

Introduction to Parallel Computing

Pangfeng Liu
National Taiwan University

February 14, 2022

Why Parallel Computing?

- Solve a problem faster.
- Solve a problem better.

Olympic Games

The modern Olympic Games are the leading international sporting event featuring summer and winter sports competitions in which thousands of athletes from around the world participate in a variety of competitions.¹

¹http://en.wikipedia.org/wiki/Olympic_Games

Olympic Motto

The Olympic motto, Citius, Altius, Fortius, a Latin expression meaning “Faster, Higher, Stronger” was proposed by Pierre de Coubertin in 1894 and has been official since 1924.

Faster

- The reigning 100 m Olympic champion is often named “the fastest man/woman in the world”.
- The world record is 9.58 seconds.

Higher

- The high jump is a track and field event in which competitors must jump over a horizontal bar placed at measured heights without the aid of certain devices.
- The world record is 2.45m.

Stronger

- Olympic-style weightlifting is an athletic discipline in the modern Olympic program in which the athlete attempts a maximum-weight single lift of a barbell loaded with weight plates.
- The world record is 305 kg, the sum of snatch and clean & jerk.

Discussion

- Why faster, higher, and stronger?

Why Faster?

- Many computations are *slow*.
- Many computations are *time critical*.

Slow Computation

- Brute force search
- Computer simulation
- Exponential time complexity
- Grand Challenge problems

Traveling Salesman

Given a set of cities and the distance between two cities, find the shortest route that goes through all cities without visiting any city twice.

- A famous NP-complete problem, i.e., a proven hard-as-hell problem in computer science.
- Easy, you permute all cities and find the shortest route.
- The number of permutations is $(n - 1)!$ where n is the number of cities.

Factorial

- Consider a traveling salesman problem with 101 cities.
- $100!$ is roughly $9.332621544 \times 10^{157}$.
- Assume that the computation of the length of a path can be done in 100 cycles.
- Assume that a computer runs at 10 GHz.
- It takes $9.332621544 \times 10^{146}$ seconds to enumerate all paths.

Computation Time

- A year has 31,556,926 seconds, so the computation takes 2.9574×10^{130} billion years.
- The age of earth is 4.54 billion years².
- The age of the universe is 13.798 billion years³.
- Do you honestly think anyone can live to see the results?

²http://en.wikipedia.org/wiki/Age_of_the_Earth

³http://en.wikipedia.org/wiki/Age_of_the_universe

Algorithm Optimization

Despite techniques like branch-and-bound and A^* search can reduce the number of cases one needs to examine, the sheer number of permutations is enormous.

Discussion

- Give an example of computation that will take a lot of time.

Time Critical Computation

- Weather forecast
- Stock market
- Radar signal processing

Worried? Me?

- We do not need to do anything because computers become faster *automatically!*
- “Moore’s law” is the observation that, over the history of computing hardware, the number of transistors in dense integrated circuit doubles approximately every two years⁴.
- If the speed of a CPU is proportional to the number of transistors, we expect an eight times speed improvement in six years.

⁴http://en.wikipedia.org/wiki/Moore%27s_law

Worried? Yes.

- We assume we have an $O(n^3)$ “efficient” algorithm to solve a problem, e.g., matrix multiplication.
- You can solve problem size twice as large as the original one if you wait six years.
- By the time you wait for the hardware to catch up, your career (e.g., as a Ph.D. student) is over.

How to be Faster?

Hardware improvement cannot help you, OK!! It would be best if you had better solutions than waiting.

Late Information

- Late is *not* better than nothing.
- Justice delayed is justice denied.
- A weather forecast for tomorrow takes three days.
- A stock recommendation for next week takes three weeks.
- The computation of an early warning radar system takes three times for an enemy missile to destroy the radar.

How to be Faster?

Having more than one CPU to work on the problem seems to be a reasonable choice.

Discussion

- Give an example of time-critical computation.
- Describe Moore's Law.

Why Better?

- People are *greedy*.
- Video quality is an increasing function of time.
 - VCD (NTSC) 352×240 .
 - DVD (25 frame rate) 720×576 .
 - Blu-ray (HD) 1920×1080 .
- The number of pixels increases 24.5 times.
- An $O(n^2)$ algorithm will have to run 600 times faster.

