

Parallel Computing - Notes - v0.6.0-dev

260236

October 2024

Preface

Every theory section in these notes has been taken from the sources:

- Course slides. [2]

About:

 [GitHub repository](#)

These notes are an unofficial resource and shouldn't replace the course material or any other book on parallel computing. It is not made for commercial purposes. I've made the following notes to help me improve my knowledge and maybe it can be helpful for everyone.

As I have highlighted, a student should choose the teacher's material or a book on the topic. These notes can only be a helpful material.

Contents

1	PRAM	5
1.1	Prerequisites	5
1.2	Definition	5
1.3	How it works	6
1.3.1	Computation	6
1.3.2	PRAM Classification	6
1.3.3	Strengths of PRAM	7
1.3.4	How to compare PRAM models	7
1.4	MVM algorithm	9
1.5	SPMD sum	11
1.6	MM algorithm	15
1.7	PRAM variants and Lemmas	16
1.8	PRAM implementation	17
1.9	Amdahl's and Gustafson's Laws	19
2	Fundamentals of architecture	22
2.1	Introduction	22
2.1.1	Simplest processor	22
2.1.2	Superscalar processor	23
2.1.3	Single Instruction, Multiple Data (SIMD) processor	24
2.1.4	Multi-Core Processor	24
2.2	Accessing Memory	25
2.2.1	What is a memory?	25
2.2.2	How to reduce processor stalls	27
2.2.2.1	Cache	27
2.2.2.2	Multi-threading	27
3	Programming models	30
3.1	Implicit SPMD Program Compiler (ISPC)	30
3.2	Shared Address Space Model	34
4	Parallel Programming Models and pthreads	35
4.1	How to create parallel algorithms and programs	35
4.2	Analyze parallel algorithms	37
4.3	Technologies	40
4.4	Threads	43
4.4.1	Flynn's taxonomy	43
4.4.2	Definition	43
4.4.3	pthreads API	45
4.4.3.1	Creation	45
4.4.3.2	Termination	46
4.4.3.3	Joining	47
4.4.3.4	Detaching	48
4.4.3.5	Joining through Barriers	49
4.4.3.6	Mutexes	50
4.4.3.7	Condition variables	50

5	OpenMP v5.2	51
5.1	Introduction	51
5.2	Basic syntax	53
	Index	57

1 PRAM

1.1 Prerequisites

Before we introduce the PRAM model, we need to cover some useful topics.

- A **Machine Model** describes a “machine”. It gives a value to the operations on the machine. It is necessary because: it makes it easy to deal with algorithms; it achieves complexity bounds; it analyses maximum parallelism.
- A **Random Access Machine (RAM)** is a model of computation that describes an abstract machine in the general class of register machines. Some features are:
 - **Unbounded** number of local memory cells;
 - Each memory cell can hold an integer of **unbounded** size;
 - Instruction set includes simple operations, data operations, comparator, branches;
 - All operations take **unit time**;
 - The definition of **time complexity** is the number of instructions executed;
 - The definition of **space complexity** is the number of memory cells used.

1.2 Definition

Definition 1: PRAM

A **parallel random-access machine (parallel RAM or PRAM)** is a **shared-memory abstract machine**. As its name indicates, the PRAM is intended as the parallel-computing analogy to the random-access machine (RAM) (not to be confused with random-access memory). In the same way that the RAM is used by sequential-algorithm designers to model algorithmic performance (such as time complexity), the **PRAM is used by parallel-algorithm designers to model parallel algorithmic performance** (such as time complexity, where the number of processors assumed is typically also stated).

The PRAM model has many interesting features:

- **Unbounded collection of RAM processors** (P_0 , P_1 , and so on);
- Processors don't have tape;
- Each processor has **unbounded registers**;
- **Unbounded collection of share memory cells**;
- All **processors can access all memory cells in unit time**;
- All **communication via shared memory**.

1.3 How it works

1.3.1 Computation

A single **processor** of the PRAM, at each computation, is **composed of 5 phases** (carried out in parallel by all the processors):

1. **Reads a value from one of the cells** $X(1), \dots, X(N)$
2. Reads one of the shared memory cells $A(1), A(2), \dots$
3. Performs some internal computation
4. **May write into one of the output cells** $Y(1), Y(2), \dots$
5. May write into one of the shared memory cells $A(1), A(2), \dots$

1.3.2 PRAM Classification

During execution, a subset of processors may remain idle. Also, some processors can read from the same cell at the same time (not really a problem), but they could also try to write to the same cell at the same time (**write conflict**). For these reasons, PRAMs are classified according to their read/write capabilities (realistic and useful):

- **Exclusive Read (ER)**. All processors can simultaneously read from distinct memory locations.
- **Exclusive Write (EW)**. All processors can simultaneously write to distinct memory locations.
- **Concurrent Read (CR)**. All processors can simultaneously read from any memory location.
- **Concurrent Write (CW)**. All processors can write to any memory location.

❓ But what value is ultimately written?

It depends on the mode we choose:

- **Priority Concurrent Write**. Processors have priority based on which value is decided, the **highest priority is allowed to complete write**.
- **Common Concurrent Write**. All processors are allowed to complete write **if and only if all the value to be written are equal**. Any **algorithm** for this model has to **make sure that this condition is satisfied**. Otherwise, the **algorithm is illegal** and the **machine state will be undefined**.
- **Arbitrary/Random Concurrent Write**. One **randomly chosen processor** is allowed to complete write.

1.3.3 Strengths of PRAM

PRAM is attractive and important model for designers of parallel algorithms because:

- It is **natural**. The number of operations executed per one cycle on P processors is at most P (equal to P is the ideal case).
- It is **strong**. Any processor can read/write any shared memory cell in unit time.
- It is **simple**. It abstracts from any communication or synchronization overhead, which makes the complexity and correctness of PRAM algorithm easier.
- It can be used as a **benchmark**. If a problem has no feasible/efficient solution on PRAM, it has no feasible/efficient solution for any parallel machine.

1.3.4 How to compare PRAM models

Consider two generic PRAMs, models A and B . Model A is **computationally stronger** than model B ($A \geq B$) **if and only if any algorithm** written for model B will **run unchanged** on model A in the **same parallel time** and with the **same basic properties**.

However, there are some useful metrics that can be used to compare models:

- **Time to solve problem of input size n on one processor, using best sequential algorithm:**

$$T^*(n) \tag{1}$$

- **Time to solve problem of input size n on p processors:**

$$T_p(n) \tag{2}$$

- **Speedup on p processors:**

$$SU_p(n) = \frac{T^*(n)}{T_p(n)} \tag{3}$$

- **Efficiency**, which is the work done by a processor to solve a problem of input size n divided by the work done by p processors:

$$E_p(n) = \frac{T_1(n)}{pT_p(n)} \tag{4}$$

- **Shortest run time** on any process p :

$$T_\infty(n) \tag{5}$$

- **Cost**, equal to processors and time:

$$C(n) = P(n) \cdot T(n) \quad (6)$$

- **Work**, equal to the total **number of operations**:

$$W(n) \quad (7)$$

Some properties on the metrics:

- The time to solve a problem of input n on a single processor using the best sequential algorithm *is not equal to* the time to solve a problem of input n in parallel using one of the p processors available. In other words, **the problem should not be solvable on a single processor on a parallel machine** (otherwise, what would be the point of using a parallel model?)

$$T^* \neq T_1$$

- $SU_P \leq P$
- $SU_P \leq \frac{T_1}{T_\infty}$
- $E_p \leq 1$
- $T_1 \geq T^* \geq T_p \geq T_\infty$
- $T^* \approx T_1 \Rightarrow E_p \approx \frac{T^*}{pT_p} = \frac{SU_p}{p}$
- $E_p = \frac{T_1}{pT_p} \leq \frac{T_1}{pT_\infty}$
- $T_1 \in O(C), T_p \in O\left(\frac{C}{p}\right)$
- $W \leq C$
- $p \approx \text{AREA} \quad W \approx \text{ENERGY} \quad \frac{W}{T_p} \approx \text{POWER}$

1.4 MVM algorithm

The **Matrix-Vector Multiply (MVM) algorithm** consists of four steps:

1. **Concurrent read of vector** X ($1 : n$) (transfer N elements);
2. **Simultaneous reads of different sections of the general matrix** A
(transfer $\frac{n^2}{p}$ elements to each processor);
3. **Compute** $\frac{n^2}{p}$ operations per processor;
4. **Simultaneous writes** (transfer $\frac{n}{p}$ elements from each processor).

