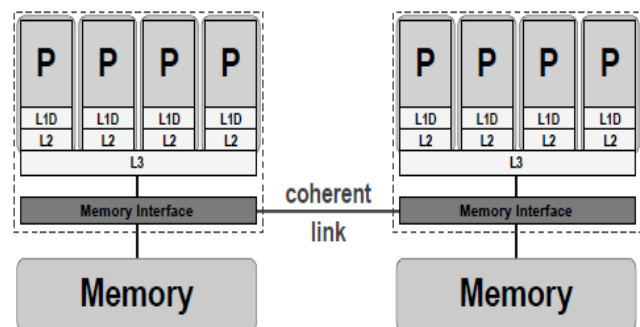


© All Rights Reserved.



NUMA Example

- Two quad-core sockets
- Each core has its own L1 and L2 caches
- The four cores of each socket share L3 cache



© All Rights Reserved.

Interconnection Networks In Shared-Memory Platforms



- The interconnect plays a significant role in performance
- Even with fast processors and memory, slow interconnect will likely degrade the overall performance
- Two widely used interconnect technologies in shared memory platforms:
 - Bus interconnect
 - Switched interconnect

© All Rights Reserved.

Bus Interconnect



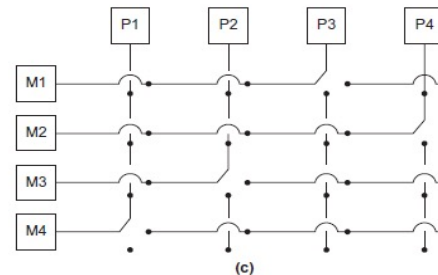
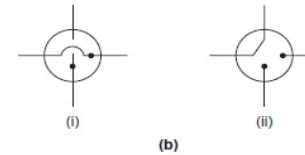
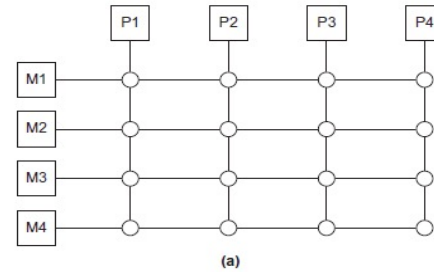
- A collection of parallel communication wires together with some hardware that controls access to the bus
- Communication wires are shared by the devices that are connected to it
- **The merit:** low cost
- **The problem:** limited scalability
 - As the number of devices connected to the bus increases, contention for use of the bus increases, and performance decreases

© All Rights Reserved.



Switched Interconnect

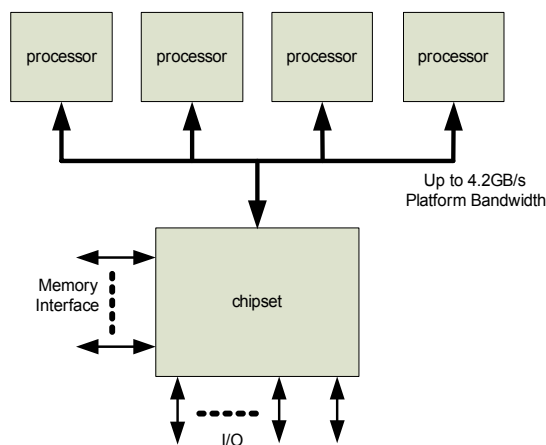
- Uses switches to control the routing of data among the connected devices
- Crossbar interconnect
 - Allows simultaneous communication among different devices
 - Faster than buses
 - But the cost of the switches and links is relatively high



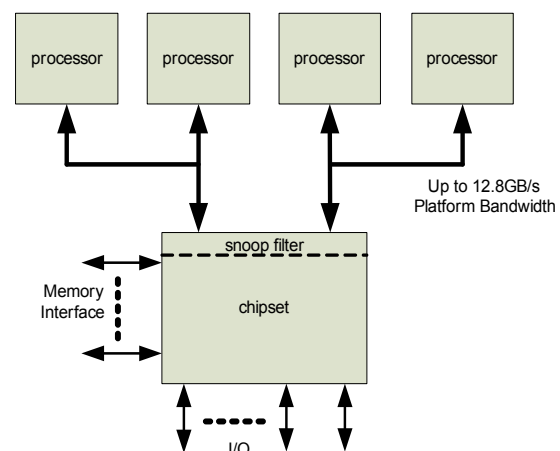
© All Rights Reserved.



Example: Intel Interconnect Evolution



Traditional Shared Frontside Bus
(until 2004)

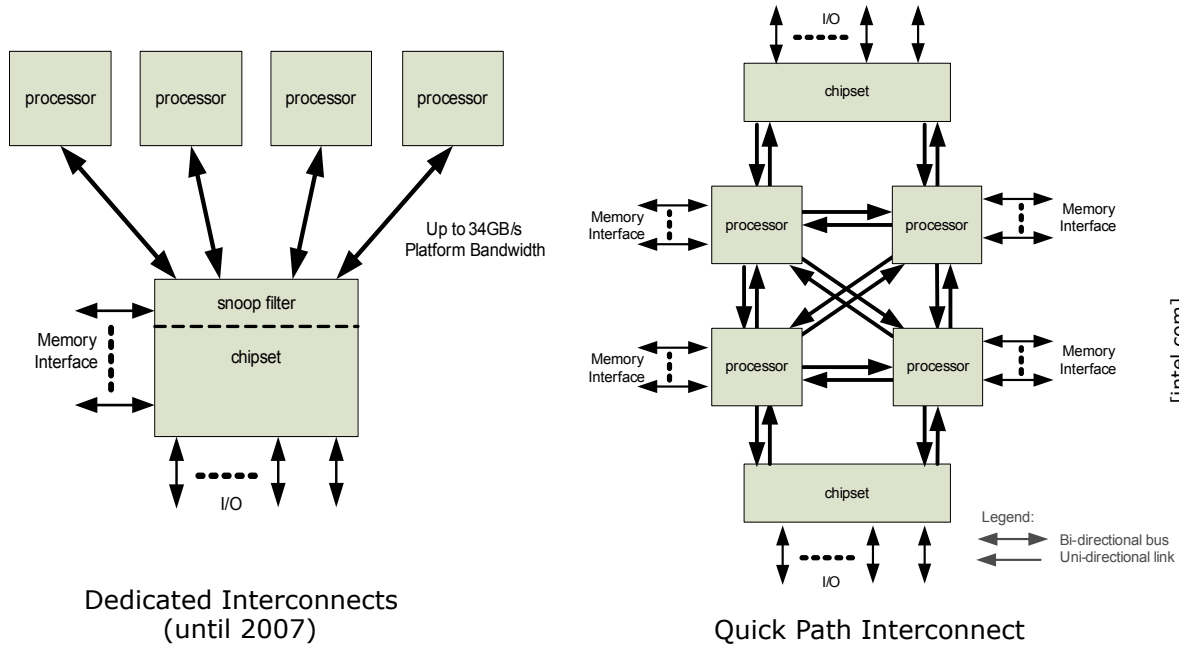


Dual Independent Buses
(until 2005)

© All Rights Reserved.



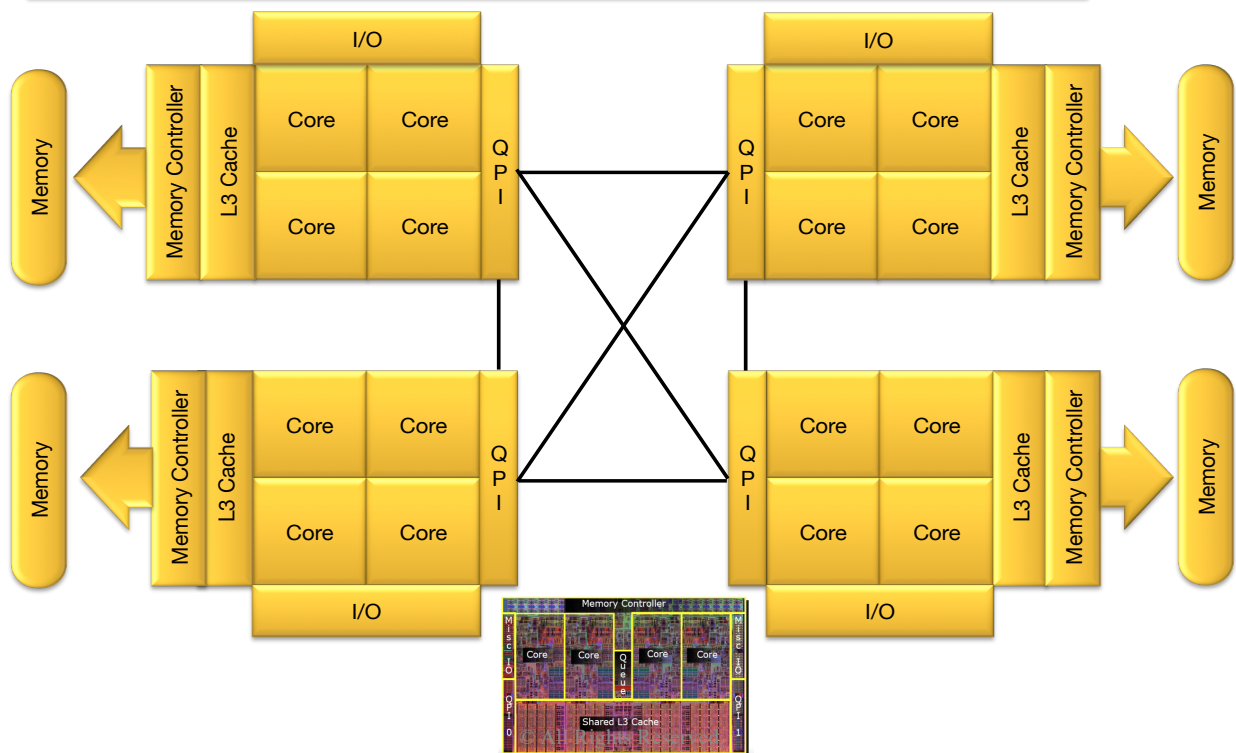
Example: Intel Interconnect Evolution



© All Rights Reserved.



Shared-Memory Processor Example: Intel Nehalem (Nov 2008)



Shared Memory Synchronization



- Synchronization is a special code that describes how to coordinate sharing among threads for correct execution
- Generally, there are three types of synchronizations:
 - **Mutual exclusion**: a code region that threads can only execute atomically, i.e., one thread at a time
 - **Barriers**: a point in execution where all threads must wait for every thread to reach this point
 - **Signal-wait mechanism**: a pairwise coordination between two threads where one thread stops its execution and waits for a signal from the other thread to continue its execution

© All Rights Reserved.

Synchronization Examples



Signal-wait mechanism

<pre>T0 A = 1 flag = 1</pre>	<pre>T1 While (flag==0) { } print A</pre>
------------------------------	---

T1 repeatedly tests a condition till it becomes true

mutual exclusion

```
while (true){
    if (LOCK(&t) == true) break ;
}
.... // critical code
UNLOCK(&t);
```

Each thread repeatedly tries to **acquire** the lock till its successful

Only one thread enters the critical section

When finishing the critical section, **release** the lock

© All Rights Reserved.

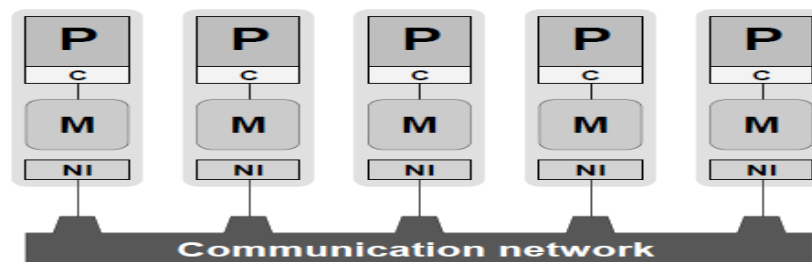
Achieving High-Performance Computing On Multicores



- Programmers write a *multithreading program*, which is a parallel program that runs on multicores
 - Programmer partitions the computation among threads
 - Programmer insert synchronization operations to coordinate data sharing between threads
 - Communication between threads is done implicitly through the memory
- Several programming models are available for shared-memory platforms, such as *Pthreads*, *OpenMP*, *Cilk*, and *Intel TBB*

© All Rights Reserved.

Distributed-Memory Platforms



- A set of computers that are connected by a network
- Each computer runs a *local program on a local memory*
- Processors **do not** share a common memory, instead they communicate by exchanging send/receive messages through the network
- Network Interfaces (NI) mediate the connections to the network

© All Rights Reserved.

Properties of Distributed-Memory Platforms



Merits:

- Allows an easy scale-out: simply add new machines to the network
- Allows simpler hardware: the complications of sharing a physical memory (such as implementing synchronization operations) do not exist
- Inexpensive: you can connect simple uniprocessors using a cheap network to build a parallel machine
- Dependable: Operational nodes can replace faulting nodes because all nodes are independent

Problems:

- Messaging has longer latency than memory access
- Programmers need to insert communication messages in their codes

© All Rights Reserved.

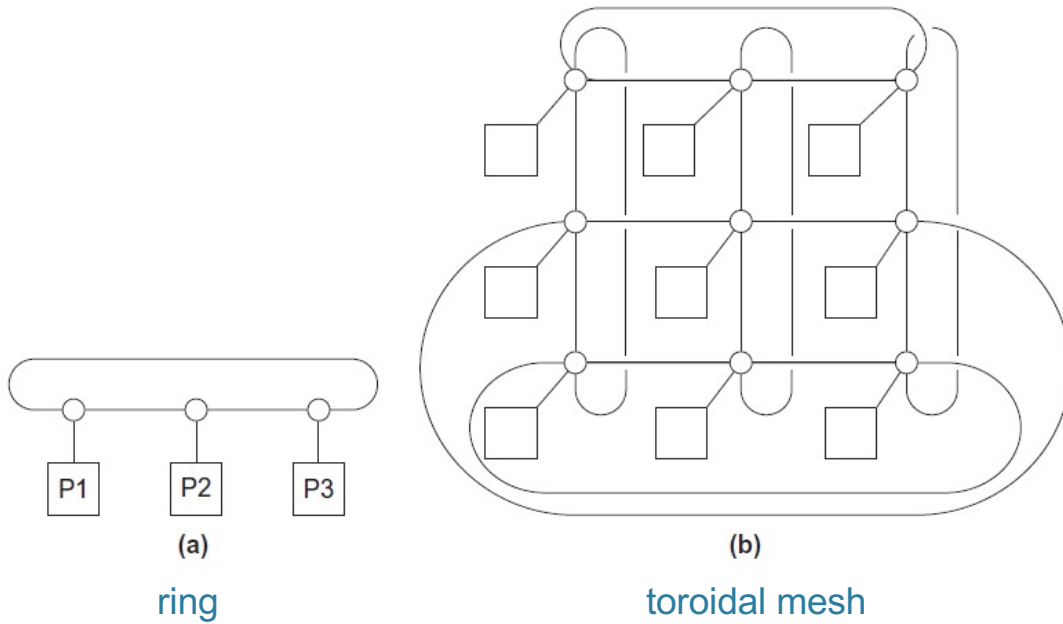
Interconnection Networks In Distributed-Memory Platforms



- Divided into two groups
 - Direct interconnect
 - Each switch is directly connected to a processor memory pair, and the switches are connected to each other
 - Indirect interconnect
 - Switches may not be directly connected to a processor

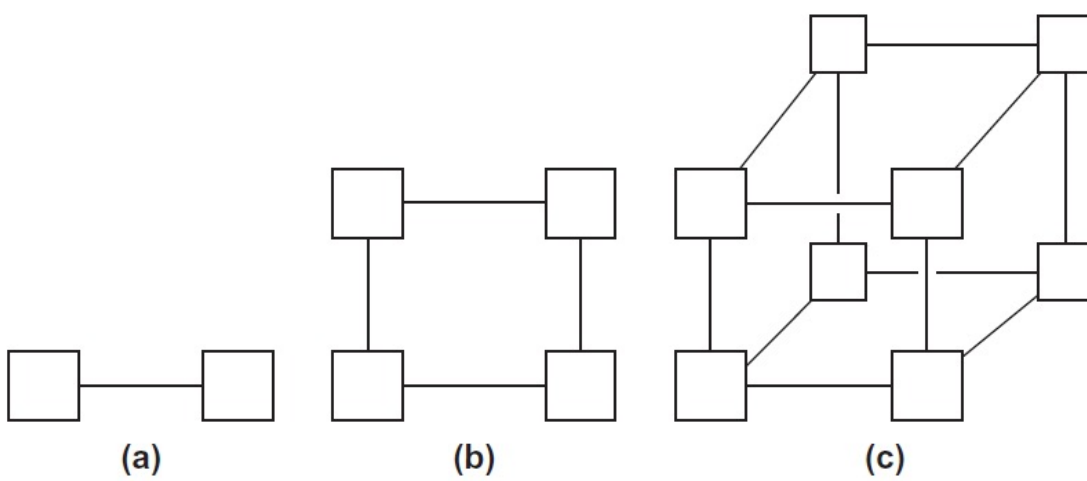
© All Rights Reserved.

Direct Interconnect



© All Rights Reserved.

Direct Interconnect (cont.)



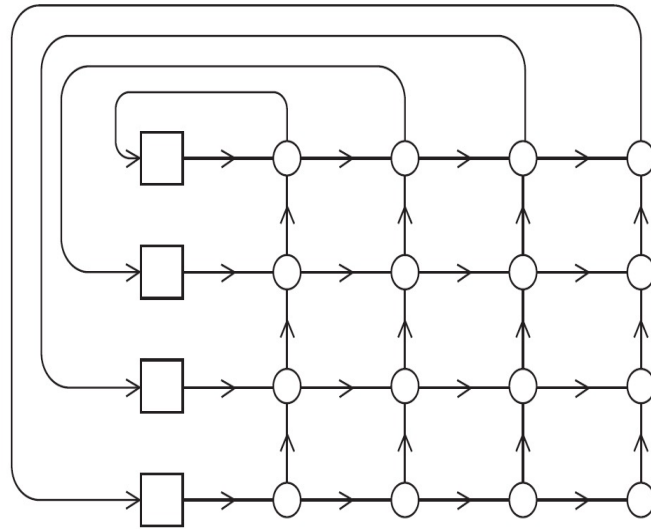
(a) One-, (b) two-, and (c) three-dimensional hypercubes

© All Rights Reserved.

Indirect Interconnect



Crossbar

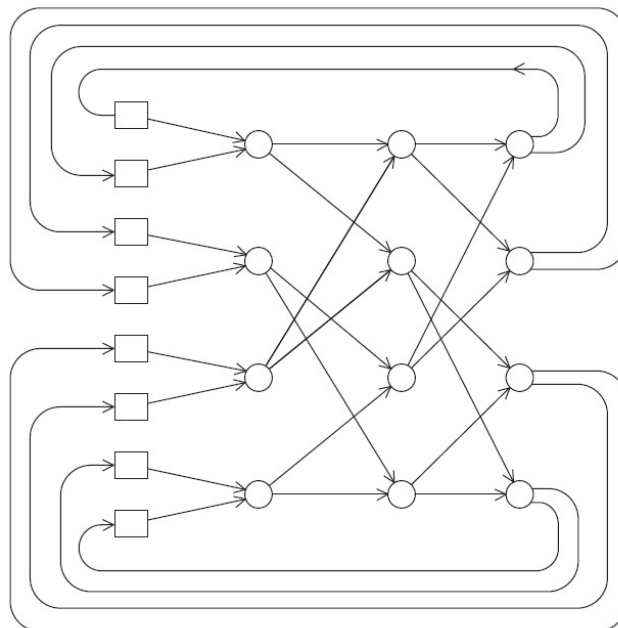


© All Rights Reserved.

Indirect Interconnect (cont.)



Omega



© All Rights Reserved.



Computer Cluster

- A set of computers (nodes) interconnected over a network that functions as a single large multiprocessor
 - Each node has its own private memory and OS
 - Usually connected using a LAN
 - E.g., Ethernet, Infiniband
- Clusters are used for compute-intensive applications (such as simulations) or IO-intensive applications (such as database systems)
- Clusters are, by far, the most common type of distributed-memory platforms

© All Rights Reserved.



Beowulf Cluster

- A simple, low-cost home-built cluster
- Commercial off-the-shelf computers
- Ethernet network (or any other network type)



© All Rights Reserved.

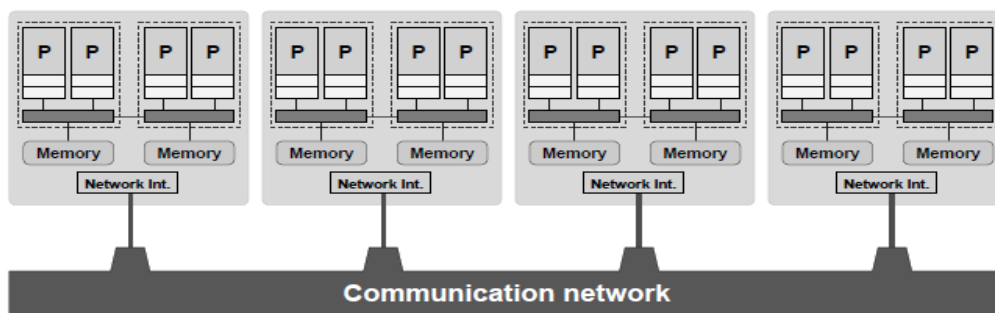
Achieving High-Performance Computing On Clusters



- Programmers write a *message-passing* program, which is a parallel program that runs on a cluster
 - Programmer partitions the work among cluster nodes
 - Programmer inserts explicit send/receive messages to exchange shared data between cluster nodes
- *MPI* is currently the dominant programming model for clusters

© All Rights Reserved.

Hybrid Clusters



- Each node consists of a multicore chip (in this case NUMA)
- The best of both worlds: scalability + faster performance per node
- Drawback: more costly nodes + might complicate programmability in some cases
- Today, most HPC clusters have multicore nodes

© All Rights Reserved.

Clusters Are Ubiquitous



- The high scalability, immunity to failure, cost effectiveness and computational power made clusters the backbone of most (if not all) HPC platforms on the planet
- Nowadays, millions of clusters are scattered around the world to make the current internet era we live in possible
- Let us consider three examples where you can find clusters:
 - Supercomputers
 - Distributed systems
 - Cloud computing

© All Rights Reserved.

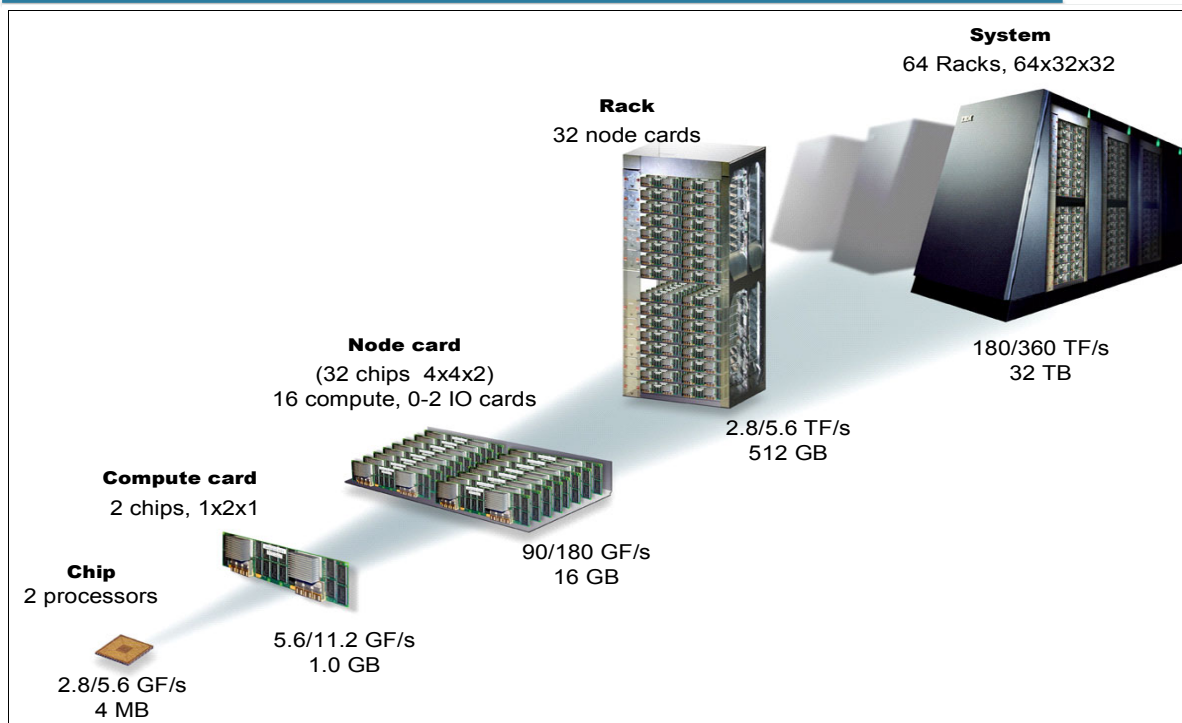
Supercomputers



- A supercomputer is basically a *HPC* cluster with a large number of nodes and a high-performance interconnection network
- A node in a supercomputer typically has a hierarchical SIMD / MIMD architecture with a lot of cores
- Performance is measured in floating-points operations per second
- Supercomputers are used for scientific computing: quantum mechanics, weather forecasting, climate research, oil and gas exploration, molecular modeling, ...etc
- Supercomputer example:
 - BlueGene/L (2007), 106.496 nodes x 2 PowerPC (700MHz)
 - IBM Sequoia (2012), 16,3 PFlops, 1.6 PB memory, 98304 compute nodes, 1.6 Million cores, 7890 kW power

© All Rights Reserved.

BlueGene/L



© All Rights Reserved.

Top500



- Annual ranking of the most powerful 500 supercomputers in the world (<https://www.top500.org/>)
- LINPACK benchmarks are used to measure performance
- Current #1 is *Fugaku*, a supercomputer located in RIKEN Center for Computational Science, Japan
 - A64FX 48C 2.2GHz
 - 7,630,848 cores
 - 5,087,232 GB RAM
 - 442,010 TFlop/s (Linpack performance)
 - <https://www.r-ccs.riken.jp/en/fugaku/project>

© All Rights Reserved.

Distributed Systems



- [Tanenbaum] “A distributed system is a collection of independent computers that appear to the users of the system as a single computer”
- [Coulouris et al.] “system in which hardware or software components located at networked computers communicate and coordinate their actions only by passing messages”
- **My definition:** A distributed system is a collection of tightly-coupled or loosely-coupled clusters that cooperate together to make a particular functionality continuously available to the user

48

Why Distributed Systems?