Why Better?

- Again, remember that people are *greedy*.
- Remember faster, higher, stronger.
- We always seek possibility just beyond our capability.
- We want to live long and prosper (show a Vulcan gesture).

Live Long and Prosper



5

⁵http://vignette4.wikia.nocookie.net/memoryalpha/images/5/52/Spock_performing_Vulcan_salute.jpg/revision/latest?cb=20090320072701&path-prefix=en

The Weather

- The weather forecast resolution is proportional to the number of cells in the simulation.
- Taiwan has a length of 394 km and a width of 144 km.
- The area to simulate is 56736 square km.
- Assume that the simulation model is 10 km in height.
- If the resolution is one cubic km, we need 567360 cells.
- If we refine the resolution to 100 m, the number of cells will be 567360000.
- An $O(n^2)$ algorithms will have to run 1000000 times faster.

Discussion

- A naive matrix multiplication algorithm multiply two matrices of size n by n in $O(n^3)$ time. If we increase the size of the matrix n by a factor of 10, the execution time of this naive algorithm will roughly be?

Scientific Process

① Observation

- Why does the apple fall on my head?

② Theory

- Every matter attracts every matter.

③ Experiment

- Cavendish's experiment to verify the theory.

Newton and Apple



6

⁶<http://www.yalcafruittrees.com.au/wp-content/uploads/Isaac-Newton.png>

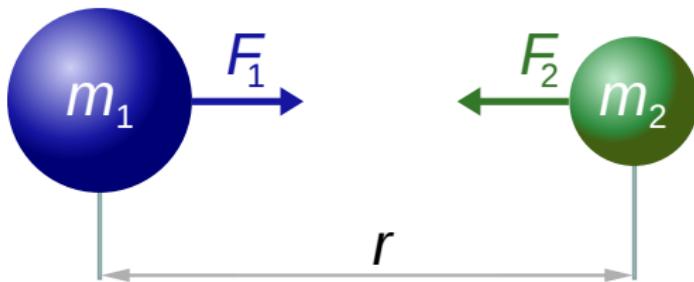


Newton's Law of Universal Gravitation

Newton's law of universal gravitation states that any two bodies in the universe attract each other with a force that is directly proportional to the product of their masses and inversely proportional to the square of the distance between them.⁷

⁷http://en.wikipedia.org/wiki/Newton%27s_law_of_universal_gravitation

Formulation

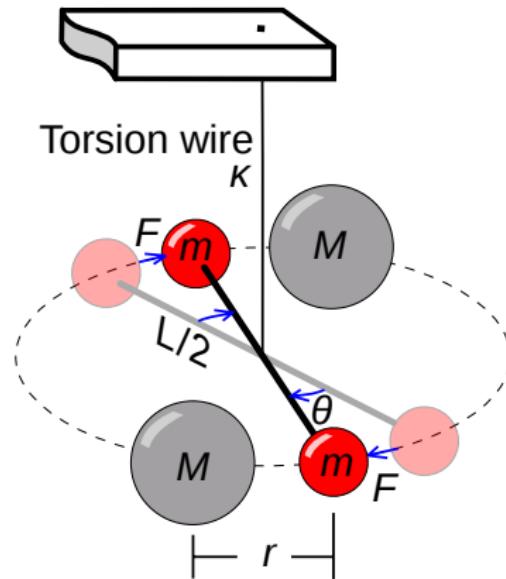


$$F_1 = F_2 = G \frac{m_1 \times m_2}{r^2}$$

8

⁸NewtonsLawOfUniversalGravitation by I, Dennis Nilsson. Licensed under CC BY 3.0 via Wikimedia Commons <http://upload.wikimedia.org/wikipedia/commons/thumb/0/0e/NewtonsLawOfUniversalGravitation.svg/200px-NewtonsLawOfUniversalGravitation.svg.png>

Cavendish's Experiment



9

⁹Cavendish Torsion Balance Diagram by Chris Burks (Chetvorno). Licensed under Public Domain via Wikimedia Commons
<http://commons.wikimedia.org/wiki/File:>

Scientific Process

① Observation

- Why does the galaxy have spirals?

② Theory

- Every matter attracts every matter.

③ Experiment

- What are you talking about?