Let i be the processor index, so the MVM algorithm is simply written as:

```

1 GLOBAL READ (Z ← X)
2 GLOBAL READ (B ← Ai)
3 COMPUTE (W := BZ)
4 GLOBAL WRITE (W → yi)

```

Algorithm 1: Matrix-Vector Multiply (MVM)

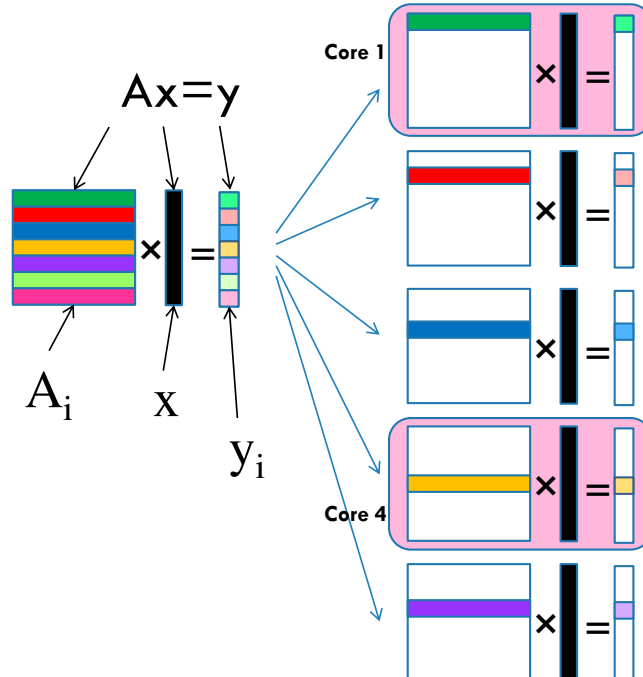


Figure 1: Example of MVM algorithm.

The performance of the MVM algorithm is as follows:

- The **time to solve** a problem of size n^2 is equal to the big O of the squared size of the problem as input divided by the number of processors available:

$$T_p(n^2) = O\left(\frac{n^2}{p}\right)$$

- The **cost** is equal to the number of processors and the time it takes to solve the problem. So it is quite trivial:

$$C = O\left(p \cdot \frac{n^2}{p}\right) = O(n^2)$$

- The **work** is equal to the cost, and the **linear power** P is equal to the ratio of work and time to solve the problem on p processors:

$$W = C \quad \frac{W}{T_p} = P$$

- The **perfect efficiency** is equal to:

$$E_p = \frac{T_1}{pT_p} = \frac{n^2}{p \frac{n^2}{p}} = 1$$

1.5 SPMD sum

The **Single Program Multiple Data (SPMD)** is a term that has been used to **describe computational models** for exploiting parallelism, where **multiple processors work together to execute a program to get results faster**.

In this section, we will see an SPMD approach on a Parallel Random Access Machine (PRAM). We will introduce one of the most common and simple mathematical operations: the sum.

The following pseudocode takes as **input an array** of size $n = 2^k$. In this case, n is a power of 2 because it ensures that the array can be evenly divided at each step of the computation. The value k is the number of iterations or levels of the summation process.

```

1 BEGIN
2   GLOBAL READ (A ← A(I))
3   GLOBAL WRITE (A → B(I))
4   FOR H = 1 : K
5     IF  $i \leq n \div 2^h$  THEN BEGIN
6       GLOBAL READ (X ← B(2i - 1))
7       GLOBAL READ (Y ← B(2i))
8       Z := X + Y
9       GLOBAL WRITE (Z → B(i))
10    END
11  IF I = 1 THEN
12    GLOBAL WRITE (Z → S)
13 END

```

Algorithm 2: Single Program Multiple Data (SPMD) sum

- First, read the entire input array A and copy the read data to another array B .
- Loop over h (1 to k). In each iteration, for each index i less than or equal to $n \div 2^h$, read values from array B at positions $2i - 1$ and $2i$; sum these values (and store the result in Z) and store the result (Z) back into $B(i)$.
- Once all iterations are complete, the final sum is stored in a variable S .

For example, if $n = 8$, then k would be 3, meaning that the algorithm will run for 3 iterations to sum all the elements in parallel.

h	i	adding
1	1	1,2
	2	3,4
	3	5,6
	4	7,8
2	1	1,2
	2	3,4
3	1	1,2

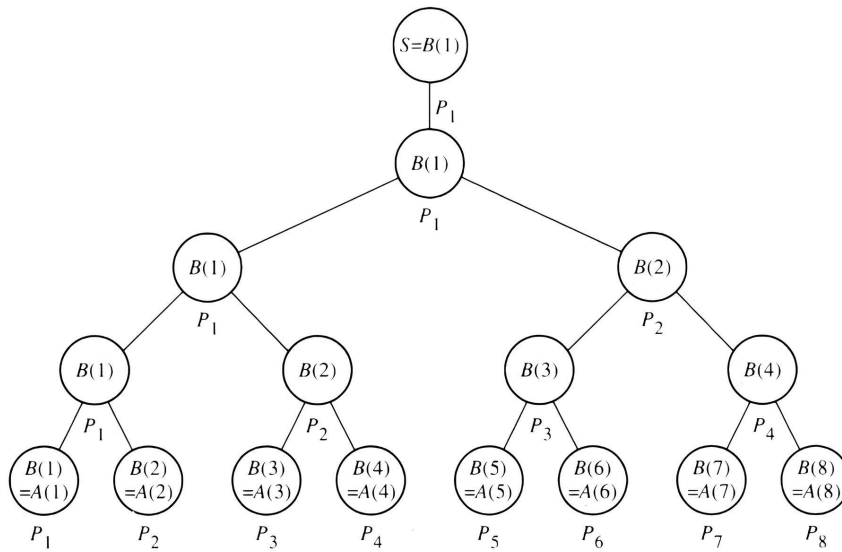
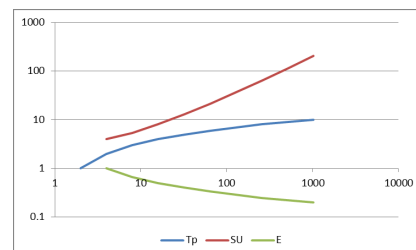
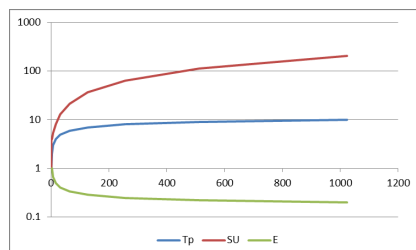


Figure 2: Computation of the sum of eight elements on a PRAM with eight processors. Each internal node represents a sum operation. The specific processor executing the operation is indicated below each node.

🕒 Performance of sum

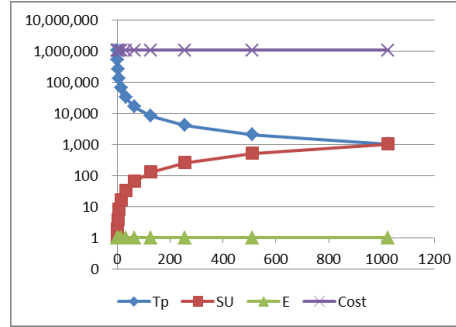
When the size of the array is equal to the number of processors ($N = P$), the **speedup and efficiency decrease**:

- $T^*(N) = T_1(N) = N$
- $T_{N=P}(N) = 2 + \log N$
- $\text{SU}_P = \frac{N}{2 + \log N}$
- $T^*(N) = P \cdot (2 + \log N) \approx N \log N$
- $E_p = \frac{T_1}{pT_p} = \frac{N}{N \log N} = \frac{1}{\log N}$



If the size of the array is much larger than the number of processors ($N \gg P$), the **speedup and power are linear**, the **cost is fixed** and the **efficiency is maximum (equal to 1)**:

- $T^*(N) = T_1(N) = N$
- $T_p(N) = \frac{N}{p} + \log p$
- $SU_P = \frac{N}{\frac{N}{p} + \log p} \approx P$
- $COST = p \left(\frac{N}{p} + \log p \right) \approx N$
- $WORK = N + P \approx N$
- $E_p = \frac{T_1}{pT_p} = \frac{N}{p \left(\frac{N}{p} + \log p \right)} \approx 1$



$n = 1'000'000$

Example 1

Refer to Figure 2 (page 12), the performance metrics are:

- $T_8 = 5$
- $C = 8 \cdot 5 = 40$ (could do 40 steps)
- $W = 2n = 16$ (16 on 40, wasted 24)
- $E_p = \frac{2}{\log n} = \frac{2}{3} = 0.67$
- $\frac{W}{C} = \frac{16}{40} = 0.4$

There is also the **Prefix Sum**, which takes **advantage of idle processors in the sum**. It computes all prefix sums:

$$S_i = \sum_{j=1}^i a_j \quad a_1, \quad a_1 + a_2, \quad a_1 + a_2 + a_3$$

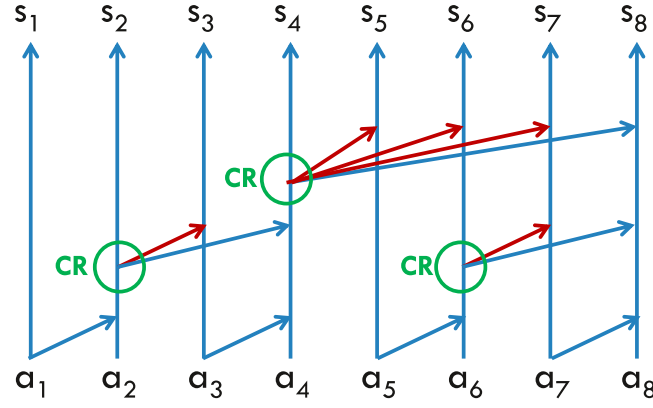


Figure 4: Prefix sum.