- Scalability
 - Many applications nowadays (e.g., web services, online games, peer-to-peer applications, social networks) have enormous number of users with dispersed geographical locations
- Parallelism
 - Users can access the resources simultaneously without interfering with each other (i.e., do independent tasks)
- Highly Reliability
 - The failure of one machine will not affect other machines
- Transparency
 - Users have little or no knowledge of where distributed systems' resources are physically located

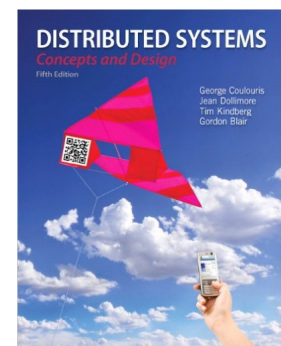
49

Distributed Systems Course



- Many universities dedicate a full course to teaching distributed systems (typically offered to graduate students)
- Discussed topics in this course include:
 - Design approaches for distributed systems
 - Networking and internet protocols
 - Distributed operating system
 - Processes communication in distributed operating systems and remote procedure call protocols
 - Distributed file systems
 - Secure Networks
 - Distributed system applications such as peer-to-peer systems and web services

Suggested textbook



50

Cloud Computing



- A relatively new computing paradigm where users have on demand access to an “infinite” pool of virtual resources
- The main principle behind cloud computing model is to offer *computing, storage, and software* as a service or as a utility *via the internet*
- My personal view: cloud computing is a fancy term for a *virtual distributed system* that offers internet services on infrastructure level, platform level and software level

51

Why Cloud Computing?



- On-demand pay-as-you-go services
 - You pay for what you use
- Ubiquitous internet access
 - You can access the cloud from anywhere as long as you have internet
- Cut capital and operational cost of buying hardware
 - Skip the pain of maintaining on-premises clusters
- Rapid elasticity
 - E.g., go from 10 machines to 100 or from 100 machines to 10
- Location-independent resource pooling
 - You work with virtual machines that could be hosted anywhere

52

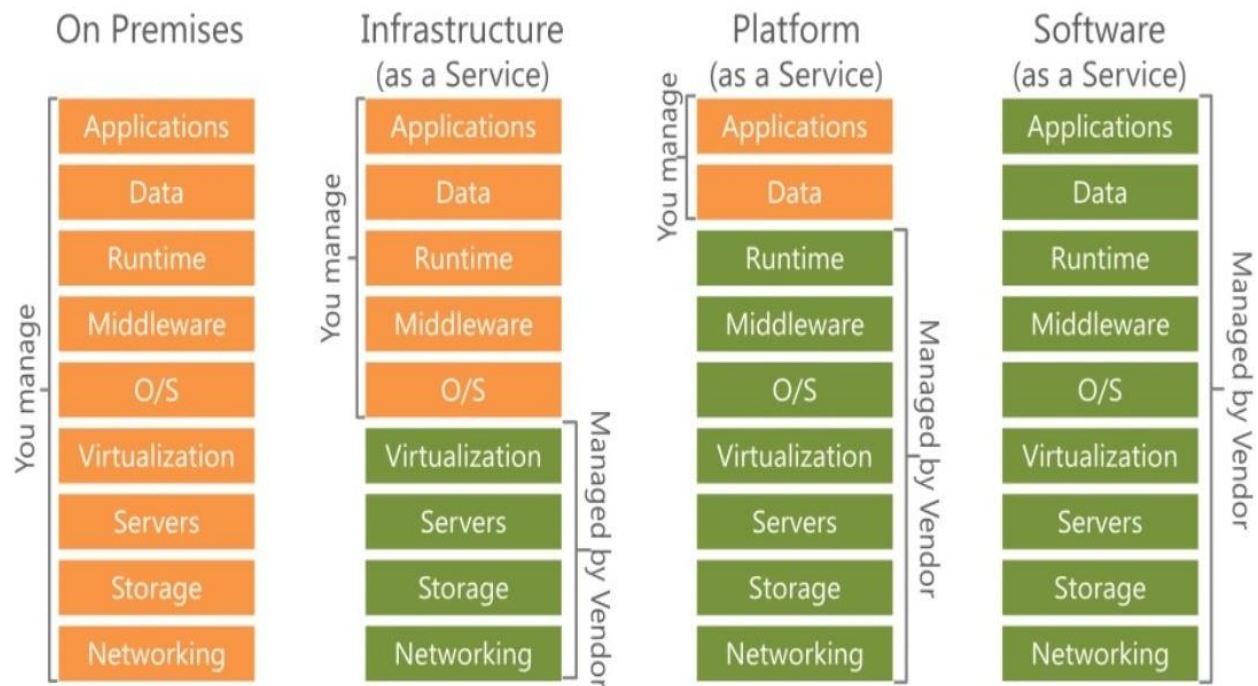
Cloud Computing Services



- *Software as a service (SaaS)*
 - Users access applications that are hosted on the cloud
- *Platform as a service (PaaS)*
 - Users build and deploy their own software applications on the cloud
- *Infrastructure as a service (IaaS)*
 - Users remotely control a collection of “virtual” machines in terms of operating system, storage, network connectivity and applications

53

IaaS, PaaS and SaaS



Cloud Platform Vendors

- Amazon Web Services (AWS)
- Windows Azure
- Google App Engine
- Rackspace



Last Remarks



- Parallel machines can be classified based on their *concurrency* and *communication* models
- Writing efficient parallel programs require exposing the programmer to the architectural details of parallel machines
- There is no standard parallel machine, and therefore there is no standard parallel programming environment
- In this course, we will learn how to write parallel programs for shared-memory and distributed-memory platforms

© All Rights Reserved.

Parallel Software

Instructor: Dr. Fahed Jubair
Computer Engineering Department
University of Jordan

Parallel Computing Quick Recap



- Parallel computing achieved by running a **parallel program** on a **parallel machine**
- Different parallel machines have different architectural organizations and therefore they require different types of parallel programs

© All Rights Reserved.

SPMD



- The programmer only needs to write a single program that is **parameterized by the thread/process number**
- At runtime, each thread/process executes the pieces of code that correspond to their thread/process number
 - E.g., using conditional branches
- SPMD is the most common style for parallel programming with MIMD architectures
 - The programmer writes a single program that gets replicated on all processors

© All Rights Reserved.



SPMD (cont.)

- Assume the number of threads/processes is n
- Below is a high-level picture of SPMD-style coding

// multithreading code

```
if (thread_id == 0)
    do task 0
else if (thread_id == 1)
    do task 1
else if (thread_id == 2)
    do task 2
...
else // thread_id = n - 1
    do task n - 1
```

// message-passing code

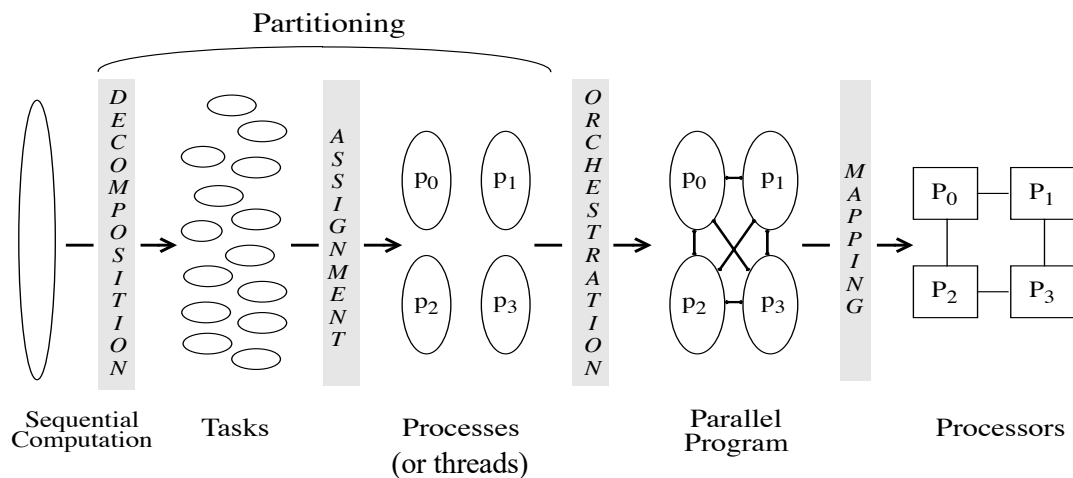
```
if (process_id == 0)
    do task 0
else if (process_id == 1)
    do task 1
else if (process_id == 2)
    do task 2
...
else // process_id = n - 1
    do task n - 1
```

© All Rights Reserved.

General Steps For Writing A Parallel Program



- [Culler et al.] described four general steps to parallelize a given sequential program:



© All Rights Reserved.

1st Step: Decomposition



- Break up computation into tasks with as much concurrency as possible
- Tasks can be created **statically** or **dynamically**
- The challenge is how to identify enough concurrency to keep the processors busy all the time, yet not so much that the overhead of managing the tasks becomes substantial compared to the useful work done
- In some cases, programmers need to perform algorithmic changes to increase parallelism

© All Rights Reserved.

2nd Step: Assignment



- Specify the mechanism by which tasks will be distributed among threads/processes
- This assignment of tasks to threads/processes can be done **statically** or **dynamically**
- Challenges:
 1. How balance the work among processes/threads
 2. How to minimize inter-process (or inter-thread) communication
 3. How to reduce the runtime overheads of managing the assignment

© All Rights Reserved.

Putting The Two Together



- The two steps together (decomposition and assignment) are called *work partitioning*
- Logically, work partitioning is algorithm-dependent and machine-independent
- However, the architectural details of parallel machines often impacts the programmer decisions on how to perform work partitioning

© All Rights Reserved.

3rd Step: Orchestration



- Specify mechanisms of synchronization and communication between threads/processes
- This step is largely dependent on the architecture and the programming model
- Goal:
 - Use efficient synchronization and communication mechanisms to reduce their cost
 - Reduce serialization of shared resources

© All Rights Reserved.

4th Step: Mapping



- Determine which threads/processes run on which processors
- Most programming models have runtime support to perform this step automatically to the programmer
- Mapping can have significant performance impact
 - For example, in NUMA architecture, mapping threads that access overlapped shared data regions to the same socket reduces remote memory accesses
 - Network topology can play a role in reducing communication latency for NUMA and distributed-memory architectures

© All Rights Reserved.

Example – Histogram



- Given an array of float numbers chosen from 0.0 – 5.0, consider a program that builds a histogram by distributing all float numbers into five “bins” where:
 - Bin[0] has the count of all floats in $0 \leq num < 1$
 - Bin[1] has the count of all floats in $1 \leq num < 2$
 - Bin[2] has the count of all floats in $2 \leq num < 3$
 - Bin[3] has the count of all floats in $3 \leq num < 4$
 - Bin[4] has the count of all floats in $4 \leq num \leq 5$

© All Rights Reserved.



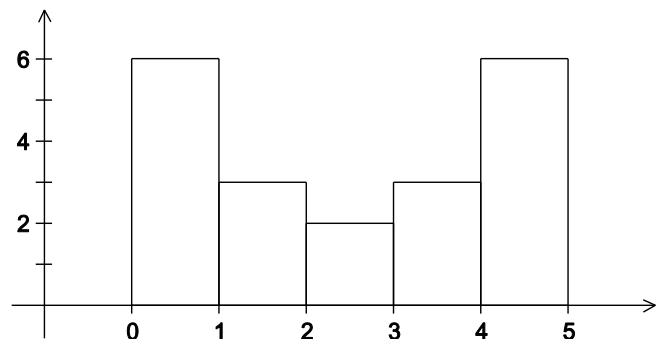
Histogram (cont.)

For example, if we have the following 20 floats:

1.3,2.9,0.4,0.3,1.3,4.4,1.7,0.4,3.2,0.3,4.9,
2.4,3.1,4.4,3.9,0.4,4.2,4.5,4.9,0.9

Then the histogram is determined as follows:

Bin[0] = 6
Bin[1] = 3
Bin[2] = 2
Bin[3] = 3
Bin[4] = 6



© All Rights Reserved.



Histogram – Sequential Code

- Let $D[N]$ be the array of floats, where N is the number of floats
- Let $\text{Bin}[5]$ be the array of bins
- Let $\text{find_bin}(\text{float } d)$ be a function that determines which Bin float d belongs to

```
for ( i = 0; i < N; i++ ) {  
    j = find_bin ( D [ i ] );  
    Bin [ j ]++ ;  
}
```

© All Rights Reserved.



Histogram – Partitioning

- Let p be the total number of threads/processes and id be the thread/process number
- Decomposition
 - Each iteration in the loop represents a single task
 - We have N tasks
- Assignment
 - N tasks and p threads/processes \Rightarrow assign N/p tasks to each process/thread

```
chunk_size = Ceiling ( N / p )
```

```
lb = id * chunk_size
```

```
ub = min ( lb + chunk_size - 1, N )
```

```
for ( i = lb ; i <= ub ; i++ ) {
```

```
    j = find_bin ( D [ i ] );
```

```
    Bin_local [ j ]++ ;
```

```
}
```

© All Rights Reserved.

Is this partitioning really even?
What if $N=124$ and $p=16$?
Can you create a better math for more balanced partitioning?

Histogram – Orchestration with Multithreading



- Each threads computes a local Bin array
- All threads update the global Bin array
- Use **mutual exclusion** to avoid race conditions

```
for ( i = lb ; i <= ub ; i++ )
```

```
    j = find_bin ( D [ i ] );
```

```
    Bin_local [ j ]++ ;
```

```
}
```

```
critical section {
```

```
    for ( i = 0 ; i < 5 ; i++ )
```

```
        Bin [ i ] = Bin [ i ] + Bin_local [ i ] ;
```

```
}
```

© All Rights Reserved.

Histogram – Orchestration with Message–Passing



- Each process computes a local Bin array
- All processes send their local Bin arrays to process 0
- Process 0 combines all local Bin arrays to compute overall Bin array

```
for (i = lb ; i <= ub ; i++)
    j = find_bin ( D [ i ] );
    Bin [ j ]++ ;
}
if ( process_id != 0 )
    send_msg ( receiver = 0 , data = Bin ) ;
else
    for (i = 1 ; i < T ; i++) {
        rcv_msg ( sender = i , & data ) ;
        for (j = 0 ; j < 5 ; j++)
            Bin [ j ] = Bin [ j ] + data [ j ] ;
    }
```

© All Rights Reserved.

Input and Output



- How do you think threads/processes in parallel programs should
 - Read inputs from user?
 - Print errors?
 - Dump text to output files?
- Input/output functionalities can potentially complicate parallel programming
- For simplicity, in this course, we will let only a single thread or process responsible for input/output operations

© All Rights Reserved.



Measuring Performance

- Recall the standard definition of performance
- For a program x , performance is measured as follows:

$$\text{Performance}_x = 1 / \text{Execution time}_x$$

- Program x is n times faster than program y if:

$$\text{Performance}_x / \text{Performance}_y = \text{Execution time}_y / \text{Execution time}_x = n$$

© All Rights Reserved.



Speedup

- Number of the threads/processes = p
- Execution time of the serial version = T_{serial}
- Execution time of the parallel version = T_{parallel}

$$\text{Speedup} = \frac{T_{\text{serial}}}{T_{\text{parallel}}}$$

- It is called “**linear speedup**” when **speedup = p**
- Achieving linear speedups is difficult due to
 - Some parts of the algorithm are purely sequential
 - Runtime overhead of parallelism

© All Rights Reserved.

Execution Time Breakdown



The sequential part of the code that can be parallelized

The sequential part of the code that is "purely" sequential

$$T_{\text{serial}} = T_{\text{parallelizable}} + T_{\text{non-parallelizable}}$$

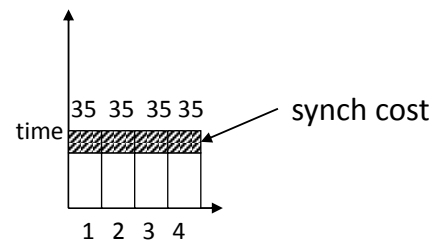
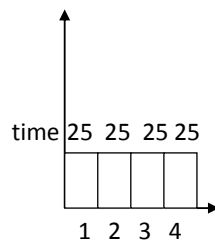
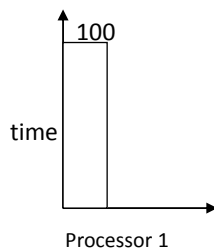
$$T_{\text{parallel}} = \frac{T_{\text{parallelizable}}}{P} + T_{\text{non-parallelizable}} + T_{\text{overhead}}$$

Assuming linear speedup

Parallelism overhead

© All Rights Reserved.

Example



$$S_p = \frac{100}{25} = 4.0,$$

Perfect parallelization!
Does it ever occur?

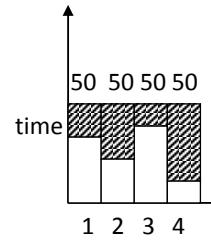
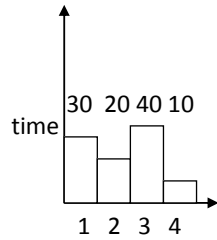
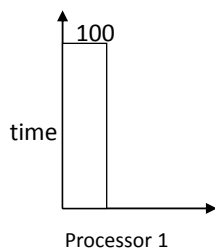
$$S_p = \frac{100}{35} = 2.85,$$

perfect load balancing

© All Rights Reserved.



Example (cont.)



closest to
real life
parallel programs

$$S_p = \frac{100}{40} = 2.5,$$

load imbalance

$$S_p = \frac{100}{50} = 2.0,$$

load imbalance
and sync cost

© All Rights Reserved.



Sources of Overhead

- Sources of runtime overhead in parallel programs:
 - Synchronization overhead (e.g., contention on shared resources)
 - Communication overhead (e.g., overhead of preparing messages)
 - Threads/processes creation and termination
 - Extra computation to perform partitioning
 - Dynamic tasks creation
 - Load imbalance

© All Rights Reserved.



Amdahl's Law

- Let f be the fraction that is “parallelizable” of the program

$$\text{Speedup} = \frac{T_{\text{serial}}}{T_{\text{parallel}}} = \frac{1}{1-f + \frac{f}{p}}$$

$$\text{Speedup} \leq \frac{1}{1-f}$$

- Amdahl's law tell us that, even if $p=\infty$, the serial part of a parallel program limits the performance to an upper bound of $1/(1-f)$
- For example, if $f=95\%$, then the maximum theoretical speedup we can get is 20!

© All Rights Reserved.



No Need To Despair

1. Amdahl's law does not take the problem size into account
 - For many problems, as we increase the problem size, the “inherently serial” fraction of the program decreases in size
2. Thousands of scientific applications have shown huge speedups on large distributed-memory systems
3. Many applications benefit from “small” speedups

© All Rights Reserved.

Scalability



- A parallel program is **scalable** if it can handle ever increasing problem sizes
- Let us define efficiency as follows:

$$\text{Efficiency} = \frac{T_{\text{serial}}}{p * T_{\text{parallel}}}$$

- A parallel program is **strongly scalable** if the efficiency, for a fixed problem size, remains fixed while increasing the number of threads/processes
- A parallel program is **weakly scalable** if the efficiency, for a fixed number of threads/processes, remains fixed while increasing the problem size

© All Rights Reserved.

Benchmarks



- A set of standard programs that measure performance
- Reasonable representation of real-world applications
- Role of benchmarks:
 - Compare different machines
 - Help exploring architectural designs
 - Identify bottlenecks
- Parallel benchmarks are benchmarks that are designed for measuring the performance of parallel machines
- Examples of parallel benchmark suites: NPB, PARSEC, Rodinia and SPLASH-2

© All Rights Reserved.

PARSEC



- Princeton Application Repository for Shared-Memory Computers
- Benchmark Suite for multicores
- Available at <http://parsec.cs.princeton.edu/>
- Objectives:
 - Multithreaded Applications: Future programs must run on multicores
 - Emerging Workloads: Increasing CPU performance enables new applications
 - Diverse: Multicores are being used for more and more tasks
 - State-of-Art Techniques: Algorithms and programming techniques evolve rapidly

© All Rights Reserved.

PARSEC (cont.)



Program	Application Domain	Parallelization
Blackscholes	Financial Analysis	Data-parallel
Bodytrack	Computer Vision	Data-parallel
Canneal	Engineering	Unstructured
Dedup	Enterprise Storage	Pipeline
Facesim	Animation	Data-parallel
Ferret	Similarity Search	Pipeline
Fluidanimate	Animation	Data-parallel
Freqmine	Data Mining	Data-parallel
Streamcluster	Data Mining	Data-parallel
Swaptions	Financial Analysis	Data-parallel
Vips	Media Processing	Data-parallel
X264	Media Processing	Pipeline

© All Rights Reserved.

Concluding Remarks



- In this course, we will focus on writing programs for homogeneous MIMD systems using the SPMD-style
- Parallel program design
 - Decompose computation into tasks
 - Assign tasks to threads/processes
 - Coordinate data sharing between threads/processes
 - Map threads/processes to processors
- Performance evaluation is very important to assess programming quality as well as the underlying architecture and how they interact
- Scalability and efficiency measure the quality of parallel programs

© All Rights Reserved.

Introduction to OpenMP

Instructor: Dr. Fahed Jubair
Computer Engineering Department
University of Jordan



Definition

- [Wikipedia] OpenMP (**Open Multi-Processing**) is an *application programming interface* (API) that supports multi-platform shared memory multiprocessing programming in C, C++, and Fortran
- In other words, OpenMP is a programming model that specifies how to **extend** a sequential programming language in order to support multithreading programming
- Major compiler vendors for OpenMP: PGI, Cray, Intel, Oracle, HP, Fujitsu, Microsoft, AMD, IBM, NEC, ... Etc
- Current standard: version 5.1, Available at <https://www.openmp.org/spec-html/5.1/openmp.html>



OpenMP Components

- **Compiler directives and clauses**
 - Interpreted when OpenMP compiler option is turned on
 - Each directive applies to the succeeding structured block
- **Runtime function calls**
 - Supported functionalities by an OpenMP runtime library
- **Environment variables**
 - User-controlled variables that affect the execution



Fork-Join-Parallelism

- OpenMP uses the **fork-join model** for its parallel execution
 - Only **master thread** executes serial regions
 - Master thread **forks** new threads at the beginning of parallel regions
 - Multiple threads share work in parallel
 - Threads **join** at the end of the parallel regions
- Each thread works on global **shared** and its own **private** variables
- Threads communicate **implicitly** by reading and writing shared variables
- OpenMP uses a relaxed memory model



4



A Starting Example

```
// sequential code
void main ( )
{
    double x[512] ;
    int i ;

    for ( i = 0 ; i < 512 ; i++ ) {
        do_some_work( x[ i ] ) ;
    }
}
```




A Starting Example

// OpenMP code

`#include <omp.h>` → OpenMP runtime library

`void main ()`

{

`double x[512];`

`int i;`

→ This directive means a parallel region is started

→ This clause means to fork eight threads

`#pragma omp parallel num_thread(8)`

{

`#pragma omp for` →

`for (i = 0; i < 512; i++) {`
`do_some_work(x[i]);`

}

}

}

- This directive means the for loop is a *work-sharing* construct
- By default, the iteration space is **evenly** partitioned between the eight threads

© All Rights Reserved.

6



Why OpenMP?

- It is simple!
 - Programmers incrementally add compiler directives to sequential programs
 - Compiler takes care of the ugly details: threads creation and termination, synchronization, data partitioning, ... etc
- Widely applicable
 - Work with both data-parallel and task-parallel applications
 - OpenMP has compiler directives/clauses for almost everything
- The same source code can be used for both sequential and OpenMP codes
 - Simply ignore directives when compiling as a sequential program

© All Rights Reserved.

7



Compiler Directives

- Assuming C/C++, all directives must start with *#pragma omp*
- Directives types:
 - `#pragma omp parallel` → Parallel region directive
 - `#pragma omp for`
 - `#pragma omp section`
 - `#pragma omp single`
 - `#pragma omp critical`
 - `#pragma omp atomic`
 - `#pragma omp barrier`
 - `#pragma omp master`
 - `#pragma omp threadprivate` → THREADPRIVATE directive
- Each directive has a set of clauses that define additional attributes (see next slides)

© All Rights Reserved.

8



Clauses

- Data attribute clauses
 - `shared`
 - `private`, `lastprivate`, `firstprivate`
 - `copyin`, `copyprivate`
- Synchronization clauses
 - `nowait`
- Work-sharing clauses
 - `schedule`
 - `reduction`
 - `collapse`

© All Rights Reserved.

9

Runtime Library Routines



- Number of threads: `omp_{set,get}_num_threads`
- Thread ID: `omp_get_thread_num`
- Wall clock time: `omp_get_wtime`
- Acquire/release lock: `omp_{set,unset}_lock`
- And more ...

Environment Variables



- `OMP_NUM_THREADS`
- `OMP_THREAD_LIMIT`
- `OMP_SCHEDULE`
- `OMP_PROC_BIND`
- And more ...

The Parallel Region Directive



- Programmers use the “`#pragma omp parallel`” directive to mark a block of code as a *parallel region*
- At runtime, when the master thread reaches a parallel region, it creates a team of threads that execute this region in parallel
 - i.e., the code of the parallel region is duplicated and all threads will execute this code simultaneously
- The master thread will also be a part of this team (as thread number 0)
- The clauses “`shared`” and “`private`” are used to determine if each variable is private or shared inside a parallel region
 - Variables declared inside a parallel region are always *private*
 - Variables that are not specified a “`private`” or a “`shared`” clause are, by default, shared

How Many Threads?



- The number of threads in a parallel region can be specified using any of the following:
 - “`num_threads`” clause
 - “`omp_set_num_threads`” library function
 - “`OMP_NUM_THREADS`” variable
 - Implementation default, which is usually the same number as the number of cores in the system



The Traditional Hello World Example

```
#include <omp.h>
#include <stdio.h>

void main () {
    int nthreads, tid;

    /* Fork a team of 8 threads with each thread having a private tid variable */
    #pragma omp parallel num_threads(8) private(tid)
    {
        /* Obtain and print thread id */
        tid = omp_get_thread_num();
        printf("Hello World from thread = %d\n", tid);

        /* Only master thread does this */
        if (tid == 0)
        {
            nthreads = omp_get_num_threads();
            printf("Number of threads = %d\n", nthreads);
        }
    } /* All threads join master thread and terminate */
}
```

© All Rights Reserved.

14



Compiling and Running Hello World Example

- Assuming we have GNU GCC compiler, we can compile the Hello World OpenMP program on a linux machine using the following command

```
gcc -fopenmp hello_world.c -o hello.out
```

- To run the program, we use the command

```
./hello.out
```

- You can set the number of threads using the OMP_NUM_THREADS environment variable, as follows:

bash shell: `export OMP_NUM_THREADS=8`

csh/tsch shell: `setenv OMP_NUM_THREADS "8"`

© All Rights Reserved.

15

Work-Sharing Constructs



- Sub-regions where the work is **partitioned** among the participating threads **inside** a parallel region
- OpenMP has three types of directives for work sharing constructs:
 1. `#pragma omp for`
 - Used for parallel *for* loops
 - Work-sharing happens by partitioning the iteration space of the *for* loop
 - Loop-level parallelism works well with data parallelism
 2. `#pragma omp section`
 - Work-sharing happens by breaking work into separate, discrete sections
 - Each section is executed by a thread
 - OpenMP sections work well with function parallelism
 3. `#pragma omp single`
 - Only one thread (can be any thread in the team) will execute the code

© All Rights Reserved.