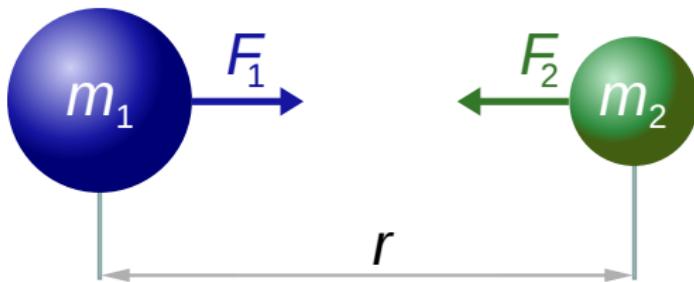
Spiral Galaxy



10

¹⁰[http://i.space.com/images/i/000/022/667/wS4/
spiral-galaxy-ngc1232-1600.jpg](http://i.space.com/images/i/000/022/667/wS4/spiral-galaxy-ngc1232-1600.jpg)

Formulation



$$F_1 = F_2 = G \frac{m_1 \times m_2}{r^2}$$

11

¹¹NewtonLawOfUniversalGravitation by I, Dennis Nilsson. Licensed under CC BY 3.0 via Wikimedia Commons <http://upload.wikimedia.org/wikipedia/commons/thumb/0/0e/NewtonLawOfUniversalGravitation.svg/200px-NewtonLawOfUniversalGravitation.svg.png>

Experiment

There are no ways we can put stars in a laboratory and observe the effects of gravity to incur spirals.

Discussion

- Give an example of scientific process, including observation, theory, and experiment.
- Give an example of experiment that cannot be conducted in a laboratory.

Difficulties

- Experiments are *expensive* – the stars may be expensive to buy.
- Experiments are *dangerous* – you do not want to have a black-hole in your laboratory.
- Experiments are *unfeasible* – my lab is not large enough.
- Experiments are *time consuming* – I do not have billions of years for observation.

Computer Simulation

- Simulation is cheap.
- Simulation is safe.
- Simulation is feasible.
- However, simulation is *slow*.

How to be Better?

Having more than one CPU to work on the problem seems to be a reasonable choice.

Discussion

- Give an example of scientific simulation.

Grand Challenge

A grand challenge is a fundamental problem in science or engineering, with broad applications, whose solution would be possible by the application of high-performance computing resources that could become available in the near future.¹²

¹²http://en.wikipedia.org/wiki/Grand_Challenge

Grand Challenge Examples

- Computational fluid dynamics
- Electronic structure calculations
- Plasma dynamics for fusion energy technology and for safe and efficient military technology
- Calculations to understand the fundamental nature of matter, including quantum chromodynamics and condensed matter theory
- Symbolic computations

How to be Better?

Having more than one CPU to work on the problem seems to be a reasonable choice.

Discussion

- Give an example of grand challenge problem.

Parallel Computing

Use multiple CPU's to solve a problem faster and/or better.

Why Parallel Computing?

- We want computation to be faster and better.
- The performance of a single computer is limited, so the only way is to have *more* computers.
- We cannot find single-core CPU anymore¹³, so it is essential to know how to get performance from a parallel computer.

¹³<http://www.intel.com/pressroom/kits/quickreffam.htm>

Limits

There are technological and physical limits to the uniprocessor performance that cannot be overcome. For example, clock times cannot be shorter than the response time of electronic circuits, which are limited by physical laws.¹⁴

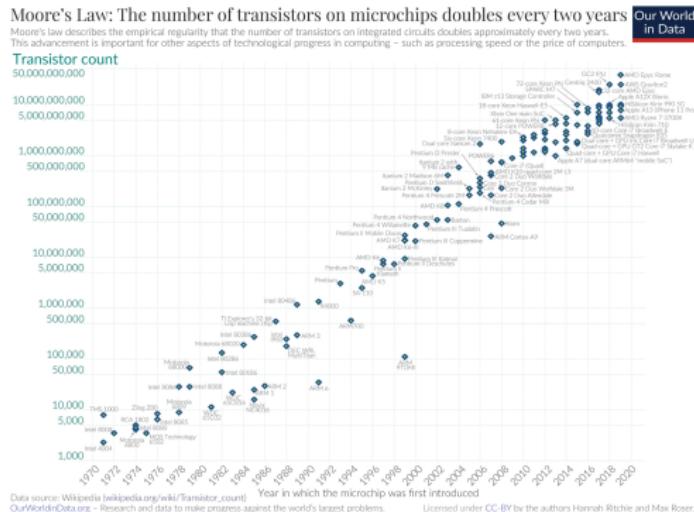
¹⁴Pangfeng Liu, The Parallel Implementation of N-body Algorithm, Ph.D. dissertation, 1994.