1.6 MM algorithm

The **Matrix Multiply (MM) algorithm** consists of three steps:

1. **Compute the two matrices** $A_{i,l}$ and $B_{l,j}$, so use the concurrent read.
2. Make the **sum**.
3. **Store** the result using exclusive write.

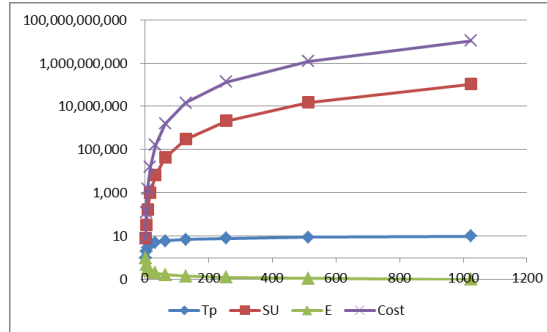
```

1 BEGIN
2    $T_{i,j,l} = A_{i,l}B_{l,j}$ 
3   FOR  $l = 1 : K$ 
4     IF  $l \leq n \div 2^h$  THEN
5        $T_{i,j,l} = T_{i,j,2l-1} + T_{i,j,2l}$ 
6   IF  $l = 1$  THEN
7      $C_{i,j} = T_{i,j,1}$ 
8 END
```

Algorithm 3: Matrix Multiply (MM)

🔗 Performance of MM

- $T_1 = n^3$
- $T_{p=n^3} = \log n$
- $SU = \frac{n^3}{\log n}$
- $\text{Cost} = n^3 \log n$
- $E_p = \frac{T_1}{pT_p} = \frac{1}{\log n}$



1.7 PRAM variants and Lemmas

The PRAM model presented here is one of the most commonly used. However, there are other important variants:

- PRAM model with a **limited number of shared memory cells** (small memory PRAM). If the input data set exceeds the capacity of the shared memory, the I/O values can be evenly distributed among the processors.
- PRAM model with **limited number of processors** (small PRAM). If the number of execution threads is higher, processors can interleave multiple threads.
- PRAM model with **limited size of one machine word**.
- PRAM model with **access conflicts handling**. These are restrictions on simultaneous access to shared memory cells.

Lemma 1. *Assume $P' < P$ and same size of shared memory. Any problem that can be solved for a P processor PRAM in T steps can be solved in a P' processor PRAM in:*

$$T' = O\left(\frac{TP}{P'}\right) \quad (8)$$

Proof. Partition P is simulated processors into P' groups of size $\frac{P}{P'}$ each. Associate each of the P' simulating processors with one of these groups. Each of the simulating processors simulates one step of its group of processors by:

- Executing all their read and local computation substeps first;
- Executing their write substeps then.

QED

Lemma 2. *Assume $M' < M$. Any problem that can be solved for a P processor and M -cell PRAM in T steps can be solved on a $\max(P, M')$ -processors M' -cell PRAM in $O\left(\frac{TM}{M'}\right)$ steps.*

Proof. Partition M simulated shared memory cells into M' continuous segments S_I , of size $\frac{M}{M'}$ each. Each simulating processor P'_I ($1 \leq I \leq P$), will simulate processor P_I of the original PRAM. Each simulating processor P'_I ($1 \leq I \leq M'$), stores the initial contents of S_I into its local memory and will use $M'[I]$ as an auxiliary memory cell for simulation of accesses to cell of S_I .

Simulation of one original read operation:

```

1 EACH  $P'_I$  ( $I = 1, \dots, \max(P, M')$ ) REPEATS FOR  $K = 1, \dots, \frac{M}{M'}$ 
2   WRITE THE VALUE OF THE  $K$ -TH CELL OF  $S_I$  INTO  $M'[I]$  ( $I = 1, \dots, M'$ )
3   READ THE VALUE WHICH THE SIMULATED PROCESSOR  $P_I$  ( $I = 1, \dots, P$ )
   WOULD READ IN THIS SIMULATED SUBSTEP, IF IT APPEARED IN THE
   SHARED MEMORY

```

The local computation substep of P_I ($I = 1, \dots, P$) is simulated in one step by P'_I . SImulation of one original write operation is analogous to that of read.

QED

1.8 PRAM implementation

The PRAM is an ideal model for creating parallel algorithms. Now we look at “*is it really implementable?*” The short answer is yes.

The longest answer is the following. There are already some examples of PRAM being converted to real machine models, such as [Explicit Multi-Threading \(XMT\)](#), Rigel, Tiler, etc. If conversion is not easy or possible, the implementation can be “*direct*”:

- The concurrent read is implemented as a detect-and-multicast technique.
- The concurrent write is implemented depending on the end result we want to achieve. Fetch-and-operate and prefix-sum are examples of serialized writing; otherwise, the CRCW technique is used:
 - Common CRCW: detect and merge
 - Priority CRCW: detect-and-priorities
 - Arbitrary CRCW: arbitrary

Example 2: Boolean DNF (sum of products) common CRCW

A logical formula is considered to be in DNF if it is a disjunction of one or more conjunctions of one or more literals.

Consider X as the sum of products of AND/OR operations:

$$X = a_1b_1 + a_2b_2 + \dots$$

The PRAM code, with X initialized to 0 and task index equal to $\$,$ is:

```
if ( $a_{\$}b_{\$}$ )  $X = 1$ ;
```

The common result is that not all processors write X and those that do write 1. The time complexity is $O(1)$. It works on common, priority and arbitrary CRCW.

Despite the previous example, exists also the PRAM SoP for the concurrent write. Let boolean X as:

$$X = a_1b_1 + a_2b_2 + \dots$$

The PRAM algorithm is:

```
if ( $a_i b_i$ )  $X = 1$ ;
```

Where all cores which write into X , **write the same value**.

✓ PRAM advantages

- Large body of algorithms.
- Easy to think about it.
- Sync version of shared memory. It eliminates sync and common issues, allows focus on algorithms, but allows adding these issues and allows conversion to async versions.
- Exists architectures for both synch (PRAM) and async (SM) model.
- PRAM algorithms can be mapped to other models.

1.9 Amdahl's and Gustafson's Laws

The **Amdahl's Law** is a formula which gives the **theoretical speedup in latency of the execution of a task at fixed workload that can be expected of a system whose resources are improved**. The law can be stated as:

Definition 2: Amdahl's Law

The overall performance improvement gained by optimizing a single part of a system is limited by the fraction of time that the improved part is actually used.

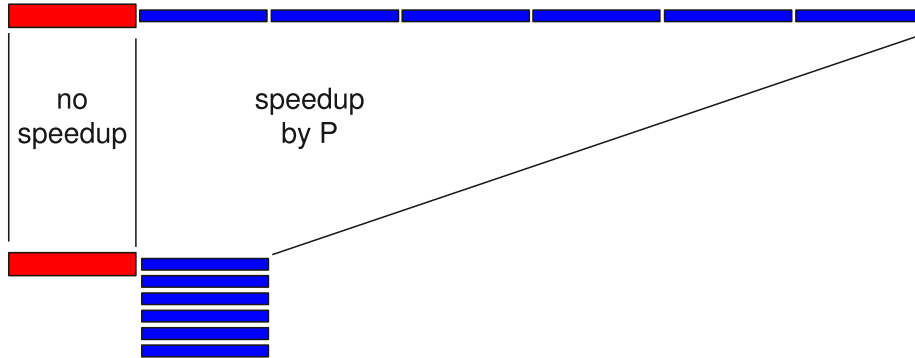
In practice, Amdahl's law says that the computation consists of interleaved segments of two types:

1. **Serial segments** (which cannot be parallelized);
2. **Parallelizable segments**.