16

Parallel *for* Loops



```
#pragma omp parallel  
{
```

```
#pragma omp for private(i) schedule(static)  
for ( i=1; i<=198; i++ )  
    C[i] = (A[i] + B[i])/2.0 ;
```

This means use *static* scheduling when dividing iterations among threads (see next slide)

```
#pragma omp for private(i) schedule(static)  
for( i=0; i<=199; i++ )  
    A[i] = sqrt(C[i]) ;
```

```
}
```

- By default, arrays A, B, and C are shared
- By default, there is an **implicit barrier** at the end of each parallel *for* loop

© All Rights Reserved.

17

The Schedule Clause



- Specify how iterations of a parallel *for* loop is partitioned
- *schedule* (*static*, [*chunk*])
 - Loop iterations are divided into pieces of size *chunk* and then **statically** assigned to threads *in round robin fashion*
 - If *chunk* size is not specified, the iterations are evenly (if possible) divided contiguously among the threads (called *block partitioning*)
 - Special case: if *chunk* size is 1, then the partitioning is called cyclic
- *schedule* (*dynamic*, [*chunk*])
 - Loop iterations are divided into pieces of size *chunk*, and dynamically scheduled among the threads
 - When a thread finishes one *chunk*, it grabs another
 - If *chunk* size is not specified, then by default, it is equal to 1

Static vs Dynamic Scheduling



- Static scheduling has lower overhead but is more prone to load imbalance
- Dynamic scheduling has higher overhead but less prone to load imbalance
- So how do we decide which scheduling scheme is better?

Iteration Scheduling Affects Thread-Data Distribution

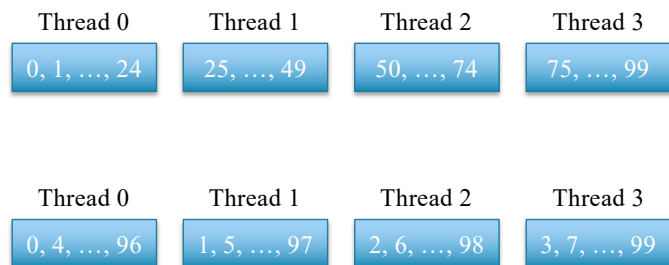


- Assuming the number of threads is 4, let us derive the distribution of *shared array A* on these four threads

```
#pragma omp for schedule(static)
for ( i = 0; i < 100; i++ )
    A[ i ] = ...
```

```
#pragma omp for schedule(static, 1)
for ( i = 0; i < 100; i++ )
    A[ i ] = ...
```

```
#pragma omp for schedule(dynamic)
for ( i = 0; i < 100; i++ )
    A[ i ] = ...
```



Only known at runtime

In general, data-thread distribution for a particular array is a function of iteration scheduling type + the array subscript function

© All Rights Reserved.

20

Exercise



- Assuming the number of threads is 4, describe the thread-data distribution of arrays A and B in the following parallel loops

- 1

```
#pragma omp for schedule(static)
for ( i = 0; i < 100; i++ )
    A[ i ] = B[ i + 1 ];
```
- 2

```
#pragma omp for schedule(static)
for ( i = 0; i < 20; i++ )
    A[ 2 * i ] = B[ 2 * i + 1 ];
```
- 3

```
#pragma omp for schedule(static)
for ( i = 99; i >= 0; i-- )
    A[ i ] = B[ i ];
```
- 4

```
#pragma omp for schedule(static)
for ( i = 0; i < 100; i++ )
    for ( j = 0; j < 100; j++ )
        A[ i ][ j ] = B [ j ][ i ];
```
- 5

```
for ( i = 0; i < 100; i++ ) {
#pragma omp for schedule(static)
for ( j = 0; j < 100; j++ )
    A[ i ][ j ] = B [ j ][ i ];
}
```

© All Rights Reserved.

21



The Nowait Clause

- Remove implicit barrier that is present at the end of parallel *for* loops
- Improve performance because it removes unnecessary synchronization
- Can cause race conditions when used incorrectly

```
#pragma omp parallel
{
#pragma omp for private(i) schedule(static) nowait
  for ( i=0; i<200; i++ )
    C[i] = (A[i] + B[i])/2.0 ;

#pragma omp for private(i) schedule(static)
  for( i=0; i<200; i++ )
    A[i] = sqrt(C[i]) ;
}
```

We can eliminate implicit barrier here because the same threads in both parallel loops produce and consume the same data (i.e., there is no cross-thread dependencies)



Which Nowait Can Cause Race Conditions?

```
#pragma omp parallel
{
#pragma omp for private(i) schedule(static, 1) nowait
  for ( i=0; i<200; i++ )
    C[i] = (A[i] + B[i])/2.0 ;

#pragma omp for private(i) schedule(static, 1)
  for( i=0; i<200; i++ )
    A[i] = sqrt(C[i]) ;
}
```



No race conditions

```
#pragma omp parallel
{
#pragma omp for private(i) schedule(static) nowait
  for ( i=0; i<200; i++ )
    C[i] = (A[i] + B[i])/2.0 ;

#pragma omp for private(i) schedule(static)
  for( i=0; i<199; i++ )
    A[i] = sqrt(C[i]) ;
}
```



A race condition might occur

```
#pragma omp parallel
{
#pragma omp for private(i) schedule(dynamic) nowait
  for ( i=0; i<200; i++ )
    C[i] = (A[i] + B[i])/2.0 ;

#pragma omp for private(i) schedule(dynamic)
  for( i=0; i<200; i++ )
    A[i] = sqrt(C[i]) ;
}
```



A race condition might occur



The Reduction Clause

- The REDUCTION clause performs a reduction on the variables that appear in its list
 - A private copy for each list variable is created for each thread
 - At the end of the reduction, the reduction variable is applied to all private copies of the shared variable
 - The final result is written to the global shared variable
- Reduction operations include $+$, $-$, $*$, $/$, $\&$, $|$, $^$
- For example, the below parallel loop computes the summation in parallel and puts the final result in the shared variable “*sum*”

```
#pragma omp for private(i) reduction(+:sum)
for(i=0;i<SIZE;i++)
    sum = sum + A[i];
```



The Reduction Clause (cont.)

- Reduction clause only supports scalar reduction
 - Variables in the reduction list cannot be arrays or structure data types
- Reduction clause works for both integer and real numbers
 - However, reduction operations may not be associative for real numbers
- Naturally, reduction clauses are only legal for shared variables
- A reduction variable must be used only in the statement where the reduction is applied

The SINGLE Directive



- The SINGLE directive specifies that the enclosed code is to be executed by only one thread in the team
- Any thread in the team can execute the code enclosed by the SINGLE directive
- By default, there is an implicit barrier at the end of the enclosed code by the SINGLE directive
 - The programmer may remove this barrier using the NOWAIT clause
- Single directives are useful when initializing global variables or dealing with I/O operations

Synchronization Directives



- OpenMP provides a set of directives to perform synchronization operations
- Examples:
 - `#pragma omp master`
 - Specifies a region that is to be executed only by the master thread
 - All other threads on the team skip this section of code
 - There is no implied barrier associated with this directive
 - `#pragma omp critical`
 - Specifies a region of code that must be executed by only one thread at a time
 - There is no implied barrier associated with this directive
 - `#pragma omp atomic`
 - Similar to CRITICAL directive but can only applies to one-statement regions
 - `#pragma omp barrier`
 - All threads must wait till all other threads reach this barrier
 - All threads then resume executing in parallel the code that follows the barrier



Example: Parallel Sum in OpenMP

```
#include <omp.h>
#include <stdio.h>
#define SIZE 1000

void main( ){

    int i, nthreads;
    int sum;
    int A[SIZE];
    double exec_time;

    for(i=0; i<SIZE; i++)
        A[i] = i%100;

    exec_time = omp_get_wtime();

    #pragma omp parallel
    {

        #pragma omp single ← Exercise: Can a nowait clause be added here?
            sum = 0;

        #pragma omp for private(i) reduction(+:sum)
            for(i=0; i<SIZE; i++)
                sum = sum + A[i];

        #pragma omp master
            printf("number of threads=%d\n", omp_get_num_threads());

    }

    exec_time = omp_get_wtime() - exec_time;
    printf("sum=%d\n", sum);
    printf("execution time=%f\n", exec_time);

}
```

© All Rights Reserved.

28



Example: Dot Product in OpenMP

```
#include <omp.h>
#include <stdio.h>
#include <stdlib.h>
#define N 256

void main( ){

    int i;
    float x[N], y[N];
    float dot_product;
    double exec_time;

    for(i=0; i<N; i++){
        x[i] = ((float)rand())/((float)RAND_MAX);
        y[i] = ((float)rand())/((float)RAND_MAX);
    }

    exec_time = omp_get_wtime();

    #pragma omp parallel shared(dot_product) private(i)
    {

        #pragma omp single
            dot_product = 0.0;

        #pragma omp for reduction(+:dot_product)
            for(i=0; i<N; i++)
                dot_product = dot_product + x[i] * y[i];

        #pragma omp master
            printf("number of threads=%d\n", omp_get_num_threads());

    }

    exec_time = omp_get_wtime() - exec_time;
    printf("dot product is %f\n", dot_product);
    printf("execution time=%f\n", exec_time);

}
```

© All Rights Reserved.

29



Example: Histogram in OpenMP

```
#include <omp.h>
#include <stdio.h>
#include <stdlib.h>
#define SIZE 20
int find_bin( int d);

void main( ){

    int i,j;
    double A[SIZE];
    int bin[5], my_bin[5];
    double exec_time;

    for(i=0; i<SIZE; i++)
        A[i] = ((float)rand()/((float)RAND_MAX)) * 5;

    exec_time = omp_get_wtime();

#pragma omp parallel private(my_bin, i, j)
    {

#pragma omp single
        for(i=0;i<5;i++)
            bin[i]=0;

        for(i=0;i<5;i++)
            my_bin[i]=0;

#pragma omp for
        for(i=0;i<SIZE;i++){
            j = find_bin(A[i]);
            my_bin[j]++;
        }

#pragma omp critical
        for(i=0;i<5;i++)
            bin[i] = bin[i] + my_bin[i];

#pragma omp master
        printf("number of threads=%d\n",omp_get_num_threads());
    }

    exec_time = omp_get_wtime() - exec_time;
    for(i=0;i<5;i++)
        printf("bin[%d] = %d\n", i, bin[i]);
    printf("execution time=%f\n",exec_time);
}
```

Exercise: Can we use atomic directive instead of critical directive?



Example: Matrix-Vector Product in OpenMP

```
#include <omp.h>
#include <stdio.h>
#include <stdlib.h>
#define N 128
#define M 256

void main( ){

    int i,j;
    float x[N][M], y[M];
    float result[N];
    double exec_time;

    for(i=0; i<N; i++)
        for(j=0; j<M; j++)
            x[i][j] = ((float)rand()/((float)RAND_MAX));

    for(i=0; i<M; i++)
        y[i] = ((float)rand()/((float)RAND_MAX));

    exec_time = omp_get_wtime();

#pragma omp parallel private(i,j)
    {

#pragma omp single
        for(i=0; i<N; i++)
            result[i] = 0.0;

#pragma omp for
        for(i=0;i<N;i++)
            for(j=0;j<M;j++)
                result[i] = result[i] + x[i][j] * y[j];

#pragma omp master
        printf("number of threads=%d\n",omp_get_num_threads());
    }

    exec_time = omp_get_wtime() - exec_time;
    printf("execution time=%f\n",exec_time);
}
```



A Nice Shortcut

- The special directive “`#pragma omp parallel loop`” can be used to start a parallel region that consists only of a parallel loop

```
#pragma omp parallel
{
    #pragma omp for private(i) schedule(static,4)
    for ( i=0; i<200; i++ )
        C[i] = (A[i] + B[i])/2.0 ;
}
```



Can be rewritten

```
#pragma omp parallel for private(i) schedule(static,4)
for ( i=0; i<200; i++ )
    C[i] = (A[i] + B[i])/2.0 ;
```



The SECTIONS Directive

- A non-iterative work-sharing construct
- It specifies that the enclosed section(s) of code are to be divided among the threads in the team
- Each SECTION is executed once by a thread in the team
- Different sections may be executed by different threads
- It is possible for a thread to execute more than one section if it is quick enough and the implementation permits such
- There is an implied barrier at the end of a SECTIONS directive, unless the *nowait* clause is used



The SECTION Directive (cont.)

- Consider the following example

```
#pragma omp parallel shared(A, B, C, D) private(i)
{
    #pragma omp sections
    {
        #pragma omp section
        for(i=0; i<N; i++)
            c[i] = a[i] + b[i];

        #pragma omp section
        for(i=0; i<N; i++)
            d[i] = a[i] - b[i];
    } /* end of sections */
} /* end of parallel region */
```

© All Rights Reserved.

34



Summary

- OpenMP is a directive-based, easy-to-use programming model for shared-memory architectures (including both UMA and NUMA)
- OpenMP relies on both a compiler (to interpret directives and clauses) and a runtime library (to invoke special routines)
- OpenMP uses the fork-join model of parallel execution
- OpenMP is implemented in Fortran, C and C++
- Using work-sharing constructs, OpenMP can exploit both types: data-level and task-level parallelism in parallel regions

© All Rights Reserved.

35

Exercises



1. Write a parallel OpenMP program that counts how many times a negative number is found in a list of 512 randomly generated floating numbers picked from the range $[-1:+1]$
2. Write a parallel OpenMP program that computes $C = A \times B$, where A, B, C are 2D matrices of size 100×100 . You may initialize arrays A and B to random floating-point numbers in the range $[1:100]$
3. Write a parallel OpenMP program that computes the XOR between two n -bit binary numbers A and B , where n, A , and B are read from an input text file.

Loop-Level Parallelism in OpenMP

Instructor: Dr. Fahed Jubair
Computer Engineering Department
University of Jordan

Loop-Level Parallelism



- Tremendous computational algorithms nowadays spend most of their computation in loops
- Therefore, when parallelizing a sequential algorithm, programmers often find themselves need to parallelize loops
- Fundamental questions:
 - How to determine if a loop is serial or parallel?
 - How to convert a serial loop into a parallel loop **without affecting correctness?**
- In this lecture, we will learn how to utilizes loop-level parallelism using OpenMP when parallelizing a sequential program

© All Rights Reserved

Motivating Question



- Can the loops on the right be run in parallel?

- i.e., can different threads run different iterations in parallel?

```
for (i = 1; i < N; i++) {  
    a[i] = b[i];  
    c[i] = a[i - 1];  
}
```

- What needs to be true for a loop to be parallelizable?

- Iterations cannot interfere with each other
- i.e., there is no **dependences** between iterations

```
for (i = 1; i < N; i++) {  
    a[i] = b[i];  
    c[i] = a[i] + b[i - 1];  
}
```

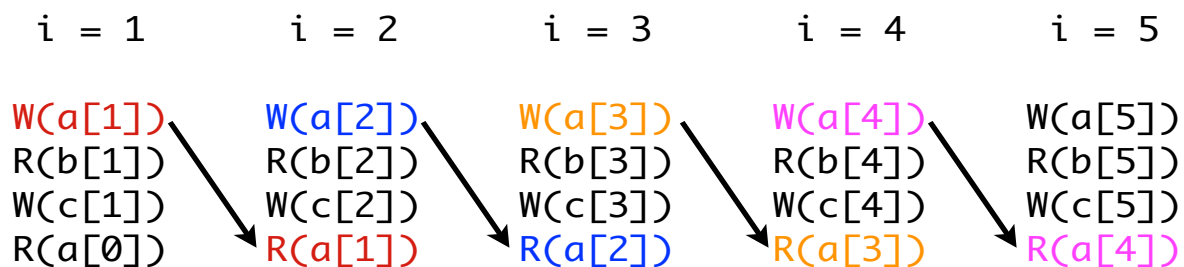
© All Rights Reserved



Flow Dependences

- Definition: A **flow dependence** occurs when one iteration writes a memory location that a later iteration reads

```
for (i = 1; i < N; i++) {
    a[i] = b[i];
    c[i] = a[i - 1];
}
```



© All Rights Reserved



Other Kind of Dependences

- *Anti dependence* – When an iteration reads a location that a later iteration writes (why is this a problem?)

```
for (i = 1; i < N; i++) {
    a[i - 1] = b[i];
    c[i] = a[i];
}
```

- *Output dependence* – When an iteration writes a location that a later iteration writes (why is this a problem?)

```
for (i = 1; i < N; i++) {
    a[i] = b[i];
    a[i + 1] = c[i];
}
```

© All Rights Reserved



Loop-Carried Dependence

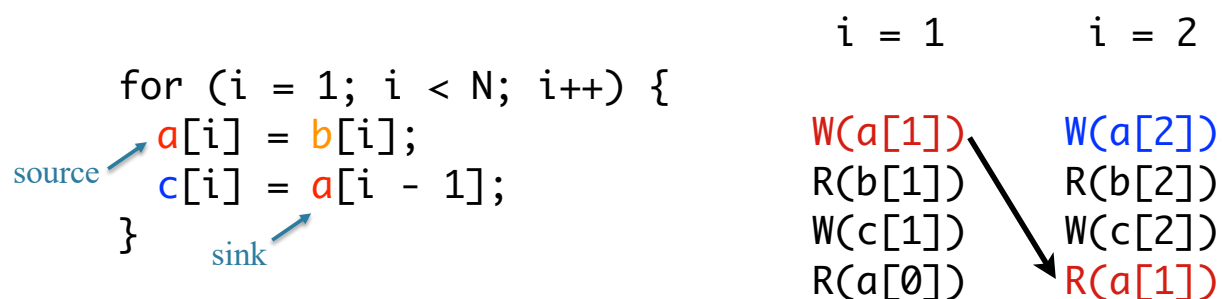
- We can run a loop in parallel if and only if there is no *loop-carried dependences* inside the loop
- **Definition:** A loop has a *loop-carried dependence* if a dependency (whether it is flow, anti- or output dependence) across two iterations is present
- **Definition:** A loop is *parallel* if it has no loop-carried dependence (otherwise, the loop is *serial*)
- To determine if a loop is parallel or serial, we will learn:
 - How to find and represent loop-carried dependences in loops
 - How to use this representation to determine if a loop is parallel or serial

© All Rights Reserved



Terminology

- Dependences can only go forward in time: always from an earlier iteration to a later iteration.
- Dependence *source* is the earlier statement (the statement at the tail of the dependence arrow)
- Dependence *sink* is the later statement (the statement at the head of the dependence arrow)

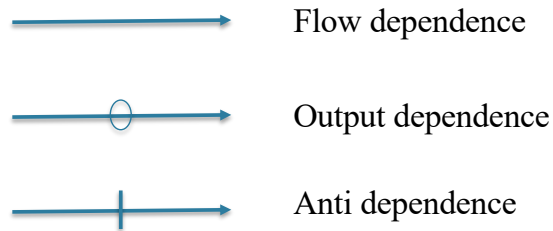


© All Rights Reserved

Representing Dependences



- We use **Loop-carried Dependence Graph** to abstractly represent dependencies inside **loops**
 - Represent iterations as nodes in the graph
 - Draw arrows from sources to sinks to represent loop-carried dependencies

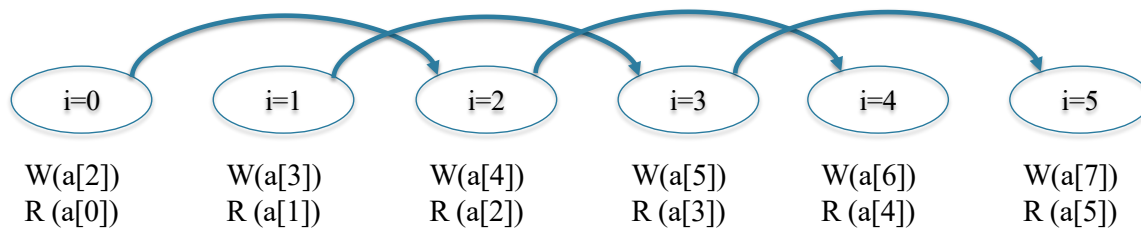


© All Rights Reserved

1D Loop-Carried Dependence Graph Example 1



```
for (i = 0; i < N; i++) {  
    a[i + 2] = a[i]  
}
```



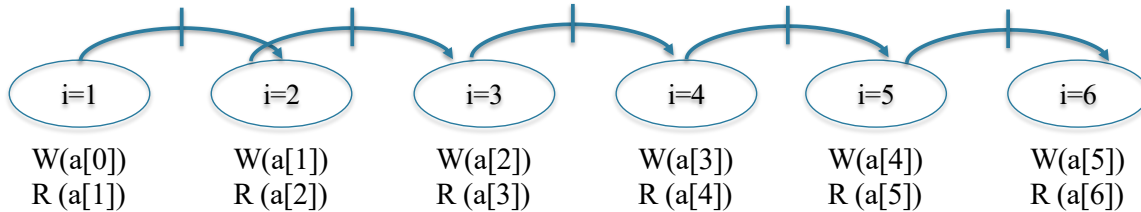
© All Rights Reserved

1D Loop-Carried Dependence Graph Example 2



```

for (i = 1; i < N; i++) {
    a[i - 1] = b[i];
    c[i] = a[i];
}
    
```



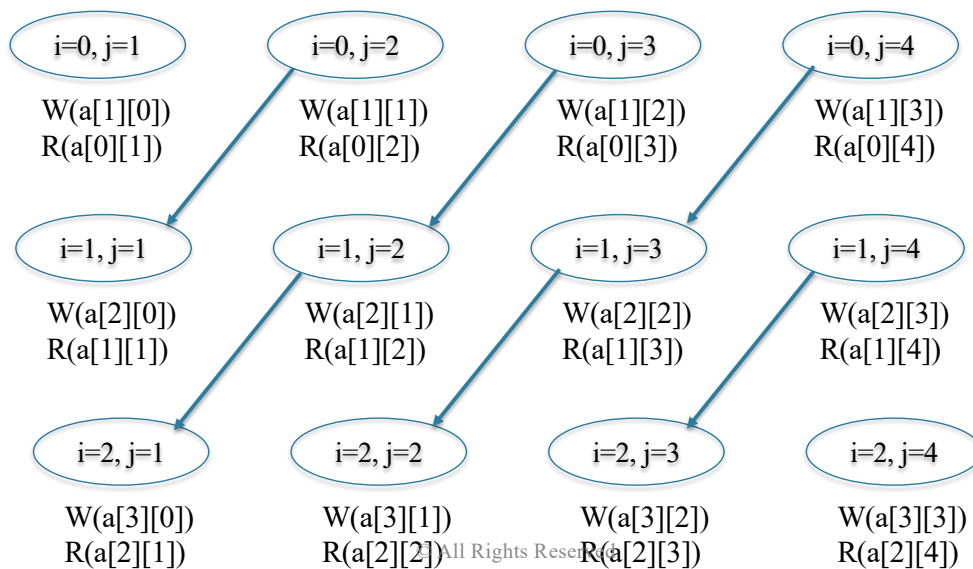
© All Rights Reserved

2D Loop-Carried Dependence Graph Example 1



```

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

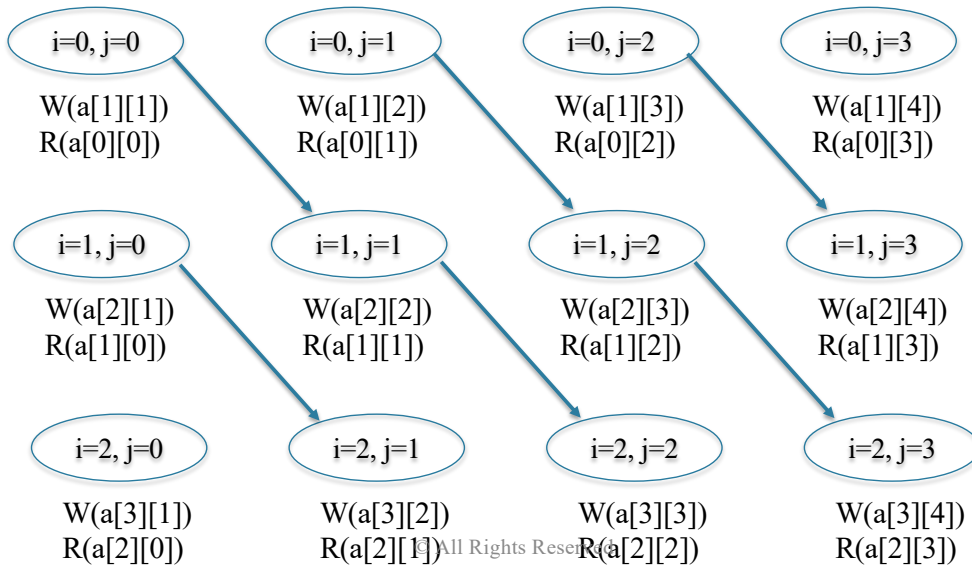


© All Rights Reserved

2D Loop-Carried Dependence Graph Example 2



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



Distance and Direction Vectors

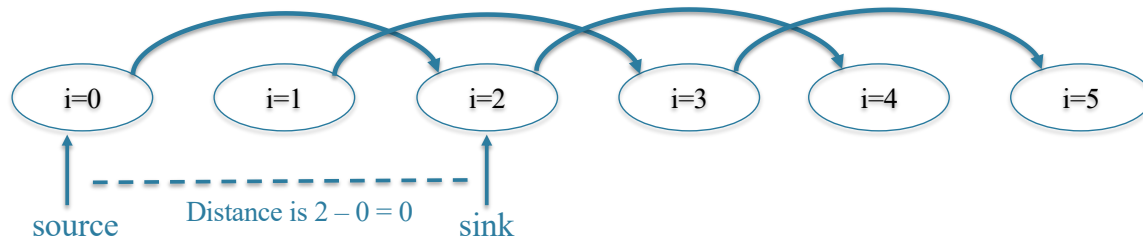


- We can use a more compact form to capture the dependence information in loop-carried dependence graphs using the following two vectors:
 - **Distance vector:** captures the “shape” of the dependence by representing the distance between the *source* and the *sink*
 - **Direction vector:** captures the “direction” of the dependence
- Each distance in the distance vector is computed by the subtraction {sink iteration – source iteration}
 - Therefore, the distance can be a positive or a negative integer (or zero)
- Direction vector can only take three values (=, <, >) which are determined by looking at the sign of the distance vector
 - Use “=” if zero, “<” if positive distance, “>” if negative distance



1D Distance and Direction Vectors