Moore's Law

"Moore's law" is the observation that, over the history of computing hardware, the number of transistors in a dense integrated circuit doubles approximately every two years.¹⁵

¹⁵http://en.wikipedia.org/wiki/Moore%27s_law

An Illustration



Fifty Years of Moore's Law

- Moore's Law is a direct consequence of the incredible and unique scaling heuristics of semiconductor manufacturing: by holding the cost per unit area of manufacturing constant¹⁶
- The economic benefits of Moore's Law come from shrinking the transistor.

¹⁶Chris A. Mack, IEEE Transactions on Semiconductor Manufacturing, Vol. 24, NO. 2, MAY 2011.

<http://ieeexplore.ieee.org/stamp/stamp.jsp?arnumber=5696765>

- Moore's Law is a learning curve by plotting minimum feature size as a function of cumulative revenue or area of silicon produced by the industry on a log-log scale.
- Moore's Law has kept on a relatively constant learning curve until about 1998—2000. The acceleration of this Moore's learning curve over the last decade is likely an unsustainable, momentum-driven attempt to recapture past revenue growth rates.
- The industry, and the world, has enjoyed 50 great years of Moore's Law. *There are unlikely to be many more years left.*

Discussion

- Give at least three data points to support the Moore's Law.

Parallel Computer

- A parallel computer is a system that has multiple processing units and supports parallel computing by *parallel programming*.
 - Multicore
 - Multiprocessor
 - Multicomputer

Multicore CPU

- A multicore CPU has multiple cores as processing units.
- The cores share the memory and have usually have their cache.
- A memory arbitrator guarantees the consistency of shared memory and cache.

Fujitsu A64FX

- The A64FX is a 64-bit ARM architecture microprocessor designed by Fujitsu.
- Fujitsu collaborated with ARM to develop the processor; it is the first processor to use the ARMv8.2-A Scalable Vector Extension SIMD instruction set with 512-bit vector implementation.
- Each A64FX processor has 4 NUMA nodes, with each NUMA node having 12 compute cores, for a total of 48 cores per processor.
- Fujitsu designed the A64FX for the Fugaku. More details later.

Discussion

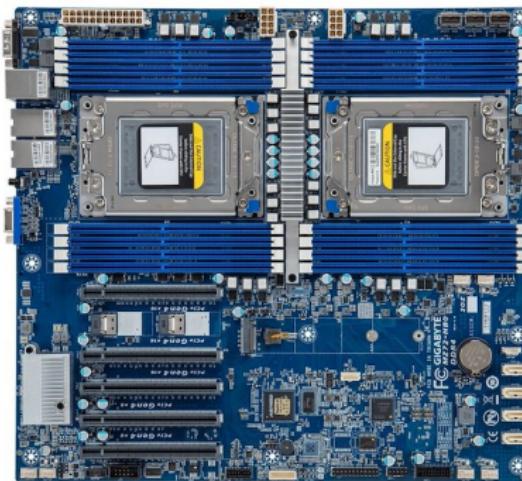
- Describe the number of cores in at least three current CPU's.

Multiprocessor

- A parallel system consists of *multiple* processors.
- Note that we usually do not distinguish processor and cores in this course, so we do not make a clear distinction between multiprocessor and multicore.
- The processors communicated with each other by reading and writing and shared memory, like a bulletin board.