Therefore, the metrics we can obtain are the time on P processors metric, that it is greater than the fraction of time on a processor divided by the processors P , and the speedup metric, that it is less than the number of processors P :

$$T_P > \frac{T_1}{P} \quad SU < P$$

Graphically, we can see a fixed part of the line, which is the **serial segment** (no speedup), and a set of instructions that can be **parallelized** (the sum of these segments is equal to the unit time 1).

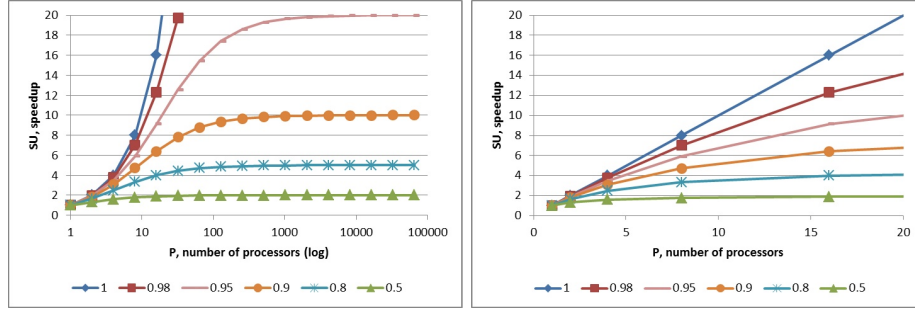


Furthermore, if we identify the parallelizable segment as f and the serial segment as $1 - f$, we obtain the following expressions:

$$SU(P, f) = \frac{T_1}{T_P} = \frac{T_1}{T_1 \cdot (1 - f) + \frac{T_1 \cdot f}{P}} = \frac{1}{(1 - f) + \frac{f}{P}} \quad (9)$$

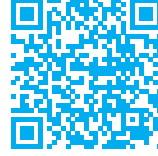
$$\lim_{P \rightarrow \infty} SU(P, f) = \frac{1}{1 - f}$$

In the following figure we can see the speedup with parameter f . Note the pessimism: for a problem with inherent $f = 90\%$, there is no point in using more than 10 processors.

Figure 5: Amdahl's law, $SU(P)$, parameter f .

The original paper presenting Amdahl's Law [1] can be viewed by clicking on the link below or by scanning the QR code.

[Amdahl's Law](#)



Amdahl's law applies only to the cases where the problem size is fixed. In practice, as more computing resources become available, they tend to get used on larger problems (larger datasets), and the time spent in the parallelizable part often grows much faster than the inherently serial work. In this case, **Gustafson's law gives a less pessimistic and more realistic assessment of the parallel performance.** [7]

Gustafson's Law gives the speedup in the execution time of a task that theoretically gains from parallel computing, using a hypothetical run of the task on a single-core machine as the baseline. To put it another way, it is the **theoretical "slowdown" of an already parallelized task if running on a serial machine.**

Against Amdahl's law, Gustafson suggests the following ideas:

- Portion f is not fixed;
- The absolute serial time is fixed;
- Parallel problem size is increased to exploit more processors;
- Fixed serial time (s of total) and fixed parallel time ($1 - s$ of total) are invariants;
- **Fixed time model** and not fixed size model (as Amdahl's law):

$$SU(P) = \frac{T_1}{T_P} = \frac{s + P \cdot (1 - s)}{s + (1 - s)} = s + P \cdot (1 - s) \quad (10)$$

Gustafson's law suggests a **linear speedup** and is **empirically applicable to highly parallel algorithms**.

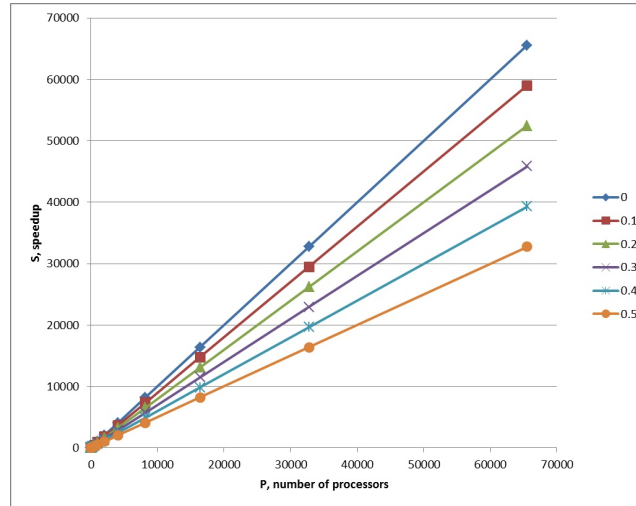


Figure 6: Gustafson's law.

Amdahl's Law states that as computing power increases, computational requirements remain the same. In other words, the **analysis of the same data will take less time with more computing power**.

Gustafson, on the other hand, argues that **more computing power leads to more careful and complete analysis of the data**. Where it would not have been possible or practical to simulate the impact of nuclear denotation on every building, car, and their contents (including furniture, structural strength, etc.) because such a calculation would have taken more time than was available to provide an answer, the increase in computing power will prompt researchers to add more data to more fully simulate more variables, giving a more accurate result.

The original paper presenting Gustafson's Law [3] can be viewed by clicking on the link below or by scanning the QR code.

Gustafson's Law



2 Fundamentals of architecture

The main purpose of this chapter is to introduce some basics of parallel computing theory. It will introduce the simplest and trivial processor and the more complex and efficient variants. The topics introduced are explained in a simple way and without any deepening, because it is only an introduction. For those who have studied computer science, this chapter might be a little boring and you might notice that some topics are explained in a simplistic way.

2.1 Introduction

2.1.1 Simplest processor

Inside a computer, a processor executes instructions.

- **Fetch/Decode:** Determine which instruction to run next;
- **ALU (execution unit):** Performs the operation described by an instruction, which may change values in the processor's registers or the computer's memory;
- **Registers:** maintain program state, store values of variables used as inputs and outputs to operations.

The simplest and most basic processor executes **one instruction per clock cycle**.

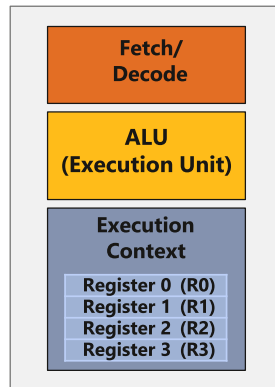


Figure 7: The simplest and most basic processor.

2.1.2 Superscalar processor

A more “complex” and realistic model is the **superscalar processor**. This **processor can decode and execute up to two instructions per clock**. The execution is slightly different from the simplest processor. The **processor automatically finds independent instructions in an instruction sequence and can execute them in parallel on multiple execution units**.

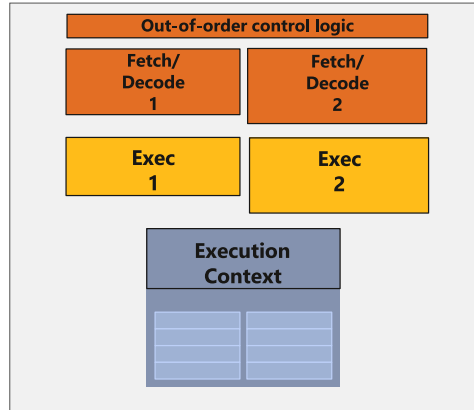


Figure 8: The superscalar processor.

The superscalar processor takes advantage of **Instruction-Level Parallelism (ILP)**¹ within an instruction stream.

- Processing **different instructions** from the same instruction stream **in parallel (within a core)**.
- **Parallelism is automatically detected by the hardware during execution**.

¹Instruction-level parallelism (ILP) is the parallel or simultaneous execution of a sequence of instructions in a computer program. More specifically ILP refers to the average number of instructions run per step of this parallel execution.

2.1.3 Single Instruction, Multiple Data (SIMD) processor

Adding execution units (ALUs) to the simplest processor can increase compute capability. Amortize the cost/complexity of managing an instruction stream across many ALUs using **Single Instruction, Multiple Data (SIMD)** processing. Therefore, the **same instruction is sent to all ALUs**. This **operation is performed in parallel on all ALUs**.

✓ Advantages

- **Efficient for data-parallel workloads:** amortize control costs over many ALUs.
- Vectorization done by:
 - Compiler (**explicit SIMD**): parallelism is explicitly requested by the programmer through intrinsics, conveyed through parallel language semantics, and inferred through loop dependency analysis by the “auto-vectorizing” compiler. In other words, the **SIMD parallelization is done at compile time, and when we inspect the program binary, we can see the SIMD instructions**.
 - At runtime by hardware (**implicit SIMD**): the **compiler generates a binary with scalar instructions**, but n instances of the program are always executed together on the processor. The **hardware** (not the compiler) is **responsible for the simultaneous execution of the same instruction by multiple program instances on different data on SIMD ALUs**.