```
for (i = 0; i < N; i++) {  
    a[i + 2] = a[i]  
}
```



- Distance vector is (+2)
- Direction vector is (<)

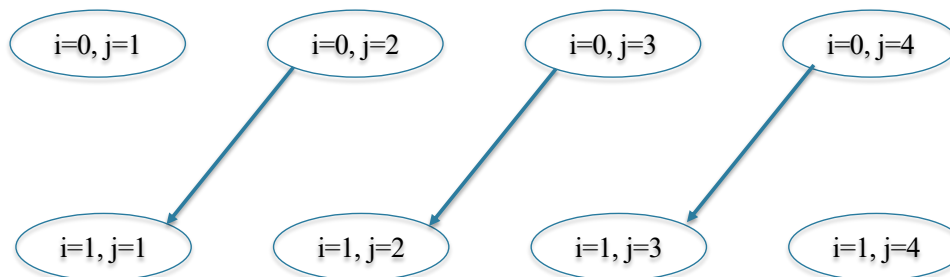
Note that "<" is used to indicate that $i_1 < i_2$, where i_1 is the source iteration and i_2 is the sink iteration

© All Rights Reserved



2D Distance and Direction Vectors

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



- Distance vector is (+1 , -1)
- Direction vector is (< , >)

© All Rights Reserved



Important Property

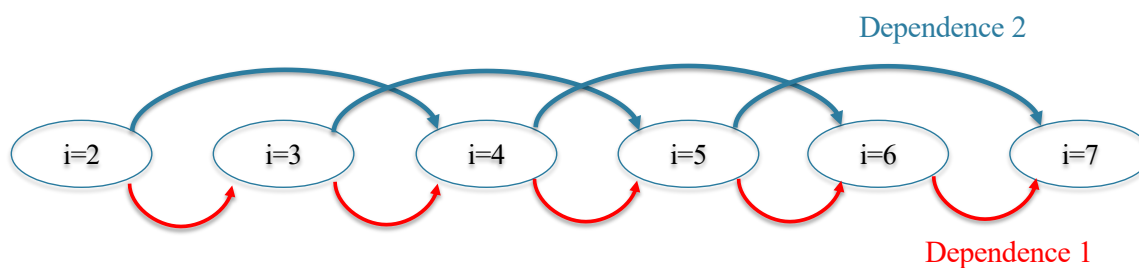
- First entry in any distance vector will always be positive (or zero)
 - Which also means that first entry in any direction vector can be either “=” or “<”
 - If your distance vector has a negative first entry, then you made a mistake in computing the vector
 - Question: why do you think this property always hold?

© All Rights Reserved



Multiple Dependencies Represented Separately

```
fibonacci[0] = fibonacci[1] = 1;  
for (i = 2; i < n; i++)  
    fibonacci[i] = fibonacci[i-1] + fibonacci[i-2];
```



- Distance vector1 is (+1), Distance vector2 is (+2)
- Direction vector1 is (<), Direction vector2 is (<)

© All Rights Reserved

Finding Loop-Carried Dependences



- Loop-Carried dependence graphs provide convenient visualization to all loop-carried dependence in graphs
- Help analyze and determine which loops are parallel.

© All Rights Reserved

Example 1

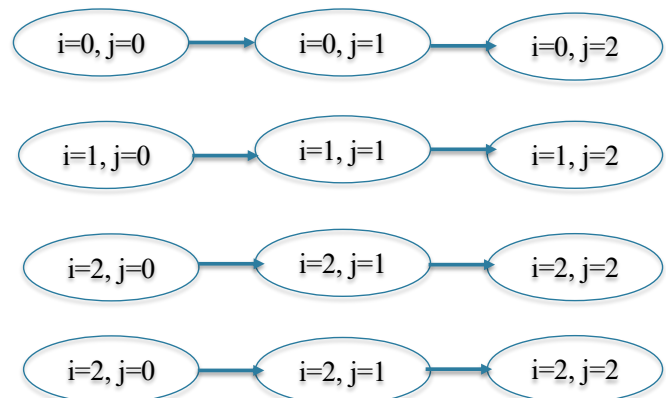


```
for (i=0; i<N; i++)  
  for(j=0; j<N; j++)  
    a [ i ][j + 1] = a [ i ][j ] + 1
```

- Distance vector is (0 , +1)
- Direction vector is (= , <)
- Can parallelize *i* loop, but not *j* loop



```
#pragma omp parallel for private (i,j)  
for (i=0; i<N; i++)  
  for(j=0; j<N; j++)  
    a [ i ][j + 1] = a [ i ][j ] + 1
```



© All Rights Reserved



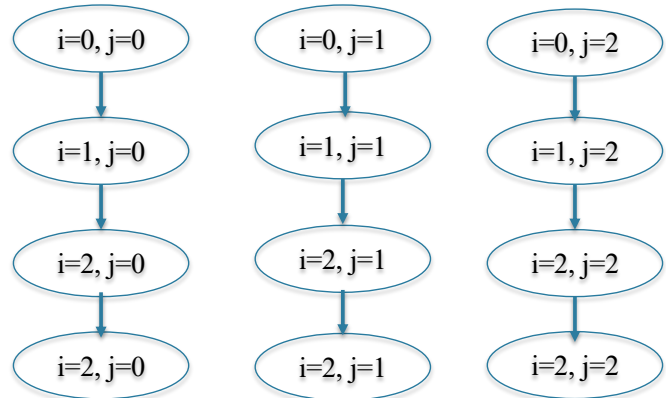
Example 2

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

- Distance vector is $(+1, 0)$
- Direction vector is $(<, =)$
- Can parallelize j loop, but not i loop



```
for (i=0; i<N; i++)
  #pragma omp parallel for private (j)
  for(j=0; j<N; j++)
    a [ i + 1 ][ j ] = a [ i ][ j ] + 1
```



© All Rights Reserved



Exercise

- Draw the iteration space graph for each of the following loops and determine distance and direction vectors
- Determine if they are serial or parallel
- If serial, determine the type of loop-carried dependencies

1

```
for ( i=1; i < n; i++) {
    x = x + 1;
    b [ i ] = a [ i ] * x;
}
```

2

```
for ( i=1; i < n; i++)
  for ( j=1; j < n; j++)
    a [ i ][ j ] = a [ i + 1 ][ j ];
```

3

```
for ( i=1; i < n; i++)
  for ( j=1; j < n; j++)
    a [ i ][ j + 1 ] = a [ i ][ j - 1 ]
                      * a [ i ][ j ];
```

© All Rights Reserved

Let us Make Things More Interesting



- Determine if the following loop is parallel or serial:

```
for ( i=0; i <= n; i++)
```

```
    A [ 2 * i ] = A [ 3 * i + 1 ] ;
```

- More complex array subscripts make using the loop-carried dependence graph method to find dependencies more difficult!

© All Rights Reserved

Problem Formulation



```
for ( i=0; i <= n; i++)
```

```
    A [ f(i) ] = ...
```

```
    ... = A [ g(i) ]
```

- A loop-carried dependence exists if and only if there are two iterations i_1 and i_2 such that:

$$f(i_1) = g(i_2)$$

$$0 \leq i_1, i_2 \leq n$$

- If $i_1 < i_2$, then it is a flow dependence
- If $i_2 < i_1$, then it is an anti dependence

© All Rights Reserved

The Data Dependence Test



- The previous slide show that we can find loop-carried dependencies by solving the below system of equations and inequalities
- Hence, to show that the below loop is parallel, we need to show that this system of equations and inequalities **has no integer solutions**
 - i.e., there is no i_1 and i_2 within $[0:n]$ such that $f(i_1) = g(i_2)$

for ($i=0; i \leq n; i++$) Data dependence test $f(i_1) = g(i_2)$
 $A[f(i)] = \dots$ \longrightarrow $0 \leq i_1, i_2 \leq n$
 $\dots = A[g(i)]$

© All Rights Reserved

Example



for ($i=0; i \leq n; i++$)
 $A[2 * i] = A[3 * i + 1];$

- The loop is parallel if there is no solution for the following system of equations and inequalities:

$$2 * i_1 = 3 * i_2 + 1$$
$$0 \leq i_1, i_2 \leq n$$

© All Rights Reserved



Loop Normalization

- Loop normalization allows us to generalize the aforementioned data dependence test for any loop

```
for (i = L; i < U; i += S)  
... a[i] ...
```



```
for (i = 0; i < (U - L)/S; i += 1)  
... a[S*i + L] ...
```

© All Rights Reserved



Multi-Dimensional Arrays

```
for (i=0; i <= n; i++)  
  for (j=0; j <= m; j++)  
    A[f1(i)][f2(j)] = ...  
    ... = A[g1(i)][g2(j)]
```

- The loop is parallel if there is no solution for the following system of equations and inequalities:

$$f_1(i_1) = g_1(i_2)$$

$$f_2(j_1) = g_2(j_2)$$

$$0 \leq i_1, i_2 \leq n$$

$$0 \leq j_1, j_2 \leq m$$

© All Rights Reserved

Exact Solution is NP-Complete



- It has been shown that providing an exact solution, in general, to the data dependence test is **NP-complete**
- Therefore, we rely on approximate solutions that are accurate for the common case but may have *false positives*
- If a test concludes that a dependence exists, while in fact it does not, then we call this result is a false positive
 - It does not affect correctness to claim a loop is serial while it is actually parallel
- However, a data dependence test should never produce a *false negative*
 - Claiming a loop is parallel while it is in fact serial is incorrect

© All Rights Reserved

The GCD Test



- One simple way to perform a data dependence test is to rewrite the it as *Diophantine equations* and perform the *GCD test*
- Our question: $f(i) = a * i + b, g(i) = c * i + d$
Does $f(i1) = g(i2)$?
- Rewrite $f(i1) = g(i2) \Rightarrow a * i1 + b = c * i2 + d \Rightarrow$
 $a * i1 - c * i2 = d - b \Rightarrow a_1 * i1 + a_2 * i2 = a_3$
- The GCD test:
A *Diophantine equation* $a_1 * i1 + a_2 * i2 = a_3$ has a solution if and only if $\text{gcd}(a_1, a_2)$ evenly divides a_3

© All Rights Reserved

Generalizing The GCD Test



- The Diophantine equation

$a_1 * i_1 + a_2 * i_2 + \dots + a_n * i_n = c$ has a solution iff $\text{gcd}(a_1, a_2, \dots, a_n)$ evenly divides c

Examples:

$15*i + 6*j - 9*k = 12$ has a solution $\text{gcd}=3$

$2*i + 7*j = 3$ has a solution $\text{gcd}=1$

$9*i + 3*j + 6*k = 5$ has no solution $\text{gcd}=3$

© All Rights Reserved

Exercise



- Use the GCD test to decide whether each loop in the following is parallel or serial:

1 for (i=0; i < 20; i++)
 A [2 * i] = A [4 * i + 1] ;

2 for (i=0; i < 20; i++)
 A [2 * i + 3] = A [2 * i] ;

3 for (i=0; i < 20; i++)
 A [3 * i] = A [7 * i] ;

© All Rights Reserved

Other Data Dependence Tests



1. **GCD test**: simple but often inaccurate because it does not account for loop bounds information
2. **Banerjee test** (Utpal Banerjee): more accurate test that takes loop bounds into consideration
3. **Omega test** (William Pugh): most accurate test for linear subscripts (uses the “fancy” idea of integer programming)
4. **Range test** (Blume and Eigenmann): unlike the above tests, this test works with non-linear subscripts

© All Rights Reserved

Quick Recap



- As stated at the beginning of this lecture, there are two fundamental issues programmers deal with when parallelizing a serial program:
 - How to determine if a loop is serial or parallel? (**we already covered this question**)
 - How to convert a serial loop into a parallel loop without affecting correctness? (**our next topic**)

© All Rights Reserved

Eliminating Loop-Carried Dependencies



- The best way is to rewrite the algorithm!
 - There are multiple ways to achieve the same outcome
 - Programmers try to rewrite the computation code in the loop so that the outcome is the same but without loop-carried dependencies

© All Rights Reserved

Parallelism Enabling Techniques



- Many techniques have been proposed in the literature that help programmers to perform appropriate algorithmic changes to eliminate loop-carried dependencies
- We cannot cover all techniques, but we will study the following five popular techniques:
 1. Scalar Privatization
 2. Array Privatization
 3. Scalar Reduction
 4. Array Reduction
 5. Induction Variable Substitution
 6. Loop Fission

© All Rights Reserved



Complications

- Consider the **correctness** of applying privatization to the below loop

```
#pragma omp parallel for private(x)
for (i = 1; i < n; i ++){
    x = f(i);
    A[x] = A[x] * 3;
}
```

```
y = x * x;
```

→ This x is equal to the x computed by the last iteration

- Privatizing scalar x will **cause an error** because the last value of x will not be seen by the code following the parallel loop
- To solve this problem, OpenMP provides **LASTPRIVATE** clause

© All Rights Reserved



LASTPRIVATE Clause

- Behaves exactly like PRIVATE clause, except that the copy computed by the last iteration is copied to the original variable upon exiting the loop

```
#pragma omp parallel for lastprivate(x)
for (i = 1; i < n; i ++){
    x = f(i);
    A[x] = A[x] * 3;
}
```

```
y = x * x;
```

Exercise: there is also **FIRSTPRIVATE** clause by OpenMP. Google it and find out what it does.

© All Rights Reserved



Privatization Exercises

- Parallelize the following loops using proper OpenMP pragmas

1

```
for ( i=0; i < n; i++) {  
    x = i + 2 ;  
    y = x * 2 ;  
    A [ y ] = A [ y ] / 2 ;  
}
```

2

```
for ( i=0; i < n; i++)  
    for ( j=0; j < n; j++) {  
        x = A [ i ][ j ] + y ;  
        B [ i ][ j ] = y + x ;  
    }
```

© All Rights Reserved



Array Privatization

- The idea of privatization can be extended to array storages
- The merit: Enables parallelization
- The problem: Increases memory footprint + Not always feasible

```
for ( i=1; i < n; i++) {  
    for ( j=0; j < 5; j++) {  
        a [ j ] = b [ i ][ j ] * 2 ;  
    }  
    for ( j=0; j < 5; j++) {  
        c [ i ][ j ] = a [ j ] * 2 ;  
    }  
}
```



```
#pragma omp for private(a, i, j)  
for ( i=1; i < n; i++) {  
    for ( j=0; j < 5; j++) {  
        a [ j ] = b [ i ][ j ] * 2 ;  
    }  
    for ( j=0; j < 5; j++) {  
        c [ i ][ j ] = a [ j ] * 2 ;  
    }  
}
```

Exercise: write a proper definition for a private array

© All Rights Reserved



Scalar Reduction

- A reduction operation inside a loop refers to a variable being iteratively updated by the iterations of the loop
- For example, the loop below has a *scalar reduction* operation because the scalar “*sum*” is being iteratively updated by the reduction statement “*sum = sum + A[i]*”
- The bad news: reduction operations cause loop-carried dependencies
- The good news: OpenMP provides the **REDUCTION** clause to enable parallelism (we already studied this clause in lecture 4)

```
for ( i = 1; i < n; i ++ ) {  
    sum = sum + A[ i ] ;  
}
```

© All Rights Reserved



Reduction Eligibility

- A variable is eligible for reduction inside a loop if and only if it is read and written by the statement where the reduction is performed

```
for ( i=1; i < n; i ++ ) {  
    A[ i ] = A[ i ] * A[ i ] ;  
    x = x * A[ i ] ;  
}
```

x is eligible for reduction
because it is read and written
by a single statement (we call
it the reduction statement)

```
for ( i=1; i < n; i ++ ) {  
    y = y / A[ i ] ;  
    B[ i ] = y ;  
}
```

y is NOT eligible for
reduction because it is being
used by two statements

© All Rights Reserved



Array Reduction

- Similar to scalars, reduction can also target arrays
- OpenMP does **not** allow using arrays with the REDUCTION clause
- Therefore, programmers are responsible for writing appropriate code to perform array reductions
- Recall the histogram example from lecture 4

```
#pragma omp parallel private(my_bin, i, j)
{
    #pragma omp single
    for(i=0; i<5; i++)
        bin[i]=0;

    for(i=0; i<5; i++)
        my_bin[i]=0;

    #pragma omp for
    for(i=0; i<SIZE; i++){
        j = find_bin(A[i]);
        my_bin[j]++;
    }

    #pragma omp critical
    for(i=0; i<5; i++)
        bin[i] = bin[i] + my_bin[i];
}
```

© All Rights Reserved



Induction Variables

- **Definition:** a variable inside a loop is called an *induction variable* if it is incremented or decremented by a fixed amount in each iteration of the loop
- For example, x and y are induction variables in the below loop but w and z are not

```
for ( i = 1; i < n; i ++ ) {
    x = x + 1;
    y = y - 5;
    w = w * 2;
    z = z + i;
}
```

Induction variables
cause loop-carried
dependencies

© All Rights Reserved

Induction Variable Substitution



- A techniques that rewrites induction variables inside a loop so that they are a *recurrence function* of the loop index
- By doing so, loop-carried dependencies will be eliminated

Has a loop-carried dependency

```
x = x0;  
for ( i = 1; i < n; i ++ ) {  
    x = x + a;  
    ...  
}
```



Has no loop-carried dependency

```
for ( i = 1; i < n; i ++ ) {  
    x = a * i + x0;  
    ...  
}
```

Induction loop substitution also results in making variable x private



© All Rights Reserved

Example



- Parallelize the following loop

```
x = 0;  
for ( i = 1; i < 100; i ++ ) {  
    x = x + 1;  
    A [ x ] = A[x] * A[x];  
}
```

Induction variable substitution



```
for ( i = 1; i < 100; i ++ ) {  
    x = i ;  
    A [ x ] = A[x] * A[x];  
}
```

Privatization

2

```
#pragma omp parallel for private (x)  
for ( i = 1; i < 100; i ++ ) {  
    x = i ;  
    A [ x ] = A[x] * A[x];  
}
```

© All Rights Reserved

Induction Variable Substitution Exercises



- Parallelize the following loops using proper OpenMP pragmas

1

```
x = -1 ;  
for ( i=1; i < n; i++ ) {  
    x = x + 2 ;  
    y = y + x * 2 ;  
}
```

2

```
z = 1 ;  
for ( i=0; i < n; i++ ) {  
    x = A[ z ] ;  
    z = z + 1 ;  
}
```

© All Rights Reserved

Loop Fission



- Splits a loop into two or more loops (see the below example)
- Also called loop distribution

```
for ( i=1; i < n; i++ ) {  
    a [ i ] = b [ i ] + 1 ;  
    c [ i ] = a [ i ] ;  
}
```



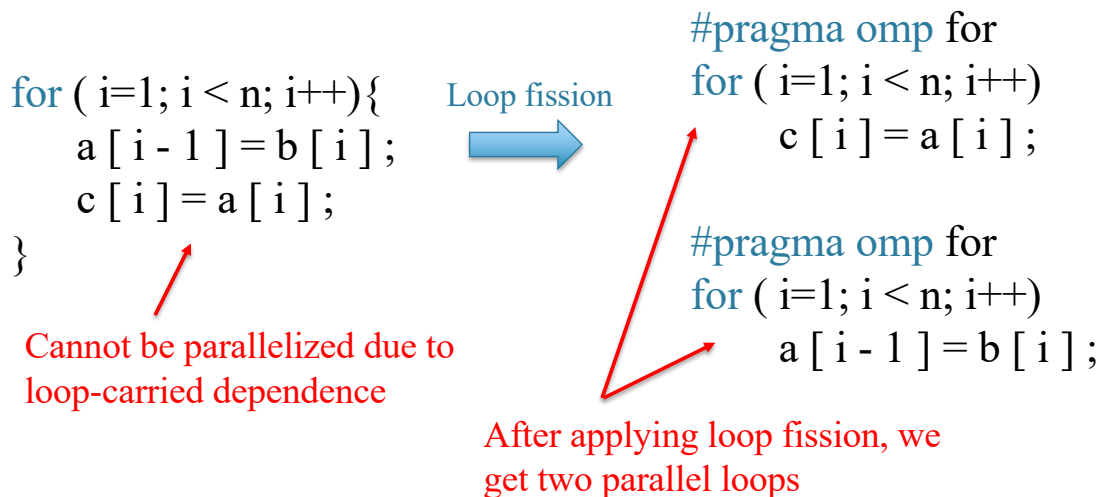
```
for ( i=1; i < n; i++ )  
    a [ i ] = b [ i ] + 1 ;
```

```
for ( i=1; i < n; i++ )  
    c [ i ] = a [ i ] ;
```

But how loop fission can be useful for parallelizing loops?

© All Rights Reserved

Splitting Sequential Loops Into Parallel Loops



Question: can you see any drawbacks of loop fission in the above example?

© All Rights Reserved

Legality of Loop Fission



- Loop fission is legal only when dependencies do NOT form a cycle inside the loop body

```
for ( i=1; i < n; i++){  
  a [ i ] = b [ i ];  
  a [ i + 1 ] = a [ i ] + c [ i ];  
}
```

Cycle exists
⇒ Loop fission is illegal

```
for ( i=1; i < n; i++){  
  b [ i ] = a [ i - 1 ];  
  a [ i ] = c [ i ];  
}
```

Cycle does not exist
⇒ Loop fission is legal

© All Rights Reserved



Loop Fission Exercises

- Parallelize the following loops using proper OpenMP pragmas

1

```
for ( i=1; i < n; i++) {  
    c [ i - 1 ] = b [ i ] ;  
    c [ i ] = a [ i ] ;  
}
```

2

```
for ( i=1; i < n; i++)  
    for ( j=1; j < n; j++) {  
        B [ i ] [ j ] = A [ i+1 ] [ j+1 ] + y ;  
        A [ i ] [ j ] = C [ i ] [ j ] ;  
    }
```

© All Rights Reserved



Using Temporaries With Loop Fission

```
for ( i=0; i < n; i++)  
    b [ i ] = b [ i + 1 ] ;
```

Has loop-carried anti-dependence

Introduce
temporary array

```
for ( i=0; i < n; i++) {  
    tmp [ i ] = b [ i + 1 ] ;  
    b [ i ] = tmp [ i ] ;  
}
```

Apply
loop fission

```
#pragma omp parallel for  
for ( i=0; i < n; i++)  
    tmp [ i ] = b [ i + 1 ] ;
```

```
#pragma omp parallel for  
for ( i=0; i < n; i++)  
    b [ i ] = tmp [ i ] ;
```

© All Rights Reserved



Bubble Sort

- Sorts N integers in ascending order
- The outer loop first finds the largest element in the list and stores it in $a[N - 1]$; it then finds the next-to-the-largest element and stores it in $a[N - 2]$, and so on

```
for(i=N-1; i>=1; i--){  
    for(j=0; j<i; j++){  
        if(a[j] > a[j+1]){  
            tmp = a[j];  
            a[j] = a[j+1];  
            a[j+1] = tmp;  
        }  
    }  
}
```

Both outer and inner loops have loop-carried dependencies

© All Rights Reserved



Odd-Even Transportation Sort

- Similar to bubble sort but with more opportunities for parallelism
- During “even” phases, each odd-subscripted element, $a[i]$, is compared to the element to its left, $a[i - 1]$, and if they’re out of order, they’re swapped
- During “odd” phases, each odd-subscripted element, $a[i]$, is compared to the element to its right, $a[i + 1]$, and if they’re out of order, they’re swapped
- A theorem guarantees that after N phases, the list will be sorted

```
for(phase=0; phase<N; phase++){  
    if(phase%2 == 0){  
        for(j=1; j<N; j+=2){  
            if(a[j-1] > a[j]){  
                tmp = a[j-1];  
                a[j-1] = a[j];  
                a[j] = tmp;  
            }  
        }  
    }  
    else {  
        for(j=1; j<N-1; j+=2){  
            if(a[j] > a[j+1]){  
                tmp = a[j];  
                a[j] = a[j+1];  
                a[j+1] = tmp;  
            }  
        }  
    }  
}
```

© All Rights Reserved

Odd-Even Transportation Sort



- Suppose $a = \{9, 7, 8, 6\}$
- The algorithm phases is shown in below

Phase	Subscript in Array			
	0	1	2	3
0	9 ↔ 7	8 ↔ 6		
	7	9	6	8
1	7	9 ↔ 6	8	
	7	6	9	8
2	7 ↔ 6	9 ↔ 8		
	6	7	8	9
3	6	7 ↔ 8	9	
	6	7	8	9

© All Rights Reserved

Odd-Even Transportation Sort OpenMP Code (Version 1)



- Outer loop is serial
- Inner loops are parallel
- The problem with this version is that the overhead of creating and terminating threads is proportional to N

```
for(phase=0; phase<N; phase++){
  if(phase%2 == 0){
    #pragma omp parallel for private(j,tmp)
    for(j=1; j<N; j+=2)
      if(a[j-1] > a[j]){
        tmp = a[j-1];
        a[j-1] = a[j];
        a[j] = tmp;
      }
  }
  else {
    #pragma omp parallel for private(j,tmp)
    for(j=1; j<N-1; j+=2)
      if(a[j] > a[j+1]){
        tmp = a[j];
        a[j] = a[j+1];
        a[j+1] = tmp;
      }
  }
}
```

© All Rights Reserved

Odd-Even Transportation Sort OpenMP Code (Version 2)



- Another OpenMP version will less parallelism overhead

```
#pragma omp parallel private(phase)
{
  for(phase=0; phase<N; phase++){
    if(phase%2 == 0) {
      #pragma omp for private(j,tmp)
      for(j=1; j<N; j+=2)
        if(a[j-1] > a[j]){
          tmp = a[j-1];
          a[j-1] = a[j];
          a[j] = tmp;
        }
    }
    else {
      #pragma omp for private(j,tmp)
      for(j=1; j<N-1; j+=2)
        if(a[j] > a[j+1]){
          tmp = a[j];
          a[j] = a[j+1];
          a[j+1] = tmp;
        }
    }
  }
}
```

© All Rights Reserved

Odd-Even Transportation Sort Performance Comparison



Odd-even sort with two parallel for directives and two for directives.
(Times are in seconds.)

thread_count	1	2	3	4
Two parallel for directives	0.770	0.453	0.358	0.305
Two for directives	0.732	0.376	0.294	0.239

© All Rights Reserved

Summary



- We determine if a loop is serial or parallel by analyzing loop-carried dependencies in the loop
- To convert a serial loop into a parallel loop, programmers need to rewrite the algorithm
- We studied six techniques that help programmers in parallelizing sequential loops

© All Rights Reserved

Exercise 1



```
for ( i=0; i < n; i++)  
  for ( j=0; j < m; j++)  
    A[ i ][ j + 2] = A[ i + 2][ j + 1] ;
```

- Draw the iteration space graph for the above loop
- Determine the distance and direction vectors

© All Rights Reserved

Exercise 2



- Normalize the following loops
- Write the dependency tests for the normalized loops
- Use the gcd test to determine if normalized loops are parallel or serial

```
for (i = 1; i < n; i += 2) {  
    A[i] = A[i + 1];  
}
```

```
for (i = 0; i < n; i += 2) {  
    A[2 * i] = A[3 * i];  
}
```

© All Rights Reserved

Exercise 3



- Parallelize the following loops using OpenMP

```
j = 1;  
for (i = 1; i < n; i++)  
{  
    j = 2 + j;  
    A[j] = A[j + 1] / 2;  
}
```

```
b = 0;  
c = 0;  
for (i = 1; i < n; i++)  
{  
    b = b + (i % 2);  
    c = c + b;  
}
```

© All Rights Reserved

Exercise 4



Question 5.8 in textbook

Consider the following loop:

```
a [ 0 ] = 0 ;  
for ( i = 1; i < n; i + + )  
    a [ i ] = a [ i - 1 ] + i ;
```

There's clearly a loop-carried dependence, as the value of $a[i]$ can't be computed without the value of $a[i - 1]$. Can you see a way to eliminate this dependence and parallelize the loop?

© All Rights Reserved

Function-Level Parallelism in OpenMP

Instructor: Dr. Fahed Jubair
Computer Engineering Department
University of Jordan

Function-Level Parallelism



- Programs with *function-level* (aka *task-level*) parallelism focuses on distributing parallel functions (i.e., tasks) on threads
- OpenMP provides two directives that are commonly used for writing task-parallel programs:
 1. The SECTIONS directive } We will focus on this directive
 2. The TASK directive } Beyond the scope of this class

© All Rights Reserved.

Quick Review The SECTIONS Directive



- A SECTIONS directive specifies that all enclosed section(s) of code are to be divided between threads such that each section is executed once by a thread

```
#pragma omp parallel sections
{
    #pragma omp section
    function1();

    #pragma omp section
    function2();

    #pragma omp section
    function3();
}
```

© All Rights Reserved.

Challenges



- There are two primary challenges that programmers need to tackle when converting a sequential algorithm into a task-parallel algorithm:
 - How to maximize the number of tasks that can run in parallel?
 - How to assign these parallel tasks to threads such that parallelism overhead is minimized?

© All Rights Reserved.

Algorithms With Function-Level Parallelism



We will study two sorting algorithms that take advantage of function-level parallelism:

- Merge Sort
- Quick Sort

© All Rights Reserved.



Merge Sort

- **Split** the unsorted list into n sublists, each containing one element
- Repeatedly **merge** sublists to produce new sorted sublists until there is only 1 sublist remaining, which will be the sorted list

// A is the input unsorted list and B is the output sorted list

```
void merge_sort (int A[ ], int begin, int end, int B[ ])
{
    if ( begin == end ) // if A has only one integer
        return A[begin] ;
    middle = (begin + end) / 2 ;
    merge_sort ( A, begin, middle - 1, left) ; // sort left half [begin : middle - 1]
    merge_sort ( A, middle, end, right) ; // sort right half [middle : end]
    merge (left, right, B) ; // merge left and right sublists into B
}
```

© All Rights Reserved.



Merge Sort Example Split Phase

99	6	86	15	58	35	86	4	0
----	---	----	----	----	----	----	---	---

99	6	86	15	58	35	86	4	0
----	---	----	----	----	----	----	---	---

99	6	86	15	58	35	86	4	0
----	---	----	----	----	----	----	---	---

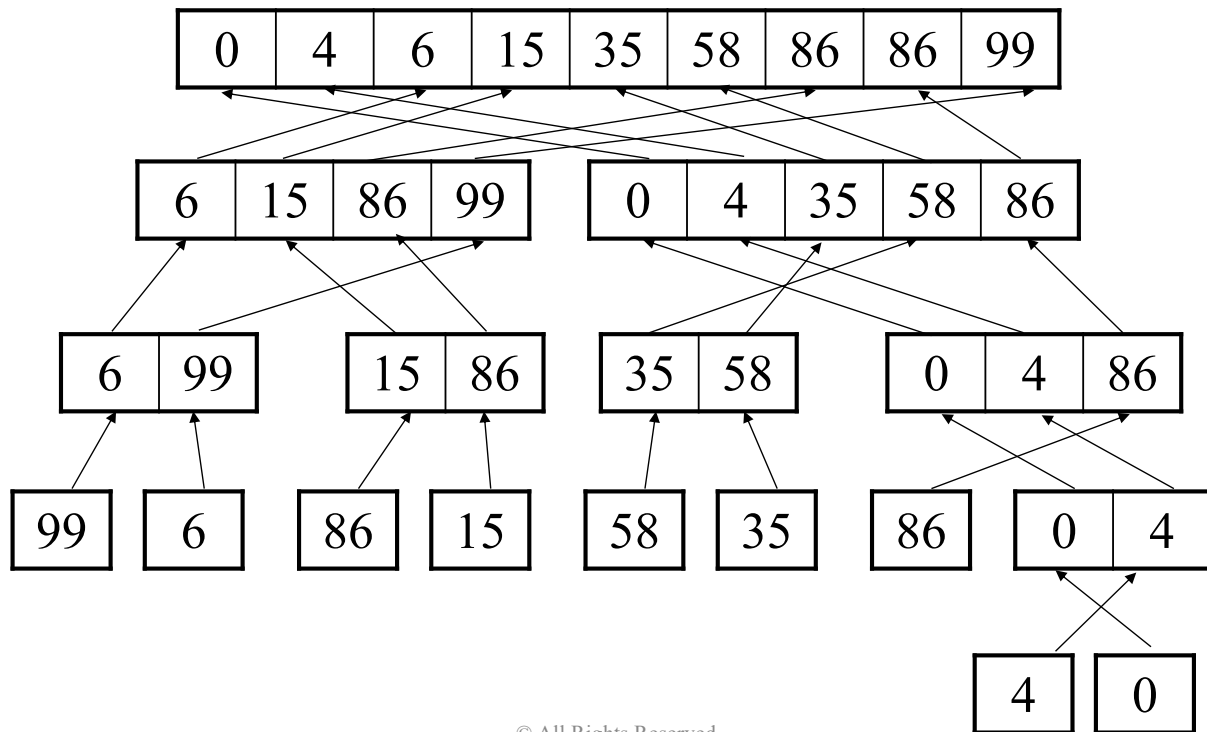
99	6	86	15	58	35	86	4	0
----	---	----	----	----	----	----	---	---

4	0
---	---

© All Rights Reserved.

Merge Sort Example

Merge Phase



Parallelism In Merge Sort



- Upon careful examine of the merge sort code, we can observe that the two recursive calls of merge sort can run in parallel

```
void merge_sort (int A[ ], int begin, int end, int B[ ])
```

```
{
```

```
    if ( begin == end )
```

```
        return A[begin] ;
```

```
    middle = (begin + end) / 2 ;
```

```
    #pragma omp parallel sections num_threads(2)
```

```
    { /* start of parallel region */
```

```
        #pragma omp section
```

```
        merge_sort ( A, begin, middle - 1, left) ;
```

```
        #pragma omp section
```

```
        merge_sort ( A, middle, end, right);
```

Run in parallel

```
    } /* end of parallel region */
```

```
    merge (left, right, B);
```

© All Rights Reserved.

Hold On!

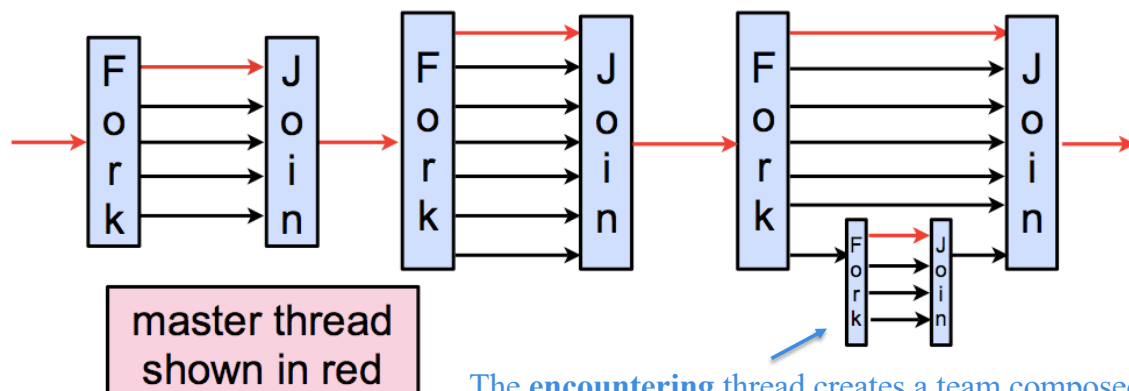
- The recursive call of merge_sort in the previous slide will cause parallel regions to be created within each other!
- Does OpenMP even allow this? If so, then how does the parallel execution occur?



© All Rights Reserved.

Nested Parallelism

- OpenMP parallel regions can be nested inside each other
- A thread within a team of threads may spawn a team of threads



The **encountering** thread creates a team composed of itself and some additional (possibly zero) number of threads (the encountering thread becomes the master thread)

© All Rights Reserved.

More about Nested Parallelism



- It is possible to enable/disable nested parallelism in OpenMP using the environment variable `OMP_NESTED` or the library routine `omp_set_nested()`
- If nested parallelism is enabled, when a thread on a team within a parallel region encounters a new parallel construct, an additional team of threads is forked off of it, and it becomes their master
- If nested parallelism is disabled, when a thread on a team within a parallel region encounters a new parallel construct, execution continues on this additional single thread only

© All Rights Reserved.

Example on Nested Parallel Regions



```
#pragma omp parallel num_threads(4)
{
    printf("I am thread %d from region 1\n",omp_get_thread_num());
    #pragma omp barrier

    #pragma omp single
    {
        #pragma omp parallel num_threads(2)
        {
            printf("I am thread %d from region 2\n",omp_get_thread_num());
            #pragma omp barrier

            #pragma omp single
            {
                #pragma omp parallel num_threads(8)
                {
                    printf("I am thread %d from region 3\n",omp_get_thread_num());
                }
            }
        }
    }
}
```

© All Rights Reserved.

Example on Nested Parallel Regions (cont.)



```
$ gcc -fopenmp nested.c -o a.out  
$ export OMP_NESTED = TRUE  
$ ./a.out
```

```
I am thread 3 from region 1  
I am thread 0 from region 1  
I am thread 2 from region 1  
I am thread 1 from region 1  
I am thread 0 from region 2  
I am thread 1 from region 2  
I am thread 7 from region 3  
I am thread 4 from region 3  
I am thread 0 from region 3  
I am thread 3 from region 3  
I am thread 1 from region 3  
I am thread 2 from region 3  
I am thread 6 from region 3  
I am thread 5 from region 3
```

© All Rights Reserved.

Remarks on Nested Parallelism



- Nesting parallel regions provides an immediate way to allow more threads to participate in the computation
- However, nesting parallel regions can easily create too many threads and oversubscribe the system
- In addition, creating nested parallel regions adds overhead
- To avoid the above problems, programmers usually try to impose some discipline on the execution using appropriate OpenMP environment variables
 - For example, “[OMP_MAX_ACTIVE_LEVELS](#)” environment variable controls the maximum allowed number of active nested parallel regions

© All Rights Reserved.

Quick Sort



Sorts an array of integers as follows:

- Pick an element, called a **pivot**, from the array
- Reorder the array so that all elements with values less than the pivot come before the pivot, while all elements with values greater than the pivot come after it (equal values can go either way)
- Recursively apply the above steps to the sub-array of elements with smaller values and separately to the sub-array of elements with greater values

Source:
Wikipedia

© All Rights Reserved.

Quick Sort (cont.)



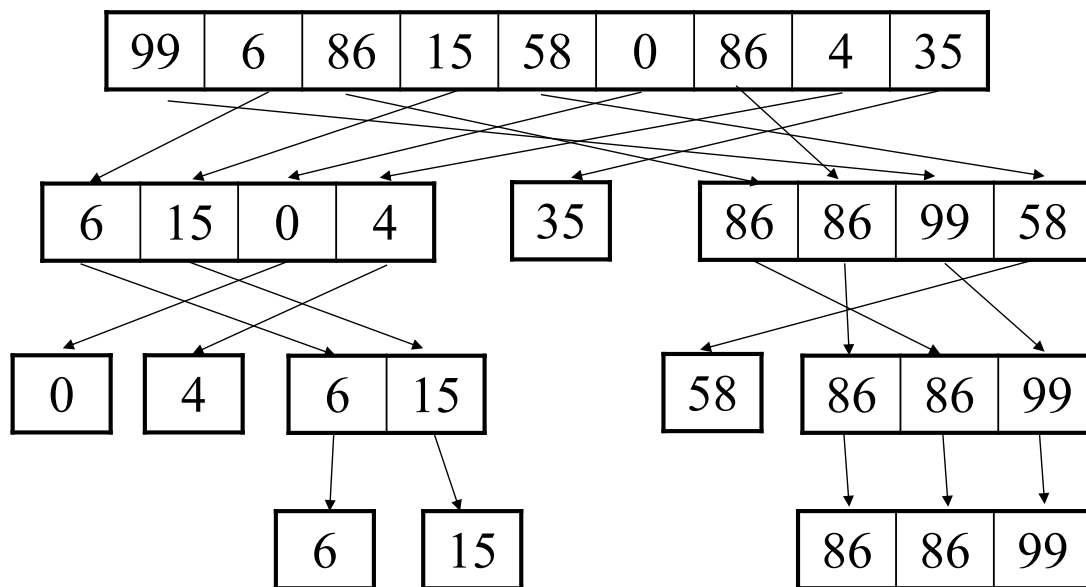
```
void quick_sort (int A[ ], int begin, int end) {
    if ( begin < end ) {
        p = partition (A, begin, end ) ;
        quick_sort ( A, begin, p - 1 ) ;
        quick_sort ( A, p + 1, end ) ;
    }
}

int partition (int A[ ], int begin, int end) {
    pivot = A[end] ; // let us pick the last element as our pivot
    i = begin ;
    for ( j = begin; j < end; j++ ) {
        if ( A[j ] <= pivot ) {
            swap ( A[j ] , A[ i ] ); // swap the content of A[ j ] and A[ i ]
            i++ ;
        }
    }
    swap ( A[ end ] , A[ i ] );
    return i ;
}
```

© All Rights Reserved.



Quick Sort Example



© All Rights Reserved.



Parallelism In Quick Sort

- Upon careful examine of the quick sort code, we can observe that the two recursive calls of quick sort can run in parallel

```

void quick_sort (int A[ ], int begin, int end) {
    if ( begin < end ) {
        p = partition (A, begin, end) ;
        #pragma omp parallel sections num_threads(2)
        {
            #pragma omp section
            quick_sort ( A, begin, p - 1 ) ;
            #pragma omp section
            quick_sort ( A, p + 1, end ) ;
        }
    }
}

```

Run in parallel

© All Rights Reserved.

Exercise



- Write a C function that returns the occurrences count of an integer x inside A , an array of randomly generated integers between $[0: 99]$. The algorithm should work as follows:
 1. Split the array A into two halves: left and right
 2. Count x occurrences in the left sub-array
 3. Count x occurrences in the right sub-array
 4. The total count is the summation of the left and right counts
 5. Do the above steps recursively
- Now, parallelize your code using OpenMP

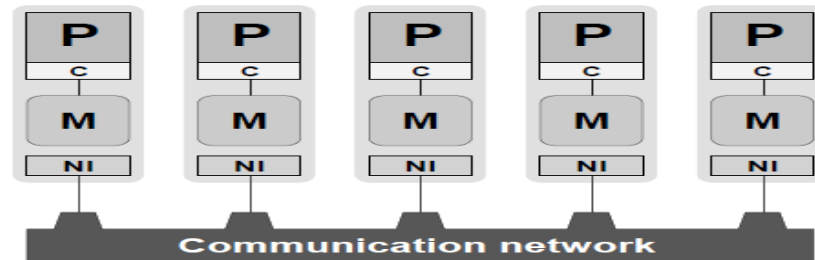
© All Rights Reserved.

Introduction to MPI

Instructor: Dr. Fahed Jubair
Computer Engineering Department
University of Jordan



Distributed-Memory Architectures



- A set of processors that are connected by a network
- Each processor runs a **local program on a local memory**
- Processors **do not** share a common memory, instead they communicate by **passing messages** through the network
- **Message passing programming models** are introduced to enable writing parallel programs that target distributed memory architectures

© All Rights Reserved.

2

MPI



- *Message Passing Interface* (MPI) is the most commonly-used standard for message passing programs
- MPI is a standard developed by MPI Forum (<http://www.mpi-forum.org>) that specifies how to implement a library that can be used with conventional sequential languages such as Fortran, C and C++ for writing SPMD-style parallel programs for distributed-memory architectures
- MPI standard versions: **1.0** (1994), **2.0** (1997), **3.0** (2012)
- Several implementations for MPI are available
 - E.g., MPICH, OpenMPI, Intel, Microsoft



© All Rights Reserved.

3

Execution Model of MPI

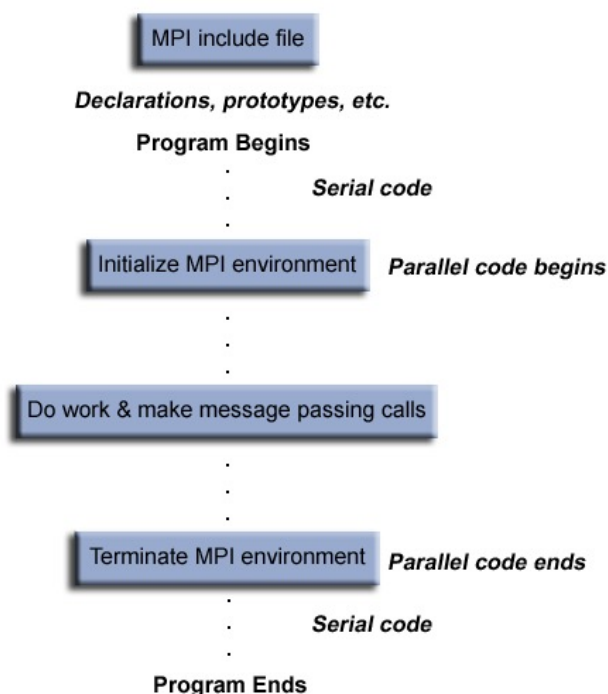


- Fixed number of processes (determined on startup) that starts the execution together and lasts till the execution terminates
- Each process is given a unique rank between 0 and (*number_of_processes* – 1)
- Each process has its instruction and data address spaces
 - No physical sharing between processes
 - All data must be explicitly partitioned and placed
 - No data races due to hardware artifacts
- To exchange data, processes use library routines to communicate with each other
- Essentially, each process runs like a sequential process except that it makes calls to the MPI library

© All Rights Reserved.

4

A Generic MPI Program Structure



Source:

<https://computing.llnl.gov/tutorials/mpi>

© All Rights Reserved.

5



MPI Library Routines

- **Environment management routines:** used for setting and interrogating the MPI execution environment
 - E.g., initializing and terminating the MPI environment
 - E.g., querying the number of processes
 - E.g., querying a process rank
- **Communication routines:** generally, divided into two types of routines:
 - **Point-to-Point:** a message is sent from a specific sending process to a specific receiving process
 - **Collective:** **all** processes participate in the communication



MPI Communicator

- MPI uses objects called *communicators*, which are a collection of processes that can send messages to each other
- MPI execution, as startup, defines a default communicator that encapsulates all the processes in the system, called `MPI_COMM_WORLD`
- Most MPI routines require you to specify a communicator as an argument



Environment Routines

- `MPI_Init (&arg1, &arg2)`
 - Initializes the MPI execution environment
 - This function must be called in every MPI program, must be called before any other MPI functions and must be called only once in an MPI program
- `MPI_Comm_size (comm, &size)`
 - Returns the total number of MPI processes in the specified communicator, such as `MPI_COMM_WORLD`
- `MPI_Comm_rank (comm, &rank)`
 - Returns the rank of the calling MPI process within the specified communicator
- `MPI_Finalize ()`
 - Terminates the MPI execution environment
 - This function should be the last MPI routine called in every MPI program - no other MPI routines may be called after it

© All Rights Reserved.

8



The Traditional Hello World Example

```
#include <mpi.h>
#include <stdio.h>

int main(int argc, char *argv[]){
    int rank, size;

    MPI_Init(&argc, &argv);

    MPI_Comm_rank(MPI_COMM_WORLD, &rank);

    MPI_Comm_size(MPI_COMM_WORLD, &size);

    printf("Hello from process %d\n", rank);
    if(rank == 0)
        printf("number of processes is %d\n", size);

    MPI_Finalize();

    return 0;
}
```

© All Rights Reserved.

9

Compiling Hello World Example



- We can compile the Hello World OpenMP program on a linux machine using the following command

wrapper script to compile
source file

```
mpicc mpi_hello.c -o mpi_hello
```

*create this executable file name
(as opposed to default a.out)*

Running Hello World Example



```
mpiexec -n <number of processes> <executable>
```

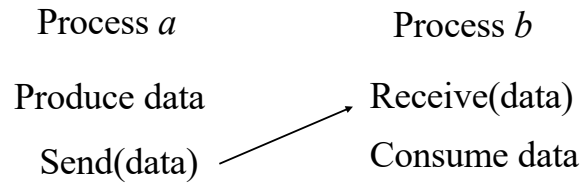
```
mpiexec -n 4 ./mpi_hello
```

run with 4 processes

One possible output

```
Hello from process 2  
Hello from process 0  
number of processes is 4  
Hello from process 1  
Hello from process 3
```

MPI Point-to-Point Communication



- Basic point-to-point communication in MPI is **two-sided**
 - Process *a* must make a call to a send routine and specify process *b* as the receiver in the send message arguments
 - Process *b* must make a call to a receive routine and specify process *a* as the sender in the receive message arguments
 - Both, send and receive routines must have matching communicator object and similar description (size and type) of the data being sent/received

Point to Point Communication Routines



- Generally, most of the MPI point-to-point routines can be used in either *blocking* or *non-blocking* mode
- Let us first consider the basic type: blocking communication
- Blocking send/receive routines have the following properties:
 - A blocking send routine will only "return" after it is safe to modify the application buffer (your send data) for reuse
 - Safe means that modifications will not affect the data intended for the receive task
 - Safe does not imply that the data was actually received - it may very well be sitting in a system buffer
 - A blocking receive only "returns" after the data has arrived and is ready for use by the program

Blocking Send/Receive Routines



`MPI_Send(buffer, count, type, dest, tag, comm)`

`MPI_Recv(buffer, count, type, source, tag, comm, status)`

- Buffer: the variable name that is being sent/received
- Count: number of data elements to be sent
- Type: type of data to be sent
- Source: rank of the sender process (needed for receive messages)
- Dest: rank of the receiver process (needed for send messages)
- Tag: arbitrary non-negative integer assigned by the programmer to uniquely identify a message (send and receive operations should match message tags)
- Comm: communicator object
- Status: contains further information about received messages

© All Rights Reserved.

14

MPI Elementary Data Types



For reasons of portability, MPI predefines its elementary data types

MPI datatype	C datatype
MPI_CHAR	signed char
MPI_SHORT	signed short int
MPI_INT	signed int
MPI_LONG	signed long int
MPI_LONG_LONG	signed long long int
MPI_UNSIGNED_CHAR	unsigned char
MPI_UNSIGNED_SHORT	unsigned short int
MPI_UNSIGNED	unsigned int
MPI_UNSIGNED_LONG	unsigned long int
MPI_FLOAT	float
MPI_DOUBLE	double
MPI_LONG_DOUBLE	long double
MPI_BYTE	
MPI_PACKED	

- Programmers may also create their own data types (called **Derived Data Types**)

© All Rights Reserved.

15

Error Handling



- Most MPI routines include a return/error code parameter that can be tested to check if a routine call succeeds/fails
- However, the default behavior of failed routines is to abort the execution
- Users can override this behavior (if necessary) to capture this error
- The code MPI_SUCCESS (zero) is returned when a routine succeeds

```
#include <mpi.h>
#include <stdio.h>

int main(int argc, char *argv[]){
    int rank, size;
    int sbuff, rbuff;
    MPI_Status stat;
    int err_code;
    int tag = 1;

    MPI_Init(&argc, &argv);

    MPI_Comm_rank(MPI_COMM_WORLD, &rank);

    MPI_Comm_size(MPI_COMM_WORLD, &size);

    if(rank == 0){
        sbuff = 100;
        err_code = MPI_Send(&sbuff, 1, MPI_INT, 1, tag, MPI_COMM_WORLD);
    }
    else{
        err_code = MPI_Recv(&rbuff, 1, MPI_INT, 0, tag, MPI_COMM_WORLD, &stat);
        printf("data=%d is received\n", rbuff);
    }

    MPI_Finalize();

    return 0;
}
```

MPI Program Example 1



Compiling and Running Example 1



```
$ mpicc example1.c -o a.out
```

```
$ mpiexec -n 2 ./a.out
```

data=100 is received

```
#include <mpi.h>
#include <stdio.h>
#include <stdlib.h>
#define N 16

int main(int argc, char *argv[]){
    int rank, size, chunk, i, tag=1;
    float *sbuff, *rbuff;
    MPI_Status stat;

    MPI_Init(&argc, &argv);

    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    MPI_Comm_size(MPI_COMM_WORLD, &size);

    chunk = N / size;

    if(rank == 0){
        sbuff = (float*) malloc(N * sizeof(float));
        for(i=0; i<N; i++)
            sbuff[i] = 10.0 * i;

        for(i=0; i<size; i++)
            MPI_Send(&sbuff[chunk*i], chunk, MPI_FLOAT, i, tag, MPI_COMM_WORLD);
    }

    rbuff = (float*) malloc(chunk * sizeof(float));
    MPI_Recv(rbuff, chunk, MPI_FLOAT, 0, tag, MPI_COMM_WORLD, &stat);
    printf("Process %d received the following data\n", rank);
    for(i=0; i<chunk; i++)
        printf("%.3f\n", rbuff[i]);

    MPI_Finalize();

    return 0;
}
```

MPI Program Example 2



Compiling and Running Example 2



```
$ mpicc example2.c -o a.out
$ mpiexec -n 8 ./a.out
```

```
Process 0 received the following data
0.000
10.000
Process 1 received the following data
20.000
30.000
Process 2 received the following data
40.000
50.000
Process 3 received the following data
60.000
70.000
Process 4 received the following data
80.000
90.000
Process 5 received the following data
100.000
110.000
Process 6 received the following data
120.000
130.000
Process 7 received the following data
140.000
150.000
```

```
#include <mpi.h>
#include <stdio.h>
#include <stdlib.h>
#define N 512

int main(int argc, char *argv[]){
    int rank, size, chunk, i, tag=1;
    int A[N];
    int *buff;
    MPI_Status stat;
    int local_sum, global_sum;

    MPI_Init(&argc, &argv);

    MPI_Comm_rank(MPI_COMM_WORLD, &rank);

