

Exam Assignments V01

chengchengguo 183090

1. Describe how **parallelism** differs from **concurrency**

concurrency: A system is **concurrent** if it can support two or more actions in progress at the same time. Concurrency means executing multiple tasks at the same time but not necessarily simultaneously.

parallelism: A system is **parallel** if it can support two or more actions executing simultaneously. Parallelism is a specific kind of concurrency where tasks are really executed simultaneously.¹

2. What is fork-join parallelism?

Fork/join parallelism is a style of parallel programming useful for exploiting the parallelism inherent in divide and conquer algorithms on shared memory multiprocessors.

The idea is quite simple: a larger task can be divided into smaller tasks whose solutions can then be combined. As long as the smaller tasks are independent, they can be executed in parallel.

¹ [Concurrency vs. Parallelism — A brief view | by Madhavan Nagarajan | Medium](#)

```
pi_numerical_integration_parallel.cpp X
pi_numerical_integration_parallel.cpp > main()
1  #include <iomanip>
2  #include <iostream>
3  #include <omp.h>
4  // g++ -fopenmp pi_numerical_integration_parallel.cpp -o piparallel
5  // OMP_NUM_THREADS=4 ./pi
6  using namespace std;
7
8  int main() {
9
10     int num_steps = 100000000; // amount of rectangles
11     double width = 1.0 / double(num_steps); // width of a rectangle
12     sequential part
13     for (int th = 4; th < 129; th++) {
14         double sum = 0.0; // for summing up all heights of rectangles
15
16         double start_time = omp_get_wtime(); // wall clock time in seconds
17         omp_set_num_threads(th);
18
19         #pragma omp parallel // parallel region starts here-----
20         {
21             int num_threads = omp_get_num_threads();
22             int thread_id = omp_get_thread_num();
23             double sum_local = 0.0;
24
25             for (int i = thread_id; i < num_steps; i += num_threads) {
26                 double x = (i + 0.5) * width; // midpoint
27                 sum_local = sum_local + (1.0 / (1.0 + x * x)); // add new height of a rectangle
28             }
29             #pragma omp atomic
30             sum += sum_local;
31         } // parallel region ends here -----
32
33         double pi = sum * 4 * width; // compute pi
34         double run_time = omp_get_wtime() - start_time;
35         sequential part
36         cout << th << ", " << num_steps << ", " << setprecision(17)
37             << pi << ", " << setprecision(6) << run_time << "\n";
38     }
39 }
```

3. Read Chapter 1 from Computer Systems: A Programmer's Perspective². Discuss one thing you find particularly interesting.
(google it to find more information)

I found the part **1.9,1 Concurrency and Parallelism** is most interesting, because we have been exposed to this concept in class. This section explains it in more detail. It introduces 3 levels of parallelism from highest to lowest.

1) Thread-Level Concurrency

This form of concurrency allows multiple users to interact with a system at the same time, such as when many people want to get pages from a single web server. It also allows a single user to

² book.dvi.cmu.edu

engage in multiple tasks concurrently, such as having a web browser in one window, a word processor in another, and streaming music playing at the same time.³

uniprocessor system.: allows a single user to engage in multiple tasks concurrently

multiprocessor system: a system consisting of multiple processors all under the control of a single operating system kernel

Multi-core processors have several CPUs (referred to as “cores”) integrated onto a single integrated-circuit chip. Each core has its own L1 and L2 caches, but sharing the higher levels of cache as well as the interface to main memory.

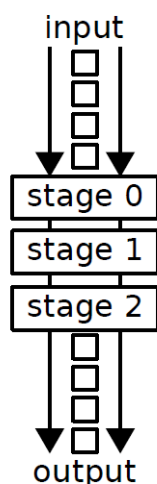
Hyperthreading, sometimes called *simultaneous multi-threading* is a technique that allows a single CPU to execute multiple flows of control. (As an example, the Intel Core i7 processor can have each core executing two threads, and so a four-core system can actually execute eight threads in parallel.)

The use of multiprocessing can **improve system performance** in two ways. **First**, it reduces the need to simulate concurrency when performing multiple tasks. **Second**, it can run a single application program faster, but only if that program is expressed in terms of multiple threads that can effectively execute in parallel. Thus, although the principles of concurrency have been formulated and studied for over 50 years, the advent of multi-core and hyperthreaded systems has greatly increased the desire to find ways to write application programs that can exploit the thread-level parallelism available with the hardware.

2) Instruction-Level Parallelism

At a much lower level of abstraction, modern processors can execute multiple instructions at one time, a property known as **instruction-level parallelism**.

it introduces the concept of **pipelining**, where the actions required to execute an instruction are partitioned into different **steps**, and the processor hardware is organized as a series of **stages**, each performing one of these steps. The stages can operate in parallel, working on different parts of different instructions.