2.1.4 Multi-Core Processor

A **Multi-Core Processor (MCP)** is a **microprocessor** on a single integrated circuit (IC) with **two or more separate central processing units (CPUs)**, called *cores* to emphasize their multiplicity (e.g., *dual-core* or *quad-core*). Each core reads and executes program instructions, specifically ordinary CPU instructions (such as **add**, **move data**, and **branch**). However, the MCP can **execute instructions on separate cores simultaneously, increasing overall speed for programs that support multithreading or other parallel computing techniques**.

✓ Advantages

- Provides **thread-level parallelism**: execute a completely different instruction stream on each core simultaneously.
- **Software creates threads to expose parallelism to hardware** (e.g., via threading API)

2.2 Accessing Memory

2.2.1 What is a memory?

A computer's memory is organized as an array of bytes. Each byte is identified by its address in memory (its position in that array).

Address	Value
0x0	16
0x1	255
0x2	14
0x3	0
0x4	0
0x5	0
0x6	6
0x7	0
0x8	32
0x9	48
0xA	255
0xB	255
0xC	255
0xD	0
0xE	0
0xF	0
0x10	128
⋮	⋮
0x1F	0

Table 1: Example illustration of the program's memory address space of 32 bytes, range from 0x0 to 0x1F.

From the processor's point of view, loading an instruction to access the contents present in memory is done with the `ld` assembly instruction. For example, `ld R0 ← mem[R2]` means “take the value from register R2 and put that value into register R0”.

Before we introduce new concepts, let us take a moment to explain some important **terminology**:

- **Memory Access Latency**, is the **time** it takes for the **memory system** to deliver data to the processor.
- **Processor Stall**. A processor stalls when it **cannot execute the next instruction in an instruction stream because of a dependency on a previous instruction that has not been completed**. Accessing memory is a major source of stalling, which is one of the main reasons why memory accesses should be limited.

For **example**, in the following three assembler instructions, the `add` has to wait for the loading of `R2` and `R3` values, making parallelization more complicated:

```
ld r0 mem[r2]
ld r1 mem[r3]
add r0, r0, r1
```

- **Memory Bandwidth**, is the rate at which the memory system can provide data to a processor.

Bandwidth is the critical resource in modern computing.

High-performance parallel programs will:

1. **Organize computation to fetch data from memory less frequently.** For example, reuse data previously loaded by the same thread (temporal locality optimizations) or share data across threads (inter-thread cooperation);
2. Prefer to **perform additional arithmetic to store/reload values**;
3. **Programs need to access memory infrequently** to take advantage of modern processors.

2.2.2 How to reduce processor stalls

2.2.2.1 Cache

One of the most common solutions is caching.

A **cache** is a hardware or software **component that stores data so that future requests for that data can be served faster**; the data stored in a cache might be the result of an earlier computation or a copy of data stored elsewhere.

- A **Cache Hit** occurs when the **requested data is found** in a cache;
- A **Cache Miss** occurs when **it cannot**.

Cache hits are served by reading data from the cache, which is faster than recomputing a result or reading from a slower data store, so the **more requests that can be served from the cache, the faster the system performs**. [9]

Many modern CPUs have logic that predicts what data will be accessed in the future and “pre-fetches” that data into caches. **Prefetching** reduces stalls because the data is resident in the cache when it is accessed. But beware, the other side of the coin is that if the **guess is wrong**, the **performance is worse than the system without prefetching!**

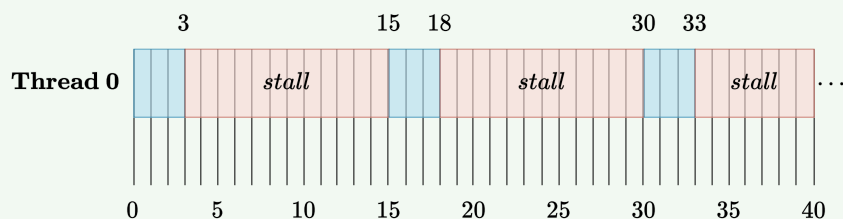
2.2.2.2 Multi-threading

A **Multithreaded Processor** is one that has the **ability to follow multiple instruction streams without software intervention**. In practice, then, this includes any machine that stores **multiple program counters (PCs) in hardware within the processor** (i.e., on chip, for microprocessor-era machines). [8]

The main idea in this architecture is to **interleave processing of multiple threads on the same core to hide stalls**. In other words, if we can’t make progress on the current thread, we work on another one.

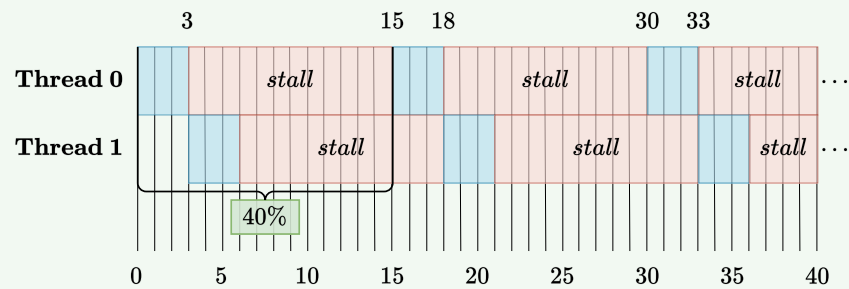
Example 1: core utilization

To better understand our explanation, suppose we are running a program where threads perform **three arithmetic instructions followed by a memory load** (with 12 cycle latency).



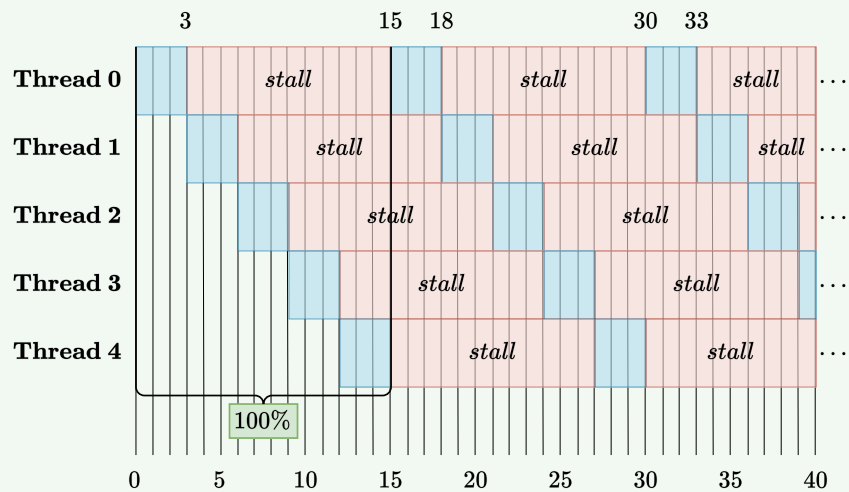
From the figure, it is clear that if we consider an arithmetic instruction and a memory stall, the core is not fully optimized to work at 100%. In

practice, we see that a single arithmetic instruction takes 3 clock cycles and the memory stall takes 12 clock cycles. This means that from this situation we are using the CPU at only 20% (3 work clock cycles on 15)! Without suggesting the final solution, try to see what happens when we add another thread.



We have gained three clock cycles, and now we are taking advantage of the 40% of the core.

Now, how many threads do we need to achieve 100% utilization? The answer is simple: the number of clock cycles of the operations to be done before the stall plus the clock cycles of the stall divided by the working operations (operations that are not stalls). In our case: $15 \div 3 = 5$.



Note that if we add more threads, there will be no benefit because the CPU is already at 100%.

✓ Multithreaded Processor benefits

- A processor with multiple hardware threads has the **ability to avoid stalls** by executing instructions from other threads when one thread must wait for a long latency operation to complete. The latency of the memory operation is not changed by multithreading, it just no longer causes reduced processor utilization.
- A multithreaded processor hides memory latency by performing

arithmetic from other threads. Program that feature more arithmetic per memory access need fewer threads to hide memory stalls.