    MPI_Comm_size(MPI_COMM_WORLD, &size);

    chunk = N / size;

    if(rank == 0){
        for(i=0; i<N; i++){
            A[i] = rand() % 100; // each integer is randmly picked from 0-99
            for(i=0; i<size; i++){
                MPI_Send(&A[chunk*i], chunk, MPI_INT, i, tag, MPI_COMM_WORLD);
            }
        }

        buff = (int*) malloc(chunk * sizeof(int));
        MPI_Recv(buff, chunk, MPI_INT, 0, tag, MPI_COMM_WORLD, &stat);
        local_sum = 0;
        for(i=0; i<chunk; i++){
            local_sum += buff[i];
        }
        MPI_Send(&local_sum, 1, MPI_INT, 0, tag, MPI_COMM_WORLD);

        if(rank==0){
            global_sum = 0;
            for(i=0; i<size; i++){
                int tmp;
                MPI_Recv(&tmp, 1, MPI_INT, i, tag, MPI_COMM_WORLD, &stat);
                global_sum += tmp;
            }
            printf("global sum is %d\n", global_sum);
        }
        MPI_Finalize();
        return 0;
    }
}
```

MPI Program Example 3 Parallel Sum

} Process 0 distributes data

} Each process computes local sum

} Process 0 collects local sums to compute global sum



Non-Blocking Communication



- The other type of point-to-point communication in MPI
- Non-blocking send/receive routines have the following properties:
 - Non-blocking send and receive routines return almost immediately, i.e., they do not wait for any communication events to complete, such as message copying from user memory to system buffer space or the actual arrival of message
 - Non-blocking operations simply "request" the MPI library to perform the operation when it is able (the user can not predict when that will happen)
 - It is unsafe to the user to modify data until you know for a fact that the requested non-blocking operation was actually performed by the library (which can be done using "wait" routines)
 - Non-blocking communications are primarily used to **overlap computation with communication and exploit possible performance gains**

Non-Blocking Send/Receive Routines



`MPI_Isend(buffer, count, type, dest, tag, comm, request)`

`MPI_Irecv(buffer, count, type, source, tag, comm, request)`

- The only new argument we did not see before is “request”
- Since non-blocking operations may return before the requested system buffer space is obtained, the system issues a unique "request number”
- The programmer uses this system assigned "handle" later (in a WAIT type routine) to determine completion of the non-blocking operation
- In C, the “request” argument is a pointer to a predefined structure `MPI_Request`



Wait Routines

`MPI_Wait (request, status)`

`MPI_Waitall (count, array_of_requests, array_of_statuses)`

- `MPI_Wait` blocks the calling process until a specified non-blocking send or receive operation has completed
- For multiple non-blocking operations, the programmer can specify any, all or some completions

```
#include <mpi.h>
#include <stdio.h>
```

```
int main(int argc, char *argv[]){
    int rank, size, tag=1;
    int dst, sndr;
    int sbuff, rbuff;
    MPI_Status stats[2];
    MPI_Request reqs[2];
```

```
    MPI_Init(&argc, &argv);
```

```
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
```

```
    MPI_Comm_size(MPI_COMM_WORLD, &size);
```

```
    dst = (rank + 1) % size ;
    sndr = (rank + size - 1) % size ;
```

```
    MPI_Irecv(&rbuff, 1, MPI_INT, sndr, tag, MPI_COMM_WORLD, &reqs[0]);
```

```
    sbuff = rank * 10 ;
    MPI_Isend(&sbuff, 1, MPI_INT, dst, tag, MPI_COMM_WORLD, &reqs[1]);
```

```
    MPI_Waitall(2, reqs, stats);
```

```
    printf("process %d rbuff = %d\n", rank, rbuff);
```

```
    MPI_Finalize();
```

```
    return 0;
```

```
}
```



MPI Program Example 4

Compiling and Running Example 4



```
$ mpicc example4.c -o a.out
$ mpiexec -n 16 ./a.out
```

```
process 1 rbuff = 0
process 2 rbuff = 10
process 3 rbuff = 20
process 4 rbuff = 30
process 5 rbuff = 40
process 8 rbuff = 70
process 9 rbuff = 80
process 10 rbuff = 90
process 11 rbuff = 100
process 12 rbuff = 110
process 14 rbuff = 130
process 15 rbuff = 140
process 0 rbuff = 150
process 13 rbuff = 120
process 6 rbuff = 50
process 7 rbuff = 60
```

Example 5



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

- Let us try to parallelize the above loop using MPI as follows:
 - Process 0 initializes arrays a and b , partitions them into equal chunks and distributes these chunks onto all processes
 - Each process locally multiplies its chunks to obtain array c
 - Process 0 collects array c from all processes



```
#include <mpi.h>
#include <stdio.h>
#include <stdlib.h>
#define N 16

int main(int argc, char *argv[]){
    int rank, size, chunk, i, m, tag1=1, tag2=2;
    double a[N], b[N], c[N];
    double *buff1, *buff2;
    MPI_Status *stats, *stats0, *stats1;
    MPI_Request *reqs, *reqs0, *reqs1;

    MPI_Init(&argc, &argv);

    MPI_Comm_rank(MPI_COMM_WORLD, &rank);

    MPI_Comm_size(MPI_COMM_WORLD, &size);

    chunk = N / size;

    if(rank == 0){
        for(i=0; i<N; i++){
            a[i] = 12.4 * (i % 100);
            b[i] = -3.33 * (i % 100);
        }
        reqs0 = (MPI_Request*) malloc(2* size * sizeof(MPI_Request));
        stats0 = (MPI_Status*) malloc(2* size * sizeof(MPI_Status));
        m = 0;
        for(i=0; i<size; i++){
            MPI_Isend(&a[chunk*i], chunk, MPI_DOUBLE, i, tag1, MPI_COMM_WORLD, &reqs0[m++]);
            MPI_Isend(&b[chunk*i], chunk, MPI_DOUBLE, i, tag2, MPI_COMM_WORLD, &reqs0[m++]);
        }
        MPI_Waitall(m, reqs0, stats0);
    }
}
```

MPI Program Example 5

Using non-blocking sends allows the waiting times of send messages to be overlapped

© All Rights Reserved.

28



MPI Program Example 5 (cont.)

```
reqs = (MPI_Request*) malloc(2 * sizeof(MPI_Request));
stats = (MPI_Status*) malloc(2 * sizeof(MPI_Status));
buff1 = (double*) malloc(chunk * sizeof(double));
MPI_Irecv(buff1, chunk, MPI_DOUBLE, 0, tag1, MPI_COMM_WORLD, &reqs[0]);
buff2 = (double*) malloc(chunk * sizeof(double));
MPI_Irecv(buff2, chunk, MPI_DOUBLE, 0, tag2, MPI_COMM_WORLD, &reqs[1]);
MPI_Waitall(2, reqs, stats);
```

```
for(i=0; i<chunk; i++)
    buff1[i] = buff1[i] * buff2[i];
```

```
MPI_Send(buff1, chunk, MPI_DOUBLE, 0, tag1, MPI_COMM_WORLD);
```

the waiting times of the two receive messages are overlapped

Each process sends its result back to process 0

© All Rights Reserved.

29



MPI Program Example 5 (cont.)

```
if(rank == 0){
    reqs1 = (MPI_Request*) malloc(size * sizeof(MPI_Request));
    stats1 = (MPI_Status*) malloc(size * sizeof(MPI_Status));
    m = 0;
    for(i=0; i<size; i++)
        MPI_Irecv(&c[chunk*i], chunk, MPI_DOUBLE, i, tag1, MPI_COMM_WORLD, &reqs1[m++]);

    MPI_Waitall(m, reqs1, stats1);

    printf("array a: ");
    for(i=0; i<N; i++)
        printf("%.4f, ", a[i]);
    printf("\narray b: ");
    for(i=0; i<N; i++)
        printf("%.4f, ", b[i]);
    printf("\narray c: ");
    for(i=0; i<N; i++)
        printf("%.4f, ", c[i]);
    printf("\n");
}

MPI_Finalize();
return 0;
}
```

the waiting times of receive messages are overlapped

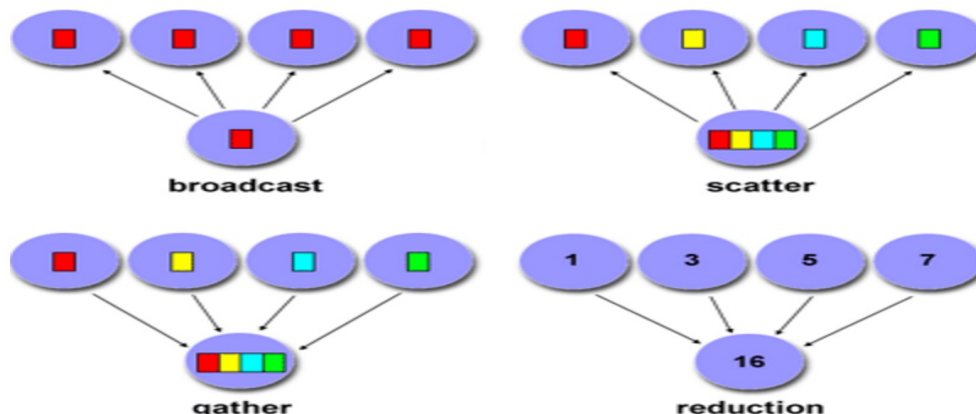
© All Rights Reserved.

30



Collective Communication Routines

- Unlike point-to-point communication, collective communication involves all processes within the scope of a communicator
- We will consider the following collective communication types:



© All Rights Reserved.

31



MPI Scatter Routine

`MPI_Scatter` (sendbuf, sendcnt, sendtype,
recvbuf, recvcnt, recvtype, root, comm)

- Sendbuf: address of send buffer
- Sendcount: number of sent data elements for a single message
- Sendtype: type of data being sent
- Recvbuf: address of receive buffer
- Recvcount: number of received data elements
- Recvtype: type of data being received
- **Root: rank of the sender process**
- Comm: communicator object



MPI Gather Routine

`MPI_Gather` (sendbuf, sendcnt, sendtype,
recvbuf, recvcnt, recvtype, root, comm)

- Sendbuf: address of send buffer
- Sendcount: number of sent data elements
- Sendtype: type of data being sent
- Recvbuf: address of receive buffer
- Recvcount: number of received data elements of a single message
- Recvtype: type of data being received
- **Root: rank of the receiver process**
- Comm: communicator object



```
#include <mpi.h>
#include <stdio.h>
#include <stdlib.h>
#define N 16

int main(int argc, char *argv[]){
    int rank, size, chunk, i;
    int *A, *tmp;

    MPI_Init(&argc, &argv);
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    MPI_Comm_size(MPI_COMM_WORLD, &size);

    chunk = N / size;
    tmp = (int*) malloc(chunk * sizeof(int));
    if(rank == 0){
        A = (int*) malloc(N * sizeof(int));
        printf("A before multiplication: ");
        for(i=0; i<N; i++){
            A[i] = i;
            printf("%d, ", A[i]);
        }
    }

    MPI_Scatter(A, chunk, MPI_INT, tmp, chunk, MPI_INT, 0, MPI_COMM_WORLD);

    for(i=0; i<chunk; i++)
        tmp[i] = tmp[i] * tmp[i] ;

    MPI_Gather(tmp, chunk, MPI_INT, A, chunk, MPI_INT, 0, MPI_COMM_WORLD);

    if(rank == 0){
        printf("\nA after multiplication: ");
        for(i=0; i<N; i++)
            printf("%d, ", A[i]);
        printf("\n");
    }
    MPI_Finalize();
    return 0;
}
```

MPI Program Example 6

© All Rights Reserved.

34

Compiling and Running Example 6



```
$ mpicc example6.c -o a.out
```

```
$ mpiexec -n 8 ./a.out
```

A before multiplication: 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15,

A after multiplication: 0, 1, 4, 9, 16, 25, 36, 49, 64, 81, 100, 121, 144, 169, 196, 225,

© All Rights Reserved.

35



MPI Reduce Routine

`MPI_Reduce (sendbuf, recvbuf, count, datatype, op, root, comm)`

- Sendbuf: address of send buffer
- Recvbuf: address of receive buffer
- Count: number of sent data elements
- Datatype: type of data being sent
- Op: reduce operation
- **Root: rank of the process where the final reduced value will be available**
- Comm: communicator object



Reduction Operation in MPI

Operation	Meaning	Datatypes
MPI_MAX	Maximum	C integers and floating point
MPI_MIN	Minimum	C integers and floating point
MPI_SUM	Sum	C integers and floating point
MPI_PROD	Product	C integers and floating point
MPI_LAND	Logical AND	C integers
MPI_BAND	Bit-wise AND	C integers and byte
MPI_LOR	Logical OR	C integers
MPI_BOR	Bit-wise OR	C integers and byte
MPI_LXOR	Logical XOR	C integers
MPI_BXOR	Bit-wise XOR	C integers and byte
MPI_MAXLOC	max-min value-location	Data-pairs
MPI_MINLOC	min-min value-location	Data-pairs



```
#include <mpi.h>
#include <stdio.h>
#include <stdlib.h>
#define N 512

int main(int argc, char *argv[]){
    int rank, size, chunk, i;
    int A[N];
    int *buff;
    int local_sum, global_sum;

    MPI_Init(&argc, &argv);
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    MPI_Comm_size(MPI_COMM_WORLD, &size);

    chunk = N / size;
    buff = (int*) malloc(chunk * sizeof(int));
    if(rank == 0){
        for(i=0; i<N; i++)
            A[i] = rand() % 100; // each integer is randomly picked from 0-99
    }

    MPI_Scatter(A, chunk, MPI_INT, buff, chunk, MPI_INT, 0, MPI_COMM_WORLD);

    local_sum = 0;
    for(i=0; i<chunk; i++)
        local_sum += buff[i];

    MPI_Reduce(&local_sum, &global_sum, 1, MPI_INT, MPI_SUM, 0, MPI_COMM_WORLD);

    if(rank==0)
        printf("global sum is %d\n", global_sum);

    MPI_Finalize();
    return 0;
}
```

MPI Program Example 7 Parallel Sum

© All Rights Reserved.

38



MPI Broadcast Routine

`MPI_Bcast (buffer, count, datatype, root, comm)`

- Buffer: address of send/receive buffer
- Count: number of sent data elements
- Datatype: type of data being sent
- **Root: rank of the process doing the broadcasting**
- Comm: communicator object

© All Rights Reserved.

39



MPI Program Example 8

```
#include <mpi.h>
#include <stdio.h>
#include <stdlib.h>
#define N 512

int main(int argc, char *argv[]){
    int rank, size, chunk, i;
    int A[N];
    int *buff;
    int local_sum, global_sum;

    MPI_Init(&argc, &argv);
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    MPI_Comm_size(MPI_COMM_WORLD, &size);

    chunk = N / size;
    buff = (int*) malloc(chunk * sizeof(int));
    if(rank == 0){
        for(i=0; i<N; i++)
            A[i] = rand() % 100; // each integer is randmly picked from 0-99
    }

    MPI_Scatter(A, chunk, MPI_INT, buff, chunk, MPI_INT, 0, MPI_COMM_WORLD);

    local_sum = 0;
    for(i=0; i<chunk; i++)
        local_sum += buff[i];

    MPI_Reduce(&local_sum, &global_sum, 1, MPI_INT, MPI_SUM, 0, MPI_COMM_WORLD);
    MPI_Bcast(&global_sum, 1, MPI_INT, 0, MPI_COMM_WORLD);
    printf("process %d, global sum is %d\n", rank, global_sum);

    MPI_Finalize();
    return 0;
}
```

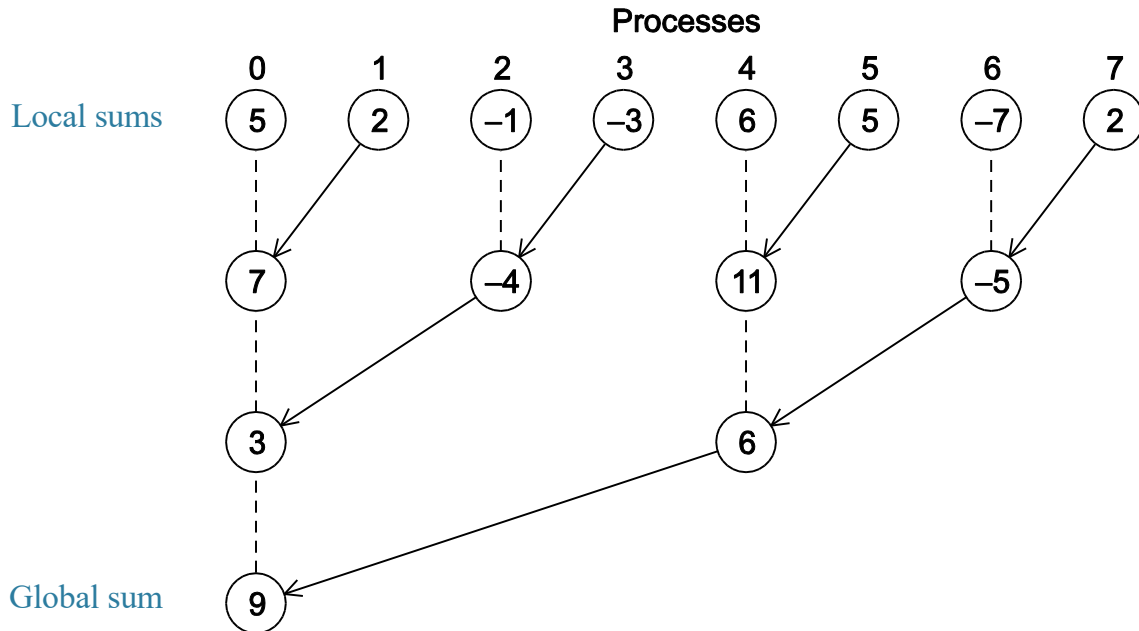
Global sum is broadcasted to every process



Wait A Minute!

- Do we really need collective communication routines? Can't we use send and receive routines for everything?
- Collective operations are introduced to give the MPI runtime environment opportunities to **optimize common communication patterns**
- For example, MPI runtime environments use advanced algorithms (such as *tree-structured* communication algorithms) to perform collective communication routines faster
- MPI runtime environments also often try to utilize the underlying network topology in reducing communication latency

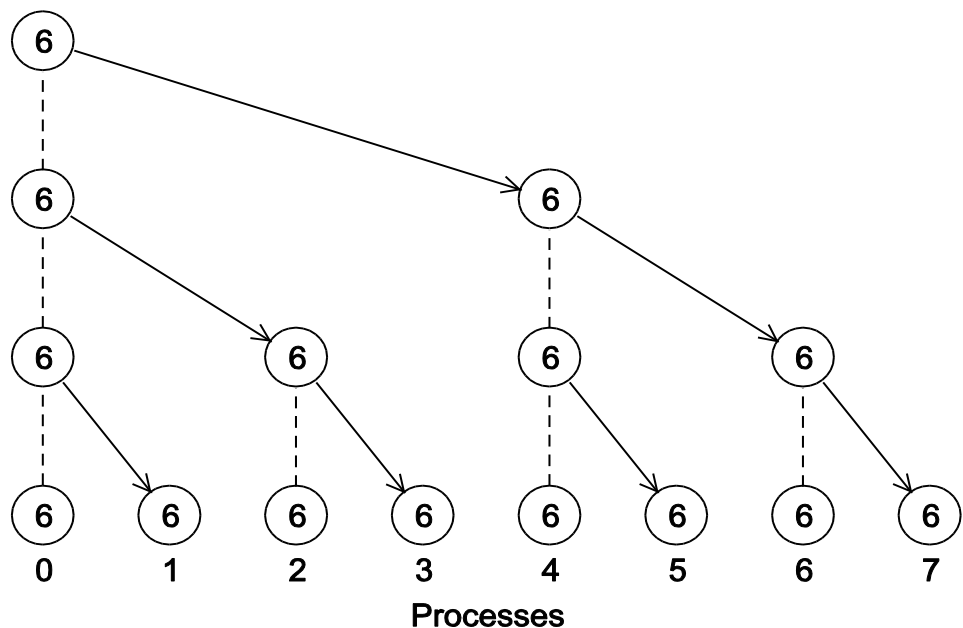
Tree-Structured MPI_Reduce



Tree-Structured MPI_Bcast



P0 broadcasts 6 to P0-P7



More Collective Communication Routines

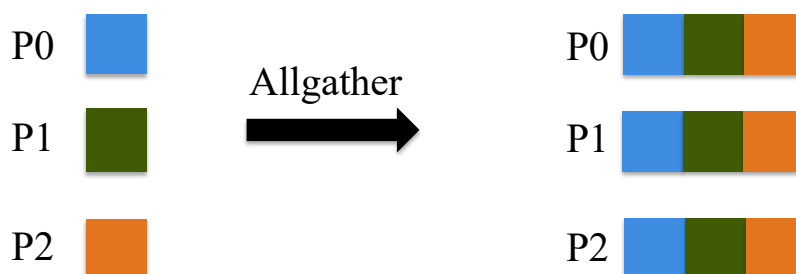


- So far, we have focused on collective routines that perform one-to-all or all-to-one communication patterns
- We will consider the following three collective routines, which perform all-to-all communication patterns:
 - **Allgather**: gathers all data elements from all processes and broadcasts to all processes
 - **Alltoall**: all processes send data to all processes (including itself)
 - **Allreduce**: combines values from all processes and broadcasts the final result to all processes

MPI Allgather Operation



- Each process gathers data from all processes
- OR equivalently, we can say that all processes will perform the same exact MPI gather routine





MPI Allgather Routine

`MPI_Allgather (sendbuf, sendcnt, sendtype,
recvbuf, recvcnt, recvtype, comm)`

- Sendbuf: address of send buffer
- Sendcount: number of sent data elements
- Sendtype: type of data being sent
- Recvbuf: address of receive buffer
- Recvcount: number of received data elements of a single message
- Recvtype: type of data being received
- Comm: communicator object

© All Rights Reserved.

46



MPI Allgather Example

Assume number of processes is 4

`MPI_Allgather (sbuff, 1, MPI_INT, rbuff, 1, MPI_INT, MPI_COMM_WORLD)`

	P0	P1	P2	P3	
sbuff (before the call)	10	11	0	99	
rbuff (after the call)	P0	10	11	0	99
	P1	10	11	0	99
	P2	10	11	0	99
	P3	10	11	0	99

© All Rights Reserved.

47



MPI Alltoall Example

Assume number of processes is 4

`MPI_Alltoall (sbuf, 1, MPI_INT, rbuf, 1, MPI_INT, MPI_COMM_WORLD)`

sbuf (before the call)	P0	1	2	3	4
	P1	11	12	13	14
	P2	21	22	23	24
	P3	31	32	33	34

© All Rights Reserved.

50



MPI Alltoall Example (cont.)

Assume number of processes is 4

`MPI_Alltoall (sbuf, 1, MPI_INT, rbuf, 1, MPI_INT, MPI_COMM_WORLD)`

rbuf (after the call)	P0	1	11	21	31
	P1	2	12	22	32
	P2	3	13	23	33
	P3	4	14	24	34

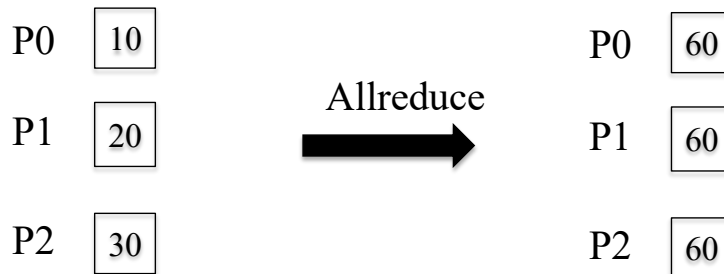
© All Rights Reserved.

51



MPI Allreduce Operation

- All processes will perform the same exact MPI reduction routine



MPI Allreduce Routine

`MPI_Allreduce (sendbuf, recvbuf, count, datatype, op, comm)`

- Sendbuf: address of send buffer
- Recvbuf: address of receive buffer
- Count: number of sent data elements
- Datatype: type of data being sent
- Op: reduce operation
- Comm: communicator object



MPI Allreduce Example

`MPI_Allreduce (sbuf, rbuf, 2, MPI_INT, MPI_SUM, MPI_COMM_WORLD)`

	P0	P1	P2	P3												
sbuf (before the call)	<table border="1"><tr><td>10</td><td>-1</td></tr></table>	10	-1	<table border="1"><tr><td>33</td><td>2</td></tr></table>	33	2	<table border="1"><tr><td>44</td><td>3</td></tr></table>	44	3	<table border="1"><tr><td>10</td><td>0</td></tr></table>	10	0				
10	-1															
33	2															
44	3															
10	0															
rbuf (after the call)	<table border="1"><tr><td>P0</td><td>97</td><td>4</td></tr></table>	P0	97	4	<table border="1"><tr><td>P1</td><td>97</td><td>4</td></tr></table>	P1	97	4	<table border="1"><tr><td>P2</td><td>97</td><td>4</td></tr></table>	P2	97	4	<table border="1"><tr><td>P3</td><td>97</td><td>4</td></tr></table>	P3	97	4
P0	97	4														
P1	97	4														
P2	97	4														
P3	97	4														

© All Rights Reserved.

54



Conclusions

- MPI is the dominant API for writing message-passing programs
- Programmers need to explicitly partition data and generate communication messages
- The main strategies for obtaining good performance with MPI: minimizing communication, maximizing concurrency, and ensuring load-balance
- When applicable, collective routines can generally obtain better performance than basic send/receive routines

© All Rights Reserved.

55

Parallel Sort in MPI

Instructor: Dr. Fahed Jubair
Computer Engineering Department
University of Jordan



Sort Algorithms



- Sort algorithms are typical targets for parallelization due to their common use in numerical applications
- We have already studied how to write the following three parallel sort algorithms in OpenMP
 - Odd-even transportation sort
 - Merge sort
 - Quick sort
- In this lecture, we will pick two of the above algorithms (odd-even transportation sort and merge sort) and study how to write MPI programs for them

Odd-Even Transportation Sort Serial Algorithm

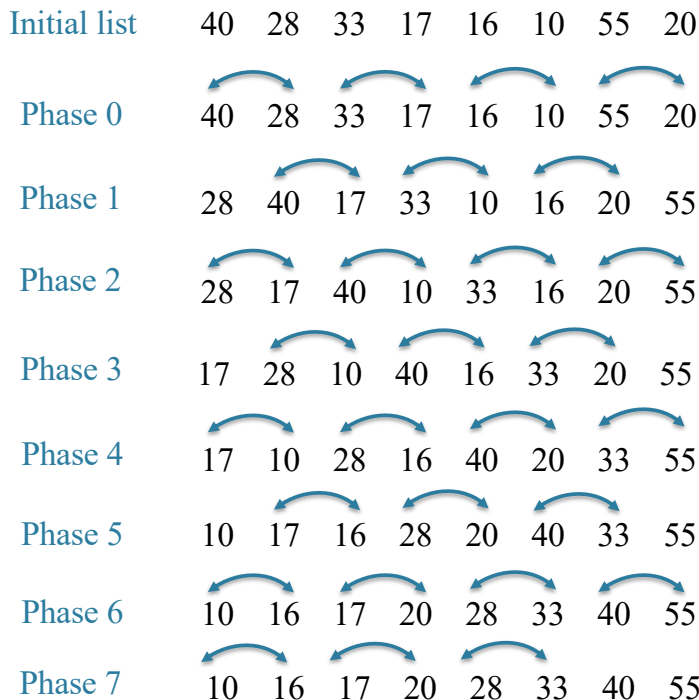


- As we previously learned in lecture 5, odd-even transportation sort works in **phases**, as follows:
 - During “even” phases, each odd-subscripted element, $a[i]$, is compared to the element to its left, $a[i - 1]$, and if they’re out of order, they’re swapped
 - During “odd” phases, each odd-subscripted element, $a[i]$, is compared to the element to its right, $a[i + 1]$, and if they’re out of order, they’re swapped
 - A theorem guarantees that after N phases, the list will be sorted