- Idea: Split complex operation into several simpler (faster) stages
- Benefit: If we process the stages in parallel, we are able to increase the throughput
- Drawback: Pipeline has to be filled to be efficient start-up latency
- Pipelining applied to an instruction implements Instruction Level Parallelism (ILP)
- Pipelining is also applied “high-level”, e.g., for the parallelization of neural networks

³ [book.dvi \(cmu.edu\)](http://book.dvi.cmu.edu)

Example: Floating Point Multiplication:

$$c = a \bullet b$$

$$a = s_1 \bullet 0.m_1 \bullet 10^{e_1}$$

$$b = s_2 \bullet 0.m_2 \bullet 10^{e_2}$$

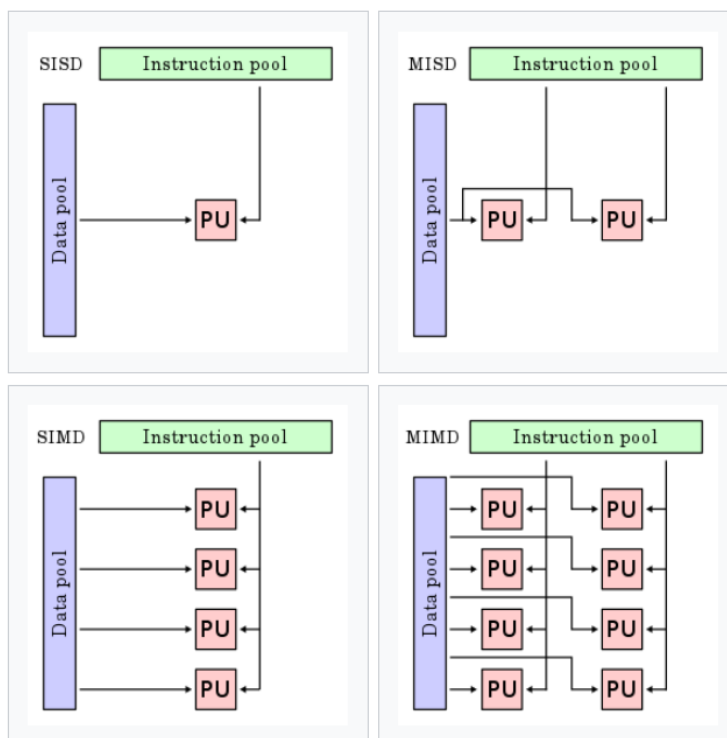
- Goal: Multiply to floating point values and given in a “normalized” representation
- Normalized:
- Sign bit (-1 or 1): s_1, s_2
 - Mantissa with non-zero leading digit: m_1, m_2
 - Exponent with positive or negative integer: e_1, e_2

3) Single-Instruction, Multiple-Data (SIMD) Parallelism

This is one of the Flynn classifications(below), it means: A single instruction is simultaneously applied to multiple different data streams. Instructions can be executed sequentially, such as by pipelining, or in parallel by multiple functional units. ⁴

Flynnsche Klassifikation

	Single Instruction	Multiple Instruction
Single Data	SISD	MISD
Multiple Data	SIMD	MIMD



⁴ [Flynn's taxonomy - Wikipedia](#)

4. Read the paper *There's plenty of room at the Top*:⁵ What will drive computer performance after Moore's law?

Explain in detail the figure *Performance gains after Moore's law ends*. (on the first page)

Moore's law is the observation that the number of transistors in a dense integrated circuit doubles about every two years.

- room at the "Top.": opportunities for growth in computing performance by improvement of software, algorithms, and hardware architecture.
- room at the "Bottom.": means semiconductor miniaturization

Moore's Law predicts exponential growth, and clearly exponential growth of computer performance can't continue forever.

"Performance gains after Moore's law ends. In the post-Moore era, improvements in computing power will increasingly come from technologies at the "Top" of the computing stack, not from those at the "Bottom", reversing the historical trend." means:

through processor simplification, where a complex processing core is replaced with a simpler core that requires fewer transistors. The freed-up transistor budget can then be redeployed in other ways. —for example, by increasing the number of processor cores running in parallel, which can lead to large efficiency gains for problems that can exploit parallelism.

In the post-Moore era, performance improvements from software, algorithms, and hardware architecture will increasingly require concurrent changes across other levels of the stack.

efficient cache utilization : restructure software to more predictable access patterns

Vectorization: use vector units for SIMD work

⁵ [Science Journals — AAAS \(microsoft.com\)](#)