Type of hardware-supported multithreading

- **Core manages execution contexts for multiple threads.** This type still has the same number of ALU resources: multi-threading only helps to use them more efficiently in the face of high latency operations such as memory access. The **processor decides which thread to run each clock cycle.**
- **Coarse-Grain Multithreading**, also called **Block Multithreading** or **Switch-On-Event Multithreading**, has multiple hardware contexts associated with each processor core. A hardware context is the program counter, register file, and other data required to enable a software thread to execute on a core. However, only one hardware context has access to the pipeline at a time. [8]
- **Fine-Grain Multithreading (FGMT)**, also known as **Interleaved Multithreading** or **Temporal Multithreading**, is the type just described on the previous pages, as in the example on page 27.
- **Simultaneous Multithreading (SMT)** has multiple hardware contexts associated with each core. In a simultaneous multithreaded processor, instructions from multiple threads are available to be issued on any cycle. Therefore, all hardware contexts, and in particular all register files, must be accessible to the pipeline and its execution resources. [8]

In other words, **each clock**, the **core selects instructions from multiple threads to execute on ALUs.**

3 Programming models

3.1 Implicit SPMD Program Compiler (ISPC)

Before introducing the ISPC compiler, we give the definition of SPMD.

Definition 1: Single Program, Multiple Data (SPMD)

Single Program, Multiple Data (SPMD) is a term that has been used to refer to computational models for exploiting parallelism, where **multiple processors work together to execute a program to achieve faster results.**

The **difference** between *SPMD* and *SIMD* (page 24) is that in SPMD parallel execution, **multiple autonomous processors simultaneously execute the same program at independent points**, rather than in SIMD it is vectorization at the instruction level so that **each CPU instruction processes multiple data elements.**

In other words:

- SPMD: is the **programming abstraction**, because the programmer has to think; the program is written in terms of this abstraction.
- SIMD: in general, the compilers (ISPC) issue special vector instructions that execute the logic performed by each parallel instance created (ISPC gang spawned). In addition, the compilers handle the mapping of conditional control flow to vector instructions.

The difference and the terminology used by ISPC will become clearer in the following pages. We suggest that finish this section and come back here in a moment.

Definition 2: Implicit SPMD Program Compiler (ISPC)

Implicit SPMD Program Compiler (ISPC) is a **compiler for a variant of the C programming language**, with extensions for *Single Program, Multiple Data (SPMD)* programming. Under the SPMD model, the programmer writes a program that generally appears to be a regular serial program, though the execution model is actually that a number of program instances execute in parallel on the hardware. In other words, the **ISPC gives the programmer some API to do parallelization on the code; it also generates high quality SIMD code to increase performance.**

The definition, implementation, and other details are explained in the official [Intel GitHub repository](#).

❓ How it works?

Let us take a general main program; when we call an `ispc` function, it causes a **spawn of gang of ISPC program instances upon return, all instances have completed**. These instances execute the same ISPC code **simultaneously**, and **each instance has its own copy of local variables**. Take the following ISPC code as an example:

```

1 export void ispc_sinx(
2     uniform int N,
3     uniform int terms,
4     uniform float* x,
5     uniform float* result
6 ){
7     // assume N % programCount = 0
8     for (uniform int i=0; i<N; i+=programCount) {
9         int idx = i + programIndex;
10        float value = x[idx];
11        float number = x[idx] * x[idx] * x[idx];
12        uniform int denom = 6; // 3!
13        uniform int sign = -1;
14        for (uniform int j=1; j<=terms; j++) {
15            value += sign * number / denom
16            number *= x[idx] * x[idx];
17            denom *= (2*j+2) * (2*j+3);
18            sign *= -1;
19        }
20        result[idx] = value;
21    }
22 }
```

In the example, the `programCount` (row 8) and `programIndex` (row 9) variables, `uniform` (row 2, and so on) data type tell us:

- `programIndex` gives the index of the SIMD-lane being used for running each program instance (in other words, it's a varying **integer value that has value zero for the first program instance, and so forth**).
- `programCount` gives the **total number of instances in the gang**.
- A variable that is declared with the `uniform` qualifier represents a **single value that is shared across the entire gang**.

Together, these can be used to uniquely map executing program instances to input data (`programIndex` and `programCount`, `uniform data type`).

With the ISPC analogy, the **SPMD programming model** should be clear:

1. **Single thread of control** (typically a main program);
2. **Invoke the SPMD function** (in the previous example, the `ispc_sinx` function);
3. **SPMD execution**, then **multiple instances of the function run in parallel** (multiple logical threads of control);
4. **Returns and resumes a single thread** of control.

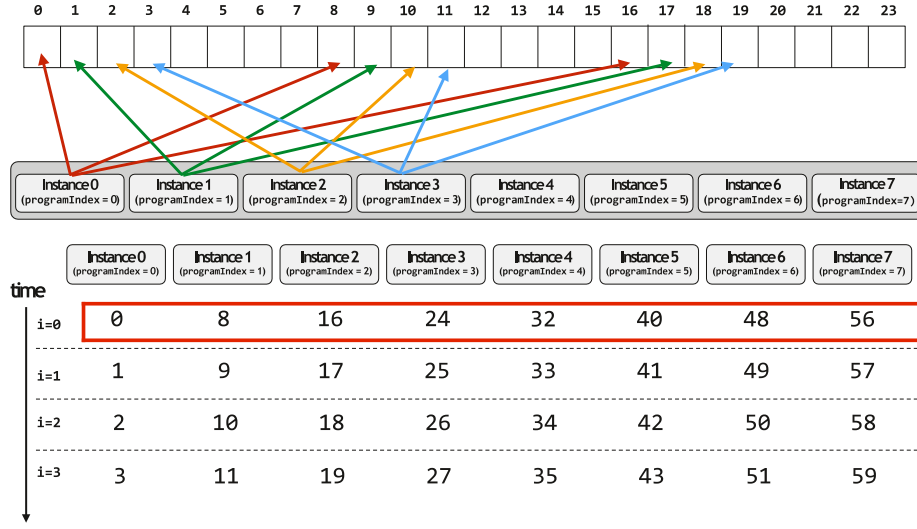


Figure 9: Example of execution with 8 instances (`programCount` equal to 8). For all program instances, there are eight non-contiguous values in memory. A special instruction called `gather` is needed to implement this, but unfortunately it is a more complex and expensive SIMD instruction rather than a contiguous implementation.

Figure 9 shows a possible execution of the ISPC function using 8 instances. The result is obtained and all is well. But there is one interesting observation. **Each ISPC instance writes each value in a non-contiguous way.** This can be done better:

```

1 export void ispc_sinx_v2(
2     uniform int N,
3     uniform int terms,
4     uniform float* x,
5     uniform float* result
6 ){
7     // assume N % programCount = 0
8     uniform int count = N / programCount;
9     int start = programIndex * count;
10    for (uniform int i=0; i<count; i++) {
11        int idx = start + i;
12        float value = x[idx];
13        float number = x[idx] * x[idx] * x[idx];
14        uniform int denom = 6; // 3!
15        uniform int sign = -1;
16        for (uniform int j=1; j<=terms; j++) {
17            value += sign * number / denom
18            number *= x[idx] * x[idx];
19            denom *= (j+3) * (j+4);
20            sign *= -1;
21        }
22        result[idx] = value;
23    }
24 }
```

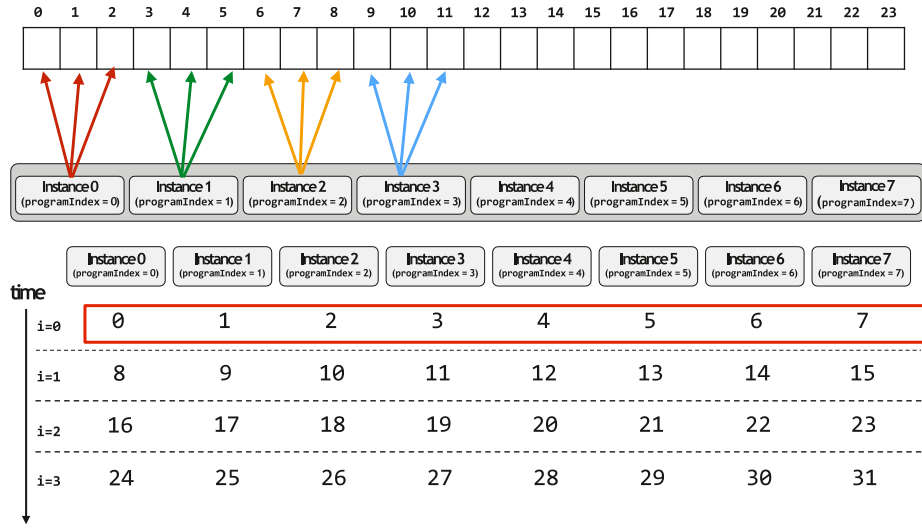


Figure 10: Example of execution with 8 instances (`programCount` equal to 8). A single “packed vector load” instruction efficiently implements this. For all program instances, since the eight values are contiguous in memory.

3.2 Shared Address Space Model

We give a general introduction to memory in the chapter 2.2. In parallel computing theory, **each thread communicates with other threads using read/write operations**. These instructions operate on a special area called the **Shared Address Space** (also called **Shared Variables**).