Dual Socket Server

- This picture has a GigaByte Dual socket motherboard.¹⁷



¹⁷[https://www.anandtech.com/show/16734/
computex-2021-gigabyte-server-updates-mz72-hb0-for-dual-socket-3rd-gen](https://www.anandtech.com/show/16734/computex-2021-gigabyte-server-updates-mz72-hb0-for-dual-socket-3rd-gen)

Top500

- The TOP500 project ranks and details the 500 most powerful non-distributed computer systems in the world.
- The project was started in 1993 and publishes an updated list of the supercomputers twice a year.
- The project aims to provide a reliable basis for tracking and detecting trends in high-performance computing and bases rankings on HPL, a portable implementation of the high-performance LINPACK benchmark written in Fortran for distributed-memory computers.

November 2021

- Fugaku continues to hold the No. 1 position that it first earned in June 2020.
- The Microsoft Azure system called Voyager-EUS2 was the only machine to shake up the top spots, claiming No. 10. Based on an AMD EPYC processor with 48 cores and 2.45GHz working together with an NVIDIA A100 GPU and 80 GB of memory,
- <https://www.top500.org/lists/top500/2021/11/>.

Discussion

- Describe the number of processors within a node from any system in the top 10 of the top 500 list.

Multicomputer

- A parallel system consists of *multiple* computers.
- The computers are connected by a *communication network*.
- Since the computers are independent and do not share memory, they communicate with each other by *messages*, like making phone calls.

Multicomputer Examples

- Cluster computing
- Massively parallel computing
- Grid computing
- Cloud computing

Cluster

- A computer cluster consists of a set of loosely or tightly connected computers that work together so that, in many respects, they can be viewed as a single system.¹⁸
- Usually the computers are *loosely connected*. i.e., they are *not* connected by fast and expensive network.
- An economical alternative to those who cannot afford expensive parallel computers.

¹⁸http://en.wikipedia.org/wiki/Computer_cluster

My Definition

- In my humble opinion, any computer system connected by a network, but not a shared memory, is a cluster.
- The point is that they do not have shared memory, so they can only communicate with the network.
- The network is not necessarily slow – some networks are extremely fast, so the term *loosely coupled* may not be true in all cases.

Supercomputer Fugaku

- The total number of nodes in Fugaku is 158,976.
- A single CPU makes up a node, and two CPUs (two nodes) are mounted on a board called the CPU Memory Unit (CMU).
- Eight CMUs make up a Bunch of Blades (BoB), which means each BoB has 16 nodes.
- Three BoBs make up a Shelf, and therefore each Shelf has 48 nodes.
- Eight Shelves (384 nodes) are installed in a computer rack.
- Fugaku is made up of 432 racks, of which 396 racks have 384 nodes, and 36 racks have 192 nodes.
- <https://www.r-ccs.riken.jp/en/fugaku/project>

Supercomputer Fugaku

- It is at RIKEN Center for Computational Science, Japan.
- The total number of cores is 7,630,848.
- The amount of memory is 5,087,232 GB.
- Linpack Performance (Rmax) is 442,010 TFlop/s.

Hierarchy

- The total number of nodes in Fugaku is 158,976.
- A single CPU makes up a node, and two CPUs (two nodes) are mounted on a board called the CPU Memory Unit (CMU).
- Eight CMUs make up a Bunch of Blades (BoB), which means each BoB has 16 nodes.
- Three BoBs make up a Shelf, and therefore each Shelf has 48 nodes.
- Eight Shelves (384 nodes) are installed in a computer rack.
- Fugaku is made up of 432 racks, of which 396 racks have 384 nodes, and 36 racks have 192 nodes.

Hierarchy

- The processor is an A64FX 48C running at 2.2GHz, which is a multicore processor.
- Fugaku is a multicomputer connected by Tofu interconnection network.
- Again, we do not intend to have a clear distinction between a multicore CPU and a multiprocessor.
- More details in the “Architecture” lecture.

Discussion

- Describe the total number of cores within a node from any system in the top 10 of the top 500 list.

Parallelism

- Instruction-level parallelism
- Data parallelism
- Task parallelism

Instruction-level Parallelism

- Instructions can be re-ordered and combined into groups which are then executed in parallel without changing the result of the program.

Data Parallelism

- Data parallelism is a form of parallelization of computing across multiple processors in parallel computing environments by distributing the *tasks* across different parallel computing nodes.
- Often in the form of loops.

Data Parallelism Example

Example 1: Data Parallelism

```
1 for (i = 0; i < 1000000; i++)  
2     array[i] = 0;
```

- All the assignment can be done in parallel.