```
for(phase=0; phase<N; phase++){
    if(phase%2 == 0){
        for(j=1; j<N; j+=2)
            if(a[j-1] > a[j]){
                tmp = a[j-1];
                a[j-1] = a[j];
                a[j] = tmp;
            }
    }
    else {
        for(j=1; j<N-1; j+=2)
            if(a[j] > a[j+1]){
                tmp = a[j];
                a[j] = a[j+1];
                a[j+1] = tmp;
            }
    }
}
```

© All Rights Reserved.

Odd-Even Sort Example

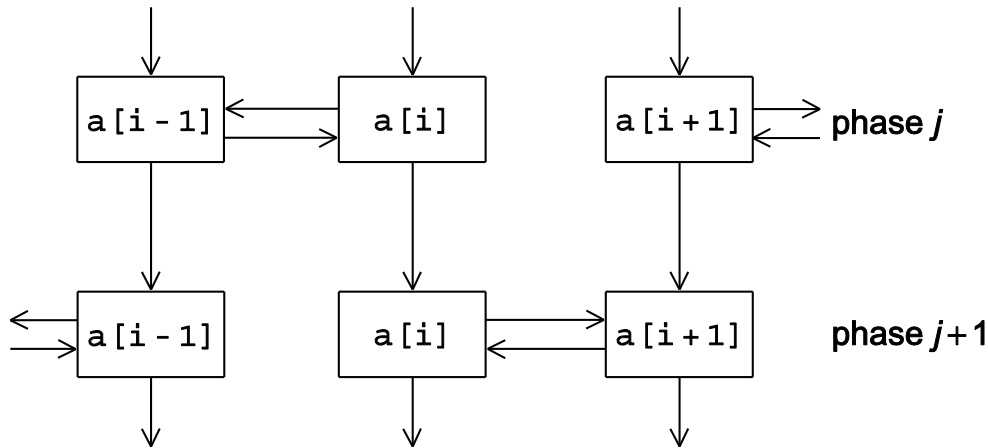


© All Rights Reserved.

Task Communication in Odd-Even Sort



In each phase, assuming each compare-swap represents a task



The task that's determining the value of $a[i]$ needs to communicate with either the task determining the value of $a[i-1]$ or $a[i+1]$

© All Rights Reserved.

Parallelizing Odd-Even Sort



- Assume n is the number of unsorted integers
- Assume p is the number of process
- If $n=p$, then the parallel algorithm is fairly straightforward
 - Depending on the phase, process i sends its current value, $a[i]$, to either process $i - 1$ or process $i + 1$
 - At the same time, process i should receive the value stored on process $i - 1$ or process $i + 1$, respectively, and then decide which of the two values it should store as $a[i]$ for the next phase
- However, the more likely scenario is that $n \gg p$
 - We need to modify the parallel odd-even sort algorithm

© All Rights Reserved.

Distributed-Memory Odd-Even Sort Assumptions



- Let us assume our algorithm will start and finish with n/p keys assigned to each process
 - Also, for simplicity, let us assume n is evenly divisible by p
- At the start, there are no restrictions on which keys are assigned to which processes
- When the algorithm terminates, the keys should be sorted in (say) increasing order such that if q and r are processes such that $q < r$, then every key in q is either less or equal to every key in r

© All Rights Reserved.

Distributed-Memory Odd-Even Sort Pseudo Code



```
Sort local keys;
for (phase = 0; phase < comm_sz; phase++) {
    partner = Compute_partner(phase, my_rank);
    if (I'm not idle) {
        Send my keys to partner;
        Receive keys from partner;
        if (my_rank < partner)
            Keep smaller keys;
        else
            Keep larger keys;
    }
}
```

Theorem: If the above parallel odd-even transposition sort is run with p processes, then after p phases, the input list will be sorted

© All Rights Reserved.

Parallel Odd-Even Sort Example



Time	Process			
	0	1	2	3
Start	15, 11, 9, 16	3, 14, 8, 7	4, 6, 12, 10	5, 2, 13, 1
After Local Sort	9, 11, 15, 16	3, 7, 8, 14	4, 6, 10, 12	1, 2, 5, 13
After Phase 0	3, 7, 8, 9	11, 14, 15, 16	1, 2, 4, 5	6, 10, 12, 13
After Phase 1	3, 7, 8, 9	1, 2, 4, 5	11, 14, 15, 16	6, 10, 12, 13
After Phase 2	1, 2, 3, 4	5, 7, 8, 9	6, 10, 11, 12	13, 14, 15, 16
After Phase 3	1, 2, 3, 4	5, 6, 7, 8	9, 10, 11, 12	13, 14, 15, 16

© All Rights Reserved.

C Implementation



```
st_time = MPI_Wtime();

qsort(my_keys, chunk, sizeof(int), Compare); // built-in sort function in C

for(phase=0; phase < comm_size; phase++){
    partner = compute_partner(phase, rank, comm_size);
    if(partner != MPI_PROC_NULL) { // not idle
        MPI_Request reqs[2];
        MPI_Status stats[2];
        MPI_Isend(my_keys, chunk, MPI_INT, partner, tag, MPI_COMM_WORLD, &reqs[0]);
        MPI_Irecv(partner_keys, chunk, MPI_INT, partner, tag, MPI_COMM_WORLD, &reqs[1]);
        MPI_Waitall(2, reqs, stats);
        if(rank < partner)
            Merge_low(my_keys, partner_keys, chunk); // keep smaller keys
        else
            Merge_high(my_keys, partner_keys, chunk); // keep larger keys
    }
}

end_time = MPI_Wtime();
sort_time = end_time - st_time;
MPI_Reduce(&sort_time, &max_sort_time, 1, MPI_DOUBLE, MPI_MAX, 0, MPI_COMM_WORLD);
```

© All Rights Reserved.

Compute_partner Function



```
int compute_partner ( int phase, int myrank, int comm_size){
    int partner;
    if(phase % 2 == 0) { /* even phase */
        if(myrank % 2 != 0) /* odd rank */
            partner = myrank - 1;
        else /* even rank */
            partner = myrank + 1;
    }
    else { /* odd phase */
        if(myrank % 2 != 0) /* odd rank */
            partner = myrank + 1;
        else /* even rank */
            partner = myrank - 1;
    }

    if(partner == -1 || partner == comm_size)
        partner = MPI_PROC_NULL;

    return partner;
}
```

© All Rights Reserved.

Merge_low Function



```
void Merge_low(int *myKeys, int *partnerKeys, int num_of_keys) {
    int ai, bi, ci;
    int *temp = (int*) malloc(num_of_keys * sizeof(int));
    ai = 0;
    bi = 0;
    ci = 0;
    while (ci < num_of_keys) {
        if (myKeys[ai] <= partnerKeys[bi]) {
            temp[ci] = myKeys[ai];
            ci++; ai++;
        } else {
            temp[ci] = partnerKeys[bi];
            ci++; bi++;
        }
    }
    memcpy(myKeys, temp, num_of_keys * sizeof(int));
    free(temp);
}
```

© All Rights Reserved.

Merge_high Function



```
void Merge_high(int *myKeys, int *partnerKeys, int num_of_keys) {
    int ai, bi, ci;
    int *temp = (int*) malloc(num_of_keys * sizeof(int));
    ai = num_of_keys-1;
    bi = num_of_keys-1;
    ci = num_of_keys-1;
    while (ci >= 0) {
        if (myKeys[ai] >= partnerKeys[bi]) {
            temp[ci] = myKeys[ai];
            ci--; ai--;
        } else {
            temp[ci] = partnerKeys[bi];
            ci--; bi--;
        }
    }
    memcpy(myKeys, temp, num_of_keys * sizeof(int));
    free(temp);
}
```

© All Rights Reserved.

Run-times of Parallel Odd-Even Sort



Processes	Number of Keys (in thousands)				
	200	400	800	1600	3200
1	88	190	390	830	1800
2	43	91	190	410	860
4	22	46	96	200	430
8	12	24	51	110	220
16	7.5	14	29	60	130

(times are in milliseconds)

© All Rights Reserved.



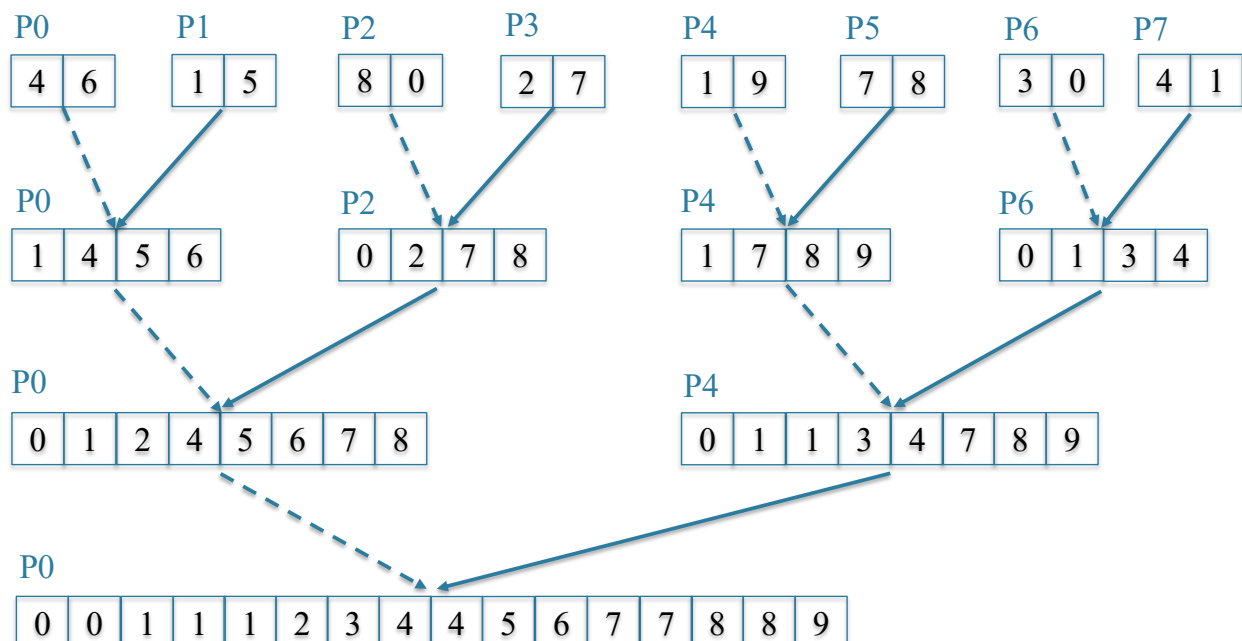
Parallel Merge Sort

- Several algorithms have been proposed to perform merge sort on distributed-memory platforms
- We will consider the following algorithm:
 - Initially, process 0 has all keys, which it distributes evenly to all processes
 - Each process locally sorts its list of keys
 - Processes communicate their lists of keys using a **tree-structured communication pattern** and merge them
 - Process 0 will have the final sorted list of keys

© All Rights Reserved.



Tree-Structured Merge Operations



© All Rights Reserved.

Distributed-Memory Odd-Even Sort Pseudo Code



```
Process 0 distributes keys evenly among processors ;
Sort local keys ;
steps = log2 (p) ;
for ( i = 1 ; i <= steps ; i++) {
    determine partner ;
    if (I'm not idle) {
        if ( I'am the receiver ) {
            receive partner keys ;
            merge partner keys with my keys ;
        }
        else { // I'am the sender
            send my keys to partner ;
        }
    }
}
```

© All Rights Reserved.

C Implementation



```
local_keys = (int*) malloc(chunk * sizeof(int));
MPI_Scatter(keys, chunk, MPI_INT, local_keys, chunk, MPI_INT, 0, MPI_COMM_WORLD);

st_time = MPI_Wtime();

qsort(local_keys, chunk, sizeof(int), Compare);

for(d=1; d<comm_size; d=d*2){
    if(rank % d == 0){ // if I'm not idle
        if(rank % (d*2) == 0){ // only true if I am the receiver
            int *rbuff, *mbuff;
            rbuff = (int*) malloc(d * chunk * sizeof(int));
            MPI_Recv(rbuff, d * chunk, MPI_INT, rank+d, tag, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
            mbuff = merge(rbuff, d * chunk, local_keys, d * chunk);
            free(local_keys); free(rbuff);
            local_keys = mbuff;
        }
        else { // I am the sender
            MPI_Send(local_keys, d * chunk, MPI_INT, rank-d, tag, MPI_COMM_WORLD);
        }
    }
}

end_time = MPI_Wtime();
sort_time = end_time - st_time;
MPI_Reduce(&sort_time, &max_sort_time, 1, MPI_DOUBLE, MPI_MAX, 0, MPI_COMM_WORLD);
```

© All Rights Reserved.



Merge Function

```
int* merge (int left[], int l_size, int right[], int r_size){
    int size = l_size + r_size;
    int* tmp = (int*) malloc(size * sizeof(int));
    int i, j, k;
    i = 0;
    j = 0;

    for(k=0; k<size; k++){
        if(i < l_size && j < r_size && left[i] <= right[j])
            tmp[k] = left[i++];
        else if(i < l_size && j < r_size && left[i] > right[j])
            tmp[k] = right[j++];
        else if(i < l_size)
            tmp[k] = left[i++];
        else
            tmp[k] = right[j++];
    }

    return tmp;
}
```

© All Rights Reserved.



Summary

- Shared-memory odd-even transportation and merge sort algorithms differ from their distributed-memory counterparts because their architectures emphasize different priorities
- Distributed-Memory sort algorithms employ “smart” communication patterns that reduce communication latency
 - In distributed-memory, minimizing communication is a priority
- The execution time of an MPI program is usually determined by measuring the execution times of all processes and then taking the longest time

© All Rights Reserved.

GPU Programming using CUDA

Instructor: Dr. Fahed Jubair
Computer Engineering Department
University of Jordan



Review: Parallel Computing Platforms



- **SIMD processors**
 - All processors execute the same instruction simultaneously, but while operating on different pieces of data
 - Work best with data-parallel applications
- **MIMD processors**
 - All processors execute in parallel but on different instructions and different pieces of data
 - Can handle both data-parallel and task-parallel applications

GPUs



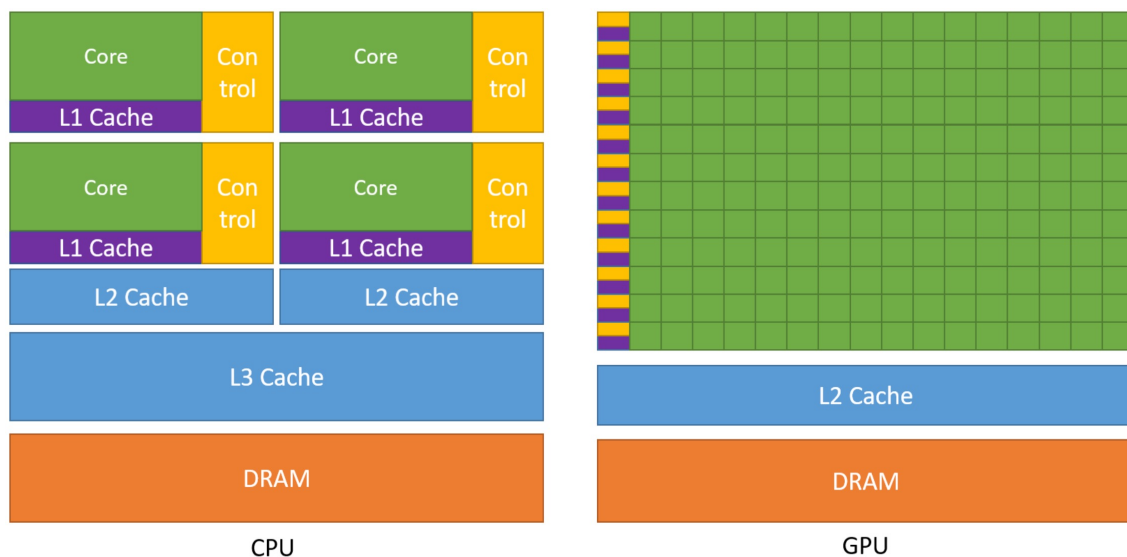
- In late 1990s, Graphical Processor Units (GPUs) have been developed to respond to the high demand of accelerating computer graphics applications
- However, many programmers have used GPUs for other applications, introducing what is known as General Purpose Computing on GPUs, or GPGPUs
- This lead to developing of several libraries and APIs that target GPGPUs, such as OpenCL and CUDA
- In this lesson, we will study CUDA C/C++ programming

© All Rights Reserved.

Why GPU Computing?



While GPUs have less programming flexibility, GPUs simply offer much higher computation throughput than CPUs



© All Rights Reserved.



What is CUDA?

- Introduced by Nvidia in 2006 for general purpose programming on Nvidia GPUs
- Supported by multiple programming languages, including C/C++
- CUDA Toolkit is available for download here:
<https://developer.nvidia.com/cuda-toolkit>
- CUDA C/C++ programming guide is available here:
<https://docs.nvidia.com/cuda/cuda-c-programming-guide/>

© All Rights Reserved.



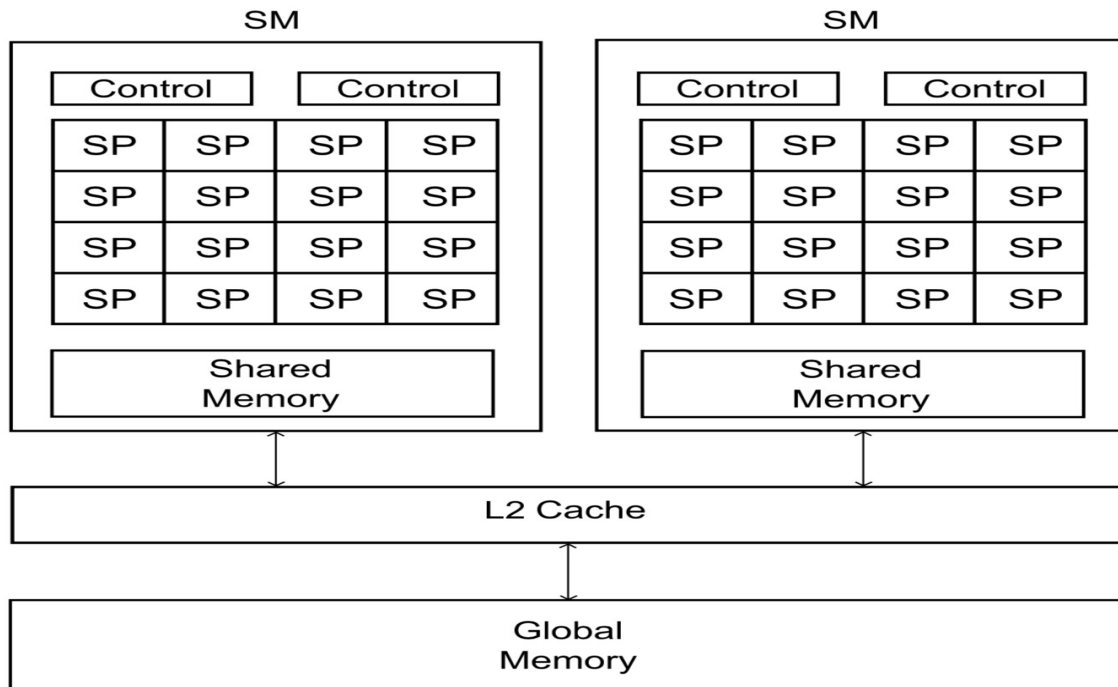
Terminology

- Nvidia GPUs are composed of multiple streaming multiprocessors, or SMs
- Each SM is composed of multiple streaming processors, or SPs
- An SP is analogous to a “core” in CPUs
- Each SP runs a “thread”
- Unlike conventional SIMD architectures, Nvidia SPs run asynchronously, i.e., threads on different SPs may execute in different speeds
- SIMT (single-instruction-multiple-thread) is a term used by Nvidia to describe their execution model

© All Rights Reserved.



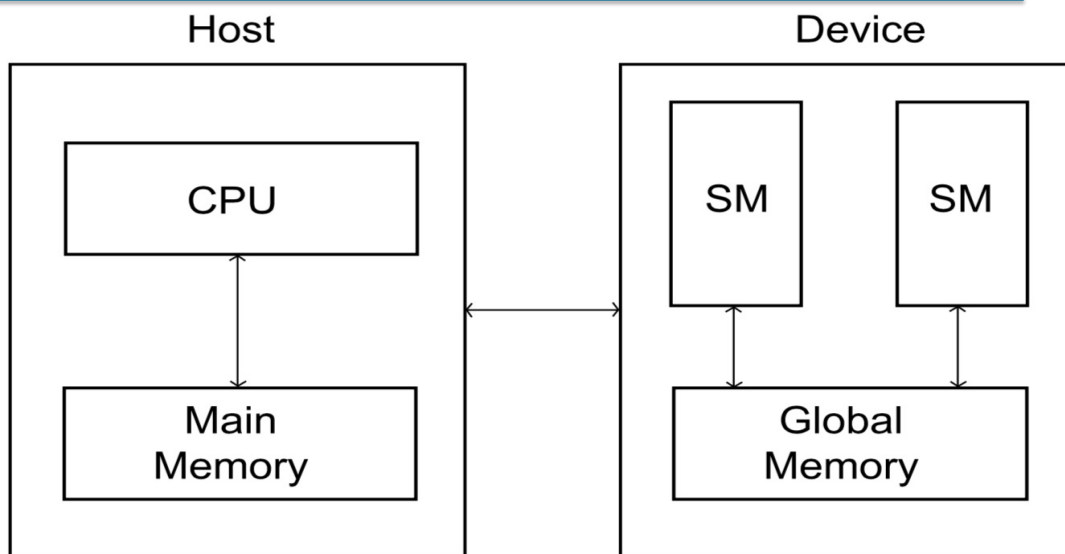
Nvidia GPU Architecture



© All Rights Reserved.



Heterogeneous Architecture



Terminology:

- **Host** The CPU and its memory (host memory)
- **Device** The GPU and its memory (device memory)

© All Rights Reserved.

Example: NVIDIA GeForce RTX 2060 SUPER



- It is the GPU that we have the lab in University of Jordan
- Released in July 2019
- Has around ten million transistors
- 34 SMs, each of which has 64 SPs (i.e., total is 2176 SPs)
- Base Clock 1470 MHz, Boost clock 1650 MHz
- 4MB L2 cache
- 8 GB memory with a bandwidth of 448 GB/s
- **Note:** as of spring 2021, one of the most powerful Nvidia processors has 128 SMs with a total of 10,496 SPs

© All Rights Reserved.

```
#include <stdio.h>
#include <cuda.h>  /* Header file for CUDA */

/* Device code: runs on GPU */
__global__ void Hello(void) {

    printf("Hello from thread %d!\n", threadIdx.x);
} /* Hello */

/* Host code: Runs on CPU */
int main(int argc, char* argv[]) {
    int thread_count; /* Number of threads to run on GPU */

    thread_count = strtol(argv[1], NULL, 10);
                    /* Get thread_count from command line */

    Hello <<<1, thread_count>>>();
                    /* Start thread_count threads on GPU, */

    cudaDeviceSynchronize(); /* Wait for GPU to finish */

    return 0;
} /* main */
```

CUDA
Hello Program



Compilation and Execution



- Compilation and execution commands

```
nvcc cuda_hello.cu -o cuda_hello
```

```
./cuda_hello 6
```

- Output:

```
Hello from thread 0!
```

```
Hello from thread 1!
```

```
Hello from thread 2!
```

```
Hello from thread 3!
```

```
Hello from thread 4!
```

```
Hello from thread 5!
```

© All Rights Reserved.

A Closer Look at the Hello Program



- Execution begins in the main by the host
- The below call to Hello starts the kernel from the device:

```
Hello <<< 1 , thread_count>>>());
```
- `threadIdx` is initialized by the system, and `threadIdx.x` returns the thread rank
- The call to the kernel from the host is asynchronous, i.e., the call returns immediately and the host resumes its execution
- `CudaDeviceSynchronize()` forces the host to wait until all threads in the device finish their execution

© All Rights Reserved.

Threads and Blocks



- CUDA organizes threads into blocks
- A thread block (or a block) is a collection of threads that run on a single SM
- Remember that each thread runs on a single SP
- When the kernel starts, each block is assigned to an SM, and threads inside this block run on that SM
- A grid is a collection of blocks

© All Rights Reserved.

```
#include <stdio.h>
#include <cuda.h>  /* Header file for CUDA */

/* Device code: runs on GPU */
__global__ void Hello(void) {

    printf("Hello from thread %d in block %d\n",
           threadIdx.x, blockIdx.x);
} /* Hello */

/* Host code: Runs on CPU */
int main(int argc, char* argv[]) {
    int blk_ct;          /* Number of thread blocks */
    int th_per_blk;     /* Number of threads in each block */

    blk_ct = strtol(argv[1], NULL, 10);
                    /* Get number of blocks from command line */
    th_per_blk = strtol(argv[2], NULL, 10);
                    /* Get number of threads per block from command line */

    Hello <<<blk_ct, th_per_blk>>>();
                    /* Start blk_ct*th_per_blk threads on GPU, */

    cudaDeviceSynchronize(); /* Wait for GPU to finish */

    return 0;
} /* main */
```

CUDA
Hello Program
Version 2



Compilation and Execution



- Compilation and execution commands

```
nvcc cuda_hello.cu -o cuda_hello
./cuda_hello 2 3
```

- Output:

```
Hello from thread 0 in block 0
Hello from thread 1 in block 0
Hello from thread 2 in block 0
Hello from thread 0 in block 1
Hello from thread 1 in block 1
Hello from thread 2 in block 1
```

© All Rights Reserved.

Grids



- In addition to threads and blocks, CUDA also allows grids
- A grid is a collection of blocks, and a block is a collection of threads
- Built-in variables:
 - `threadIdx`: the rank of a thread
 - `blockDim`: the dimensions of a thread block
 - `blockIdx`: the rank of a block
 - `gridDim`: the dimensions of a grid

All of these variables are structs with three integer fields: x, y, z

© All Rights Reserved.



Initializing Grids

```
dim3 grid_dims, block_dims;
grid_dims.x = 2;
grid_dims.y = 3;
grid_dims.z = 1;
block_dims.x = 4;
block_dims.y = 4;
block_dims.z = 4;
...
Kernel <<<grid_dims, block_dims>>> (...);
```



This kernel will have $2 \times 3 \times 1 = 6$ blocks, each of which has $4 \times 4 \times 4 = 64$ threads

© All Rights Reserved.



Vector Addition

- GPUs are designed to be effective with data-parallel applications
- Let us consider the CUDA version of performing simple addition of two n -integer vectors
- Below is the serial version of the code

```
for (i = 0; i < n; i++)
    c[i] = a[i] + b[i];
```

© All Rights Reserved.



Vector Addition Kernel

```
__global__ void vec_add(const int x[], const int y[], int sum[], int n)
{
    int my_element = blockDim.x * blockIdx.x + threadIdx.x ;

    if (my_element < n)
        sum[my_element] = x[my_element] + y[my_element] ;
}
```

- Each thread will perform the addition on exactly one element
- The index of the element is the same as the global rank of the thread (see next slide)

© All Rights Reserved.



Element Index

- `blockdim.x` returns the number of threads

Table 6.4 Global thread ranks or indexes in a grid with 4 blocks and 5 threads per block.

blockIdx.x	threadIdx.x				
	0	1	2	3	4
0	0	1	2	3	4
1	5	6	7	8	9
2	10	11	12	13	14
3	15	16	17	18	19

© All Rights Reserved.