Definition 3: Shared Address Space

The **Shared Address Space** view of a parallel platform supports a common data space that is accessible to all processors. Processors interact by modifying data objects stored in this shared-address-space. [6]

Now the first and trivial question should be: a powerful tool is the possibility to allow communication between threads, but *how can we guarantee that two or more threads accessing the same resource do not create well known problems, such as [race condition](#)?*

This property, commonly called **mutual exclusion** or **atomic operation**, can be guaranteed with some techniques:

- **Lock/Unlock mutex around a critical section:**

```
1 Lock lock_variable;
2
3 // some operations, such as spawn of threads
4
5 lock_variable.lock();
6 // critical section
7 lock_variable.unlock();
```

- Some languages have first-class support for atomicity of code blocks:

```
1 atomic {
2     // critical section
3 }
```

- Intrinsic for hardware-supported atomic read-modify-write operations:

```
1 atomicAdd(x, 10);
```

The shared address space **requires hardware support to be efficiently implemented**. The main idea is that **each processor can directly reference the contents of any memory location**. Some interesting examples that can be explored in depth are: [SUN Niagara 2](#), [Knights landing \(KNL\)](#); [2nd Generation Intel Xeon Phi processor](#).

4 Parallel Programming Models and pthreads

4.1 How to create parallel algorithms and programs

Although *parallel algorithms* and *parallel programs* are in the same father set, the parallel computing topic, these two arguments are a little different.

Definition 1: Parallel Algorithms

A **parallel algorithm**, as opposed to a traditional serial algorithm, is an **algorithm which can do multiple operations in a given time**.

Definition 2: Parallel Programs

A **parallel program** is a **program that uses multiple CPU cores, with each core performing a task independently**.

However, designing *parallel algorithms* is not an easy task because there is no heuristic for designing *parallel algorithms*. There are some rules that help in the design. The same reasoning applies to *parallel programs*, because they depend on the chosen language and architecture.

Furthermore, there is no single correct solution, but several possible parallel solutions. A **good first approach is to start with machine-independent issues** (concurrency) and **delay target-specific issues as much as possible**.

Design a parallel algorithm	Design a parallel program
Understand the problem to be solved	Analyze the target architecture(s)
Analyze data dependencies	Choose the best parallel programming model and language
Partition the solution	Analyze the communications (cost, latency, bandwidth, visibility, synchronization, etc.)

Table 2: Design parallel algorithms and parallel programs.

The **PCAM (Partitioning, Communication, Agglomeration, Mapping)** methodology described by [Argonne National Laboratory](#) is intended to promote an exploratory **approach to design in which machine independent issues**, such as *concurrency*, are **considered early** and **machine specific aspects of design are deferred until late in the design process**. In other words, we immediately consider the machine-independent issues (e.g., concurrency) at the beginning of the design approach, and all machine-specific aspects are postponed to an advanced stage of the design process.

This methodology structures the design process into **four distinct stages**:

1. **Partitioning.** The **computation** that is to be performed and the **data** operated on by this computation are **decomposed into small tasks**. Practical issues such as the number of processors in the target computer are ignored, and attention is focused on recognizing opportunities for parallel execution.
2. **Communication.** The communication required to **coordinate task execution** is determined, and appropriate **communication structures and algorithms** are defined.
3. **Agglomeration.** The task and **communication structures** defined in the first two stages of a design are **evaluated with respect to performance requirements and implementation costs**. If necessary, tasks are combined into larger tasks to improve performance or to reduce development costs.
4. **Mapping.** **Each task is assigned to a processor** in a manner that attempts to satisfy the competing goals of **maximizing processor utilization** and **minimizing communication costs**. Mapping can be specified statically or determined at runtime by load-balancing algorithms.

In the **first two stages**, we focus on **concurrency** and **scalability** and seek to **discover algorithms with these qualities**. In the **third and fourth stages**, attention shifts to **locality** and **other performance-related issues**.

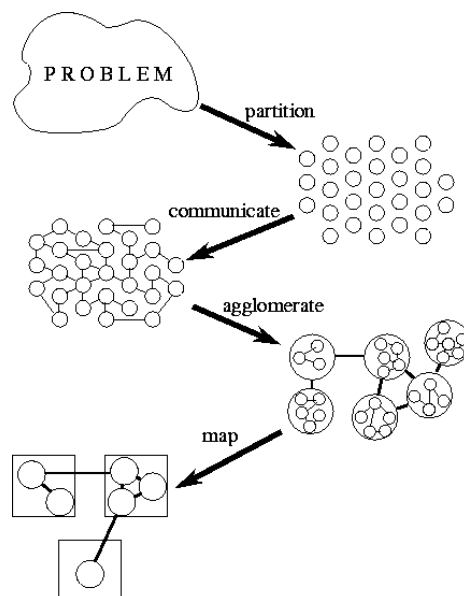


Figure 11: PCAM design methodology for parallel programs. Starting with a problem specification, develop a partition, determine communication requirements, agglomerate tasks, and finally map tasks to processors.

4.2 Analyze parallel algorithms

Whether we want to analyze our parallel algorithm created with the PCAM model or evaluate a general parallel algorithm, we need some metrics.

The classical **metrics** needed to evaluate a parallel algorithm are:

- **Time complexity:** quantifies the amount of **time required to produce a solution**.
- **Resource complexity:** quantifies how many **resources are needed to produce the solution in that time**.

In general, to analyze a **parallel algorithm**, we can consider its **structure as a directed acyclic graph (DAG)**², where the nodes are the task and the edges are the data dependencies.

Parallel Algorithm Terminology and Metrics

- **Concurrent tasks**, each task is executed *independently*.
- **Parallel tasks**, each task is executed at the *same time* (because multiple computing resources are available).
- **Work W** is the *number of operations executed*. It may be higher than the sequential version of the algorithm due to communication overhead, etc.
- **Span S** is the *longest chain of dependencies* (i.e., the critical path) that determines the *minimum time required to execute the algorithm*. This is a *lower bound* on the running time, regardless of the number of processors. The range indicates the ability of an algorithm to get better performance on more processors.

How do we calculate the Span metric?

1. As we just said, we **represent a parallel algorithm as a DAG** graph, where nodes represent tasks and edges represent dependencies between tasks;
2. We **assign weights to each node** that represent the **time required to perform the corresponding task**;
3. We try to **find the Critical Path**. In other words, we determine the path from the start node to the end node that has the *maximum cumulative weight*;
4. Finally, the **sum of the weights of the nodes on the critical path** gives us the span value!

²A **Directed Acyclic Graph (DAG)** is a directed graph, i.e. with oriented edges, without cycles.

- **Parallelism** P is the *measure of efficiency in the use of resources*. Trivially, it is the number of operations performed divided by the longest chain of dependencies:

$$P = \frac{W}{S} \quad (11)$$

It indicates **how many processors can be effectively used by the computation**. If the work is equal to the span, the parallelism is 1 and the computation is sequential. Ideally, but not necessarily, we win with polylogarithmic span, because if the work is $O(n \log n)$ and the span is $O(\log^2 n)$, then the parallelism is $O\left(\frac{n}{\log n}\right)$, which is actually quite high (and unlikely to be a bottleneck on most machines in the next 25 years). [4]

This measure is *one of the most important*. It indicates the **number of processors that are not idle**. It is obvious that a **good parallel algorithm is designed to have the lowest possible work** (less operation, then less resource usage, then less cost, and so on) **and the highest possible parallelism** (achievable by reducing the span, and this should be trivial, since the metric P is given by work divided by the span, so reducing the denominator, you can get a higher value).

As in all things, there is a **trade-off** between the lowest possible “work” and the highest possible “parallelism”. Reducing the work too much could eliminate the possibility of parallelizing our algorithm, and on the other hand, reducing the span too aggressively could cause communication/synchronization overhead.

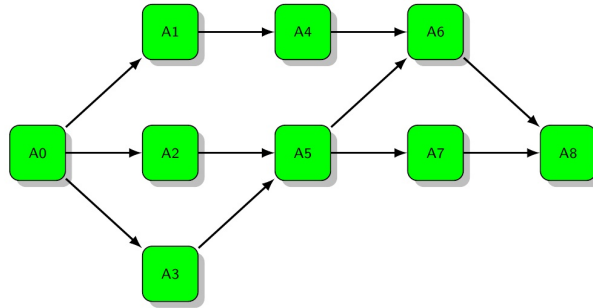


Figure 12: Example of DAG implementation with work equal to 9, span equal to 5, and parallelism equal to 1.8 ($9 \div 5$). The span calculus is not well known, has been calculated *a priori*.

Finally, we use a mathematical annotation and not only a graphical (DAG) annotation. The **Composition Rules** help determine how to combine smaller parallel tasks into a larger algorithm, while analyzing the Work and Span of the combined algorithm:

- **Single operation.** An operation takes 1 unit of work and 1 unit of span time.

$$W(op) = 1 \quad S(op) = 1$$

- **Sequential Composition.**

- The total work of executing e_1 and e_2 **sequentially** is the sum of their individual work.

$$W(e_1, e_2) = W(e_1) + W(e_2)$$

- The total span of executing e_1 and e_2 **sequentially** is the sum of their individual spans.

$$S(e_1, e_2) = S(e_1) + S(e_2)$$

- **Parallel Composition.**

- The total work of executing e_1 and e_2 in **parallel** is still the sum of their individual works.

$$W(e_1 || e_2) = W(e_1) + W(e_2)$$

- The total span of executing e_1 and e_2 in **parallel** is the *maximum of their individual spans*, since they can be executed simultaneously.

$$S(e_1 || e_2) = \max(W(e_1), W(e_2))$$

4.3 Technologies

Some famous architecture to work with parallel programming:

- **Verilog/VHDL** are *hardware description languages*. The target architectures are *ASIC and FPGA*. The **parallelism** and the **communication** are **explicit**.

✓ Pros

- Complete control on computation and memory
- No overhead introduced in the computation
- Provides access to potentially large computational power

⚠ Cons

- Requires specific hardware (e.g., ASIC or FPGA) to implement functionality
- Difficult to learn: completely different programming language and programming paradigm
- Depends on the chosen target architecture

- **MPI** is a *library*. The target architectures are *Multi CPUs*. The **parallelism** is **implicit** and the **communication** is **explicit**.

✓ Pros

- Can be adopted on different types of architecture
- Scalable solutions
- Synchronization and data communication are explicitly managed

⚠ Cons

- Communication can introduce significant overhead
- Programming paradigm more difficult than shared memory-based ones
- Standard does not reflect immediately advances in architecture characteristics

- **PThread** is a *library*. The target architectures are *Multi-core CPUs*. The **parallelism** is **explicit** and the **communication** is **implicit**.

✓ Pros

- Can be adopted on different types of architecture
- Explicit parallelism and full control over application

⚠ Cons

- Task management overhead can be significant
- Not easily scalable solutions
- Low level API

- **OpenMP** is a *C/Fortran extensions*. The target architectures are *Multi-core CPUs*. The **parallelism** is **explicit** and the **communication** is **implicit**.

✓ Pros

- Easy to learn
- Scalable solution
- Parallel applications can also be executed sequentially

⚠ Cons

- Mainly focused on shared memory homogeneous systems
- Requires small interaction between tasks

- **CUDA** is a *C extensions*. The target architectures are *CPU plus GPU(s)*. The **parallelism** is **implicit/explicit** and the **communication** is **implicit/explicit**.

✓ Pros

- Provides access to the computational power of GPUs
- Writing a CUDA kernel is quite easy
- Already optimized libraries

⚠ Cons

- Targets only NVIDIA GPUs
- Difficult to extract massive parallelism from application
- Difficult to optimize CUDA kernel

- **OpenCL** is a *C/C++ extensions and API*. The target architectures are *heterogeneous architecture*. The **parallelism** is **implicit/explicit** and the **communication** is **implicit/explicit**.

✓ Pros

- Target-independent standard
- Hides architecture details
- Same programming infrastructure for every heterogeneous architecture: CPU + GPU (and FPGA)

⚠ Cons

- Difficult programming paradigm for its heterogeneity
- Hiding of architecture details makes difficult to obtain best performances
- Gradually abandoned
- **Apache Spark** is an *API*. The target architectures are *multi CPUs*. The **parallelism** is **implicit** and the **communication** is **implicit**.

✓ Pros

- API for different languages
- Explicit parallelization and communication are not required
- Preinstalled on cloud provide VMs

⚠ Cons

- Suitable only for big data applications
- Does not (yet) fully support GPUs

Regardless of these technologies, it is quite common to mix some of them:

- **OpenMP + CUDA**: allows to exploit multi-core CPU and GPU. CUDA is used to parallelize GPU code and OpenMP is used to parallelize CPU code.
- **MPI + OpenMP**: the most common scenario are:
 1. MPI used to express coarser parallelism (multi CPU) and OpenMP used to express finer parallelism (multi core).
 2. MPI used to implement communication and OpenMP used to parallelize computation.
- **OpenCL + Verilog or VHDL**: in principle, hardware kernels (implemented for example on FPGA) can be used as accelerators; OpenMP used to describe parallelism among different processing elements; Verilog/VHDL used to describe hardware kernel. An example of target: Intel Xeon Scalable.

4.4 Threads

4.4.1 Flynn's taxonomy

Flynn's taxonomy is a **classification of computer architectures**, proposed by Michael J. Flynn. The classification system has been used as a tool in the design of modern processors and their functionalities. Since the rise of multiprocessing central processing units (CPUs), a multiprogramming context has evolved as an extension of the classification system.

The four initial classifications defined by Flynn are based upon the number of concurrent instruction (or control) streams and data streams available in the architecture:

- **Single Instruction stream, Single Data stream (SISD)**
- **Single Instruction stream, Multiple Data streams (SIMD)**
- **Multiple Instruction stream, Single Data stream (MISD)**
- **Multiple Instruction stream, Multiple Data stream (MISD)**

It is important to quote it because it is the basis for the development of many advanced technologies.

4.4.2 Definition

A UNIX process can be created by the operating system and contains information about program resources and program execution status.

Definition 3: Thread

A **thread** is an **independent stream of instructions within a process**. Threads can be scheduled by the operating system, and each thread can run concurrently with other threads. A thread also has local resources and can access the shared process resources.

In other words, a thread can be thought of as any **procedure that runs independently of its main program**. We can create each thread dynamically during execution. A good point is that a multi-threaded program is lighter than a multi-process program.

When a thread exists within a process, it shares most of the process resources, for example:

- Changes made by one thread to shared system resources (such as closing a file) will be seen by all other threads.
- Two pointers having the same value point to the same data.
- **Implicit communication** by reading and writing shared variables.
- **Reading and writing to the same memory locations requires explicit synchronization by the programmer.** If this rule is not followed, the code may suffer from a data race or race condition ³ problem.

The most common models for threaded programs are the manager / worker model⁴ and pipeline.

This chapter introduces the POSIX threads model.

Definition 4: pthreads

POSIX Threads, commonly known as **pthreads**, is an **execution model** that exists independently from a programming language, as well as a parallel execution model. It **allows a program to control multiple different flows of work that overlap in time**. Each flow of work is referred to as a thread, and creation and control over these flows is achieved by making calls to the POSIX Threads API.

POSIX threads and OpenMP are two **implementations of a shared memory parallel programming model using threads**. The **programmer is responsible for handling parallelism and synchronization**, usually through a library of subroutines or a set of compiler directives. Typically, hardware vendors have implemented their own proprietary versions of threads, but in this course we will look at POSIX threads (pthreads) and OpenMP.

³In parallel computing, a **Data Race** or **Race Condition** is a software problem that occurs when two threads (or processes) access the same variables, and at least one does a write. They can finish in a different order than expected.

⁴The manager/worker pattern is described as follows. The idea is that the work that needs to be done can be divided by a “manager” into separate pieces and the pieces can be assigned to individual “worker” processes. Thus the manager executes a different algorithm from that of the workers, but all of the workers execute the same algorithm. Most implementations of MPI allow MPI processes to be running different programs (executable files), but it is often convenient (and in some cases required) to combine the manager and worker code into a single program.

4.4.3 pthreads API

In 1995, the IEEE POSIX 1003.1c standard specified the API for explicitly managing threads. An **API is a set of C language programming types and procedure calls**.

- **Header file** to include in the main file: `pthread.h`.
- To **compile** and use it, it is necessary to add the **flag** `-pthread` to the gcc (or g++) options.

The API are divided by what we want to do. In general, there are two sets: thread management and thread synchronization.

- Thread Management
 - Creation (page 45)
 - Termination (page 46)
 - Joining (page 47)
 - Detaching (page 48)
 - Joining through Barriers (page 49)
- Thread Synchronization
 - Mutexes (page 50)
 - Condition variables (page 50)

4.4.3.1 Creation

Once threads are created, they are peers, and **may create other threads**. There is **no implied hierarchy or dependency between threads**. The **maximum number of threads depends on the implementation**.

[Doc.](#) 

pthread API: pthread_create

```
1 int pthread_create(  
2     pthread_t * thread,  
3     const pthread_attr_t * attr,  
4     void * (* start_routine) (void *),  
5     void * arg  
6 )
```

- **Return value:** on success, `pthread_create()` returns 0