Task Parallelism

- Task parallelism is a form of parallelization of computing across multiple processors in parallel computing environments by distributing the *tasks* across different parallel computing nodes.
- Also called functional parallelism. Often in the form of function calls.

Data Parallelism Example

Example 2: Task Parallelism

```
1 parallel do {  
2     {  
3         cook_dinner();  
4         eat_dinner();  
5     }  
6     {  
7         turn_on_radio();  
8         listen_to_music();  
9     }  
10    } /* end of parallel do */  
11    go_to_bed();
```

Discussion

- Give an example of data parallelism and task parallelism.

Dependency Graph

- Every node is a task.
- Every edge is a dependency. This dependency is related to data or synchronization.
- A task can only start when all tasks that precede it finish, i.e., the tasks proceed in topological sort order.

Dependency Graph

- We can eat dinner only after we cook it.
- We can listen to music only after turning on the radio.
- We can go to sleep only after dinner and music.
- We can listen to music while having dinner.

Wavefront

- Task parallelism must respect the dependency in the dependency graph.
- One can imagine that task parallelism is a series of wavefronts in the dependency graph.
- We want to *increase* the number of tasks per wavefront to *reduce* the number of the wavefront. That is, we want to increase the *parallelism* of our algorithm so that it takes less time to complete.

Discussion

- Draw the dependency graph of the previous task parallel example. You may add a *start* node and an *end* node to indicate the beginning and the end of execution. Also, point out the *wavefront* and *dependency* in your drawing.

Quantities

- P The number of processors
- N The number of tasks
- L The longest path in the dependency graph
- W The maximum number of tasks in a frontend
- T The execution time, assuming that it takes one unit of time to finish a task.

Execution Time

$$N \geq T \geq \max(L, \frac{N}{\min(P, W)}) \quad (1)$$

- Increase P .
- Increase W , i.e., parallelism.
- Decrease L .

Discussion

- Explain the inequality in the previous page.

Parallel Programming

- Program a parallel system to perform parallel computing.
- That is why we are here.
- more about this in the “Programming Model” lecture.

Automatic Parallel Programming?

- Can a program convert our sequential programs into parallel programs automatically?
- No, otherwise we will not be here.
- Computer scientists have been trying to invent *smart* compiler that automatically does it, but only have limited success.
- More details in the Programming Model lecture.

- It is very often that people talk about parallel and distributed computing.
- For example, the name of my laboratory is the Laboratory of Parallel and Distributed Computing.
- Many conferences and journals have both in their titles.
 - IEEE International Parallel and Distributed Processing Symposium
 - IEEE International Conference on Parallel and Distributed Systems
 - Journal of Parallel and Distributed Computing
 - IEEE Transactions on Parallel and Distributed Systems

Difference

Parallel Computing A set of processing unit to work on a job *at the same time*, i.e., the focus is that the computations is *temporarily* in parallel.

Distributed Computing A set of processing unit to work on a job *at different locations*, i.e., the focus is that the computations is *geographically* distributed.

Difference

- The focus of parallel processing is *performance*. We care very much about the speed that we can finish a job.
- The focus of distributed processing is *reliability*. We care very much about can we finish a job no matter what happens – network down, hardware failure, earthquake, Godzilla attacks, etc.

Race-car



19

¹⁹http://upload.wikimedia.org/wikipedia/commons/4/4a/Formel3_racing_car_amk.jpg

Tank



20

²⁰By Boevaya Mashina - Own work, CC BY-SA 4.0,
<https://commons.wikimedia.org/w/index.php?curid=43925342>

In this Course

- This course will focus on performance, so we mostly discuss parallel processing, and will briefly discuss distributed processing when necessary.
- Distributed computing has become increasingly popular because of cloud computing and big data processing.
- Nevertheless, the role of parallel processing is still crucial since the processing speed is still essential.
- The focus is to combine the speed of parallel processing and the reliability of distributed computing into a “parallel and distributed system”.

Armored Vehicle



21

²¹By Dino246 at English Wikipedia - Transferred from en.wikipedia to Commons., Public Domain,
<https://commons.wikimedia.org/w/index.php?curid=3423910>

Discussion

- Explain and describe the difference between parallel and distributed computing. Give examples to illustrate your points.