```
int main(int argc, char* argv[]){
```

```
int *a, *b, *c;  
int i;
```

The Main Function

```
//allocate vectors on unified memory  
cudaMallocManaged(&a, n*sizeof(int));  
cudaMallocManaged(&b, n*sizeof(int));  
cudaMallocManaged(&c, n*sizeof(int));
```



Allocate a, b & c on
“unified memory” (see next slide)

```
// initialize a and b with random integers  
initialize_arrays (a, b);
```

```
/// perform vec_add on GPU  
vec_add <<<blk_cnt, thrd_per_blk>>> (a, b, c, n);
```

```
cudaDeviceSynchronize();
```

```
//print content of c  
print(c);
```

```
//free memory  
cudaFree(a);  
cudaFree(b);  
cudaFree(c);
```

```
}
```



Unified Memory



- Host and Device have separate physical memories
- Unified memory is a new addition to CUDA that allows programmers to develop programs as if the host and the device share a physical memory
- Underneath, the system takes care of data transfer between the host memory and the device memory
- Programmers may choose to allocate memory specifically on the host or the device and to do the data transfer between memories by hand
 - Useful for performing optimization

Rewriting Vector Addition with Explicit Memory Transfer



- No changes on the kernel

```
__global__ void vec_add(const int x[], const int y[], int sum[], int n)
{
    int my_element = blockDim.x * blockIdx.x + threadIdx.x ;

    if (my_element < n)
        sum[my_element] = x[my_element] + y[my_element] ;
}
```

© All Rights Reserved.

```
int main(int argc, char* argv){
    int *ha, *hb, *hc; // host vectors
    int *da, *db, *dc; // device vectors
    int i;

    //allocate vectors on host memory
    ha = (int*) malloc(n*sizeof(int));
    hb = (int*) malloc(n*sizeof(int));
    hc = (int*) malloc(n*sizeof(int));

    //allocate vectors on device memory
    cudaMalloc(&da, n*sizeof(int));
    cudaMalloc(&db, n*sizeof(int));
    cudaMalloc(&dc, n*sizeof(int));

    // initialize host vectors a and b with random integers
    initialize_arrays(ha, hb);

    // copy content from host to device
    cudaMemcpy(da, ha, n*sizeof(int), cudaMemcpyHostToDevice);
    cudaMemcpy(db, hb, n*sizeof(int), cudaMemcpyHostToDevice);

    /// perform vec_add on GPU
    vec_add <<<blk_cnt, thrd_per_blk>>> (da, db, dc, n);

    // copy content back from device to host
    cudaMemcpy(hc, dc, n*sizeof(int), cudaMemcpyDeviceToHost);

    //print content of host vector c
    print(hc);

    //free memory on host and device
    cudaFree(da);
    cudaFree(db);
    cudaFree(dc);
    free(ha);
    free(hb);
    free(hc);
}
```

The Main Function



Summary



- CUDA is designed to deliver high performance computing on Nvidia GPUs for data-parallel applications
- We covered basics concepts of CUDA
- We will learn more about CUDA in the lab

© All Rights Reserved.

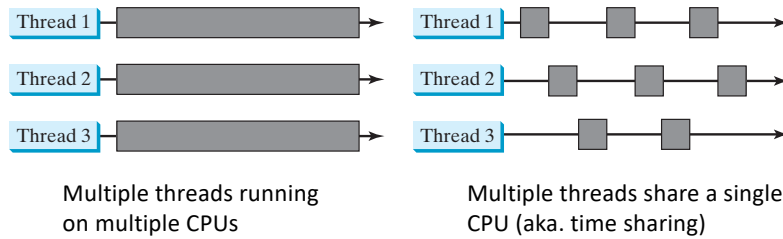
Multithreading Programming in Java

Instructor: Dr. Fahed Jubair
Computer Engineering Department
University of Jordan



Motivation

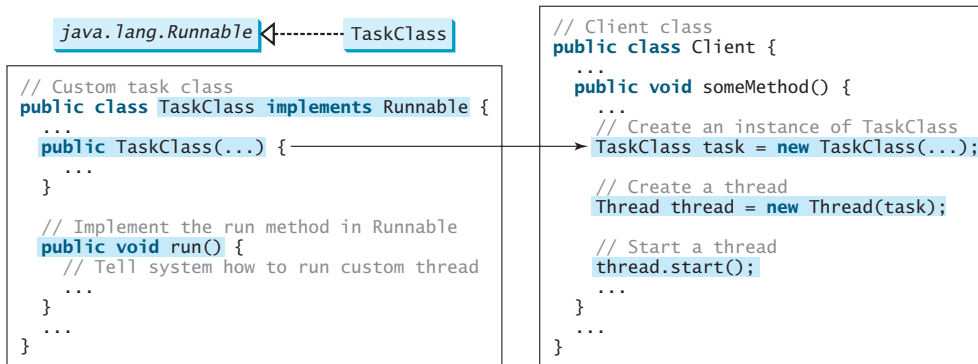
- Multithreading enables multiple tasks in a program to be executed concurrently
- A thread is the flow of execution, from beginning to end, of a specific task
- The operating system is responsible for scheduling and allocating resources to threads



Source for this lecture: Introduction to Java Programming, 10th edition

© All rights reserved.

Multithreading in Java



© All rights reserved.

Startup Example

```
public class TaskThreadDemo {
    public static void main(String[] args) {
        // Create tasks
        Runnable printA = new PrintChar('a', 100);
        Runnable printB = new PrintChar('b', 100);
        Runnable print100 = new PrintNum(100);

        // Create threads
        Thread thread1 = new Thread(printA);
        Thread thread2 = new Thread(printB);
        Thread thread3 = new Thread(print100);

        // Start threads
        thread1.start();
        thread2.start();
        thread3.start();
    }
}
```

```
class PrintChar implements Runnable {
    private char charToPrint; // The character to print
    private int times; // The times to repeat
    public PrintChar(char c, int t) {
        charToPrint = c;
        times = t;
    }

    @Override
    public void run() {
        for (int i = 0; i < times; i++) {
            System.out.print(charToPrint);
        }
    }
}
```

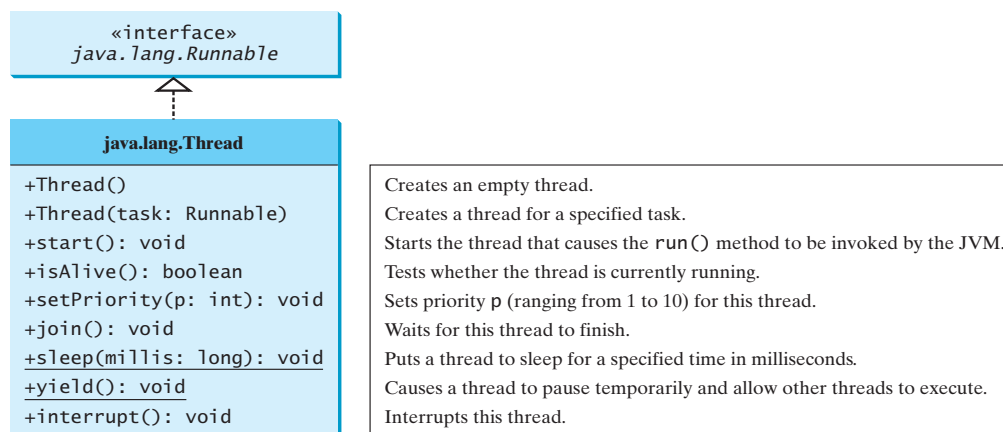
```
class PrintNum implements Runnable {
    private int lastNum;
    public PrintNum(int n) {
        lastNum = n;
    }

    @Override
    public void run() {
        for (int i = 1; i <= lastNum; i++) {
            System.out.print(" " + i);
        }
    }
}
```

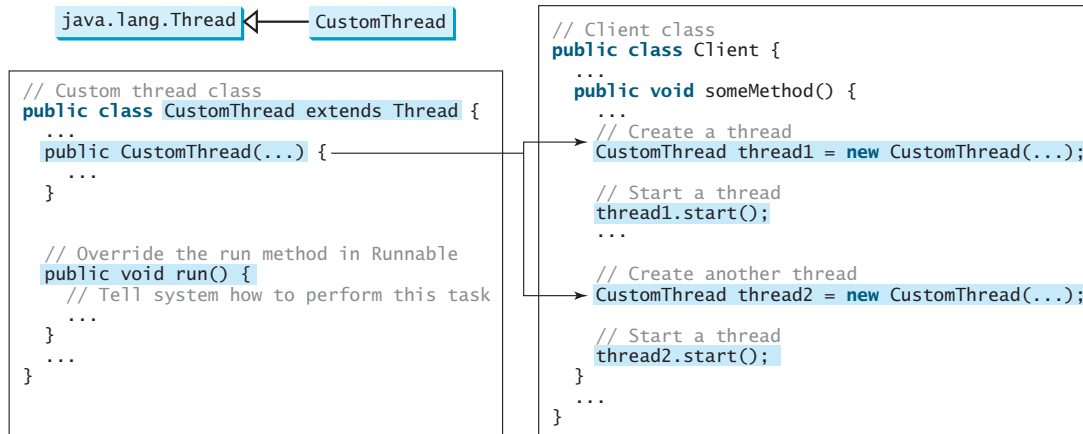
© All rights reserved.

The Thread Class

- The Thread class contains the constructors for creating threads for tasks and methods for controlling threads.



Alternative Method For Running Threads



© All rights reserved.

```

class myThread extends Thread {
    private String threadName;
    myThread( String name) {
        threadName = name;
        System.out.println("Creating " + threadName );
    }

    public void run() {
        System.out.println("Running " + threadName );
        try {
            for(int i = 4; i > 0; i--) {
                System.out.println("Thread: " + threadName + ", " + i);
                // Let the thread sleep for a while.
                Thread.sleep( millis: 50);
            }
        } catch (InterruptedException e) {
            System.out.println("Thread " + threadName + " interrupted.");
        }
        System.out.println("Thread " + threadName + " exiting.");
    }
}

public class TestThread {
    public static void main(String args[]) {
        myThread R1 = new myThread( name: "Thread-1");
        R1.start();
        myThread R2 = new myThread( name: "Thread-2");
        R2.start();
    }
}

```

Example Thread Test

```

// output
Creating Thread-1
Creating Thread-2
Running Thread-1
Thread: Thread-1, 4
Running Thread-2
Thread: Thread-2, 4
Thread: Thread-1, 3
Thread: Thread-2, 3
Thread: Thread-1, 2
Thread: Thread-2, 2
Thread: Thread-1, 1
Thread: Thread-2, 1
Thread Thread-1 exiting.
Thread Thread-2 exiting.

```

© All rights reserved.

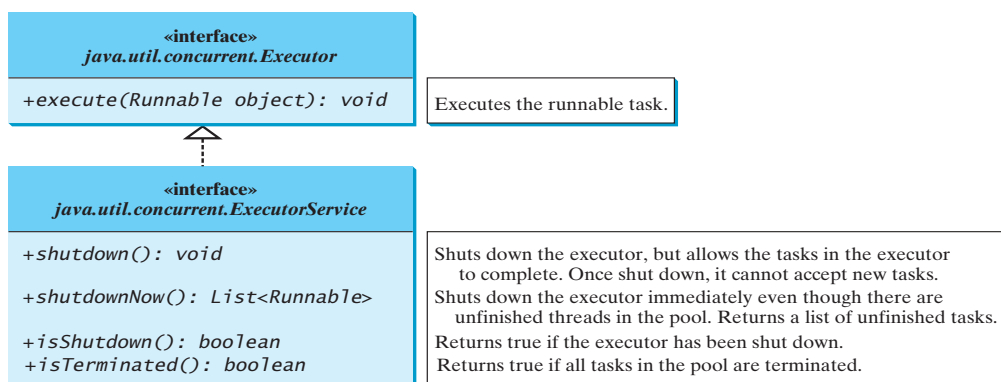
Thread Pool

- In many programs, many tasks need to be executed in parallel
- Creating a thread for each task is not efficient for a large number of tasks because the large number of threads will oversubscribe the system
- Java provides a mechanism for creating a pool of specific number of threads for executing large number of tasks
- Each thread starts by executing a task, when done, it grabs a task from the Task Queue, and so on

© All rights reserved.

The Executor and ExecutorService Interfaces

- Java provides the Executor interface for executing tasks in a thread pool and the ExecutorService interface for managing and controlling tasks
- ExecutorService is a subinterface of Executor



The Executors Class

- The Executors class has static methods useful for creating an Executor object
- There are two static methods for creating two different types of thread pools:
 - *newFixedThreadPool(int)*
 - Creates a fixed number of threads in a pool
 - If a thread completes executing a task, it can be reused to execute another task
 - If a thread terminates due to a failure prior to shutdown, a new thread will be created to replace it if all the threads in the pool are not idle and there are tasks waiting for execution
 - *newCachedThreadPool()*
 - Creates a new thread if all the threads in the pool are not idle and there are tasks waiting for execution
 - A thread in a cached pool will be terminated if it has not been used for 60 seconds
 - Efficient for many short tasks

© All rights reserved.

```
import java.util.concurrent.ExecutorService;
import java.util.concurrent.Executors;

public class ExecuterDemo
{
    public static void main(String[] args)
    {
        // pool creation
        ExecutorService executor = Executors.newFixedThreadPool( nThreads: 8);

        // tasks' execution
        for(int i=0; i<1000; i++)
            executor.execute(new Task(i));

        // pool shutdown
        executor.shutdown();
    }
}

class Task implements Runnable {
    private int taskNumber;
    public Task(int taskNumber) {
        this.taskNumber = taskNumber;
    }

    @Override
    public void run() {
        System.out.println("task " + taskNumber + " is executed.");
    }
}
```

Example Executor Demo

© All rights reserved.

Race Condition

- During concurrent execution, multiple threads may share data
- Reading and writing to shared data by multiple threads simultaneously may lead to inconsistent update of data (i.e., threads may incorrectly override each other)
- Such a situation when data is corrupted due to concurrent thread execution is called a **race condition**
- Race conditions cause incorrectness in concurrent execution

© All rights reserved.

Example: Race Condition Demo

```
import java.util.concurrent.*;

public class AccountWithoutSync {
    private static Account account = new Account();

    public static void main(String[] args) {
        ExecutorService executor = Executors.newCachedThreadPool();

        // Create and launch 100 threads
        for (int i = 0; i < 100; i++) {
            executor.execute(new AddAPennyTask());
        }

        executor.shutdown();

        // Wait until all tasks are finished
        while (!executor.isTerminated()) { }

        System.out.println("What is balance? " + account.getBalance());
    }

    // A thread for adding a penny to the account
    private static class AddAPennyTask implements Runnable {
        public void run() {
            account.deposit(1);
        }
    }
}

// An inner class for account
private static class Account {
    private int balance = 0;

    public int getBalance() {
        return balance;
    }

    public void deposit(int amount) {
        int newBalance = balance + amount;

        // This delay is deliberately added to magnify the
        // race condition problem and make it easy to see.
        try {
            Thread.sleep(5);
        } catch (InterruptedException ex) {
        }

        balance = newBalance;
    }
}

© All rights reserved.
```

Thread Synchronization and Critical Regions

- **Thread synchronization** is a mechanism for controlling access to shared data between multiple threads so that race conditions are avoided
- Example of a thread synchronization mechanism is a **critical region**, which is a region of code that allows one thread to execute this region at a time
- Critical regions are typically used whenever a specific region in the code needs to be executed serially

© All rights reserved.

The *Synchronized* Keyword

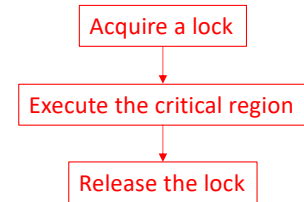
- In Java, the synchronized keyword is provided to specify a critical regions
- Using the synchronized keyword with a method will make this entire method executed as a critical region
 - Example: in the previous example, use “`public synchronized void deposit(int amount)`” to declare the deposit method as a critical region
- Using the synchronized keyword with statements make them executed inside a critical region
 - Example, in the previous example, write the below code to make a critical region

```
synchronized (account) {  
    account.deposit(1);  
}
```

© All rights reserved.

Synchronization Using Locks

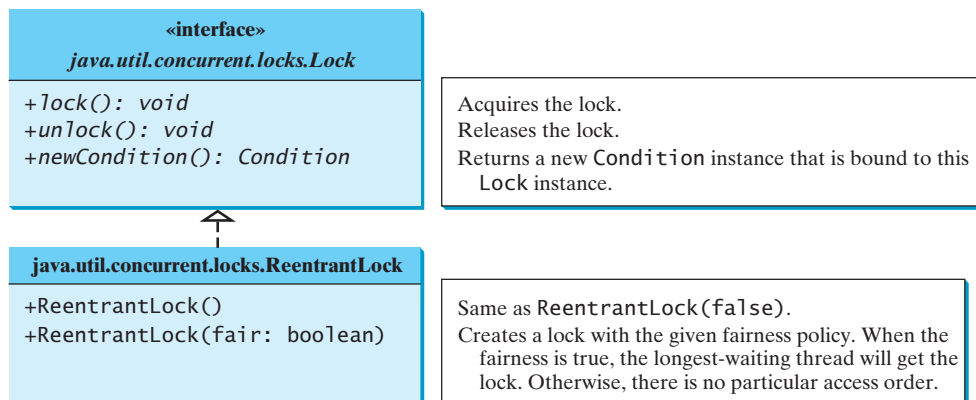
- Locks are mechanisms for implementing critical regions
- Only a thread that has the lock is allowed to enter and execute the critical region
- The *synchronized* keyword uses the lock/unlock mechanism for implementing a critical region
- However, Java also enables acquiring and releasing locks explicitly to allow users more thread control



© All rights reserved.

The Lock Interface and The ReentrantLock Class

- ReentrantLock is a concrete implementation of Lock for creating mutually exclusive locks



© All rights reserved.

Example: Lock/Unlock Demo

```
import java.util.concurrent.*;
import java.util.concurrent.locks.*;

public class AccountWithSyncUsingLock {
    private static Account account = new Account();

    public static void main(String[] args) {
        ExecutorService executor = Executors.newCachedThreadPool();

        // Create and launch 100 threads
        for (int i = 0; i < 100; i++) {
            executor.execute(new AddAPennyTask());
        }

        executor.shutdown();

        // Wait until all tasks are finished
        while (!executor.isTerminated()) {}

        System.out.println("What is balance ? " + account.getBalance());
    }

    // A thread for adding a penny to the account
    public static class AddAPennyTask implements Runnable {
        public void run() {
            account.deposit(1);
        }
    }
}
```

```
// An inner class for account
public static class Account {
    private static Lock lock = new ReentrantLock(); // Create a lock
    private int balance = 0;

    public int getBalance() {
        return balance;
    }

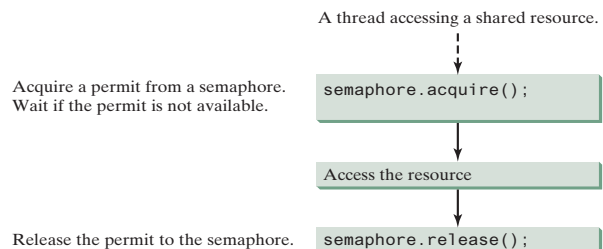
    public void deposit(int amount) {
        lock.lock(); // Acquire the lock

        try {
            int newBalance = balance + amount;
            Thread.sleep(5);
            balance = newBalance;
        }
        catch (InterruptedException ex) {}
        finally {
            lock.unlock(); // Release the lock
        }
    }
}
```

© All rights reserved.

Semaphores

- Semaphores are mechanisms that can control the number of threads that can enter a critical region
- Before entering a critical region, a thread must acquire a permit from a semaphore
- After finishing the critical region, a thread returns the permit to the semaphore object



© All rights reserved.

The Semaphore Class

<code>java.util.concurrent.Semaphore</code>	
<code>+Semaphore(numberOfPermits: int)</code>	Creates a semaphore with the specified number of permits. The fairness policy is false.
<code>+Semaphore(numberOfPermits: int, fair: boolean)</code>	Creates a semaphore with the specified number of permits and the fairness policy.
<code>+acquire(): void</code>	Acquires a permit from this semaphore. If no permit is available, the thread is blocked until one is available.
<code>+release(): void</code>	Releases a permit back to the semaphore.

© All rights reserved.

```
import java.util.concurrent.ExecutorService;
import java.util.concurrent.Executors;
import java.util.concurrent.Semaphore;

public class AccountWithSyncUsingSemaphores {
    private static Account account = new Account();

    public static void main(String[] args) {
        ExecutorService executor = Executors.newCachedThreadPool();

        // Create and launch 100 threads
        for (int i = 0; i < 100; i++) {
            executor.execute(new AddAPennyTask());
        }

        executor.shutdown();

        // Wait until all tasks are finished
        while (!executor.isTerminated()) {}

        System.out.println("What is balance ? " + account.getBalance());
    }

    // A thread for adding a penny to the account
    public static class AddAPennyTask implements Runnable {
        public void run() {
            account.deposit(1);
        }
    }
}
```

Example: Semaphore Demo

```
// An inner class for account
public static class Account {
    private static Semaphore semaphore = new Semaphore(1);
    private int balance = 0;

    public int getBalance() {
        return balance;
    }

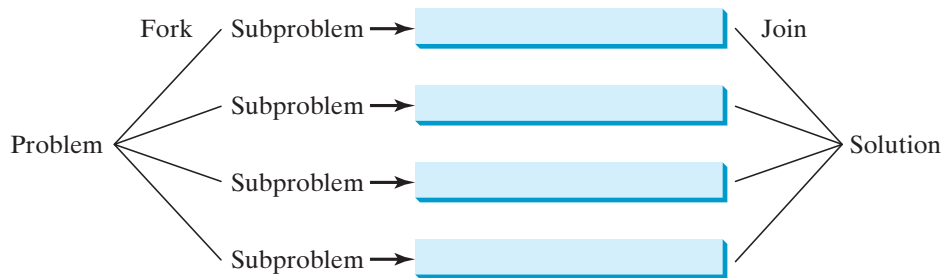
    public void deposit(int amount) {

        try {
            semaphore.acquire();
            int newBalance = balance + amount;
            Thread.sleep(5);
            balance = newBalance;
        }
        catch (InterruptedException ex) {
        }
        finally {
            semaphore.release();
        }
    }
}
```

© All rights reserved.

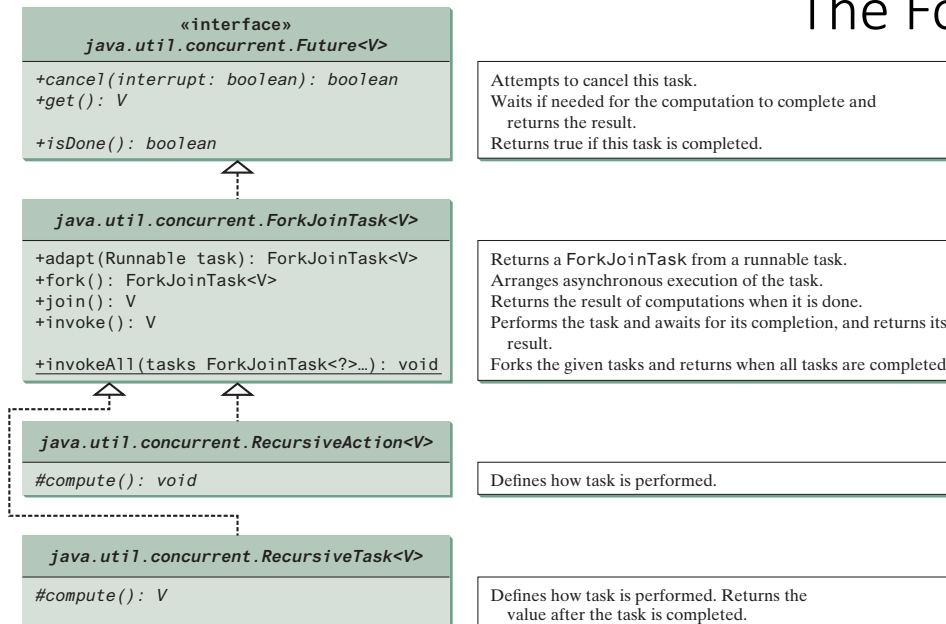
The Fork/Join Model

- In Java, the fork/join model is used for parallel programming
- This is achieved by defining the ForkJoinTask and ForkJoinPool classes



© All rights reserved.

The ForkJoinTask Class



© All rights reserved.

The ForkJoinPool Class

«interface»
`java.util.concurrent.ExecutorService`

`java.util.concurrent.ForkJoinPool`

+`ForkJoinPool()`
+`ForkJoinPool(parallelism: int)`
+`invoke(ForkJoinTask<T>): T`

Creates a `ForkJoinPool` with all available processors.
Creates a `ForkJoinPool` with the specified number of processors.
Performs the task and returns its result upon completion.

© All rights reserved.

```
import java.util.concurrent.*;

public class ParallelMax {
    public static void main(String[] args) {
        // Create a list
        final int N = 9000000;
        int[] list = new int[N];
        for (int i = 0; i < list.length; i++)
            list[i] = i;

        long startTime = System.currentTimeMillis();
        System.out.println("\nThe maximal number is " + max(list));
        long endTime = System.currentTimeMillis();
        System.out.println("Number of processors is " +
            Runtime.getRuntime().availableProcessors());
        System.out.println("Time with " + (endTime - startTime)
            + " milliseconds");
    }

    public static int max(int[] list) {
        RecursiveTask<Integer> task = new MaxTask(list, 0, list.length);
        ForkJoinPool pool = new ForkJoinPool();
        return pool.invoke(task);
    }
}
```

```
private static class MaxTask extends RecursiveTask<Integer> {
    private final static int THRESHOLD = 1000;
    private int[] list;
    private int low, high;
    public MaxTask(int[] list, int low, int high) {
        this.list = list;
        this.low = low;
        this.high = high;
    }

    @Override
    public Integer compute() {
        if (high - low < THRESHOLD) {
            int max = list[0];
            for (int i = low; i < high; i++)
                if (list[i] > max)
                    max = list[i];
            return new Integer(max);
        }
        else {
            int mid = (low + high) / 2;
            RecursiveTask<Integer> left = new MaxTask(list, low, mid);
            RecursiveTask<Integer> right = new MaxTask(list, mid, high);

            right.fork();
            left.fork();
            return new Integer(Math.max(left.join().intValue(),
                right.join().intValue()));
        }
    }
}
```

Example Parallel Max

Exercise

Parallel Sum

- Write a parallel program in Java that computes the sum of all integers between 1 and N
- Parallelize by dividing the list of integers into two chunks, and then recursively compute the sum of each chunk, then merge the summation of all chunks
- Compare the execution time of serial and parallel execution when
 - N=8000
 - N=800000
 - N=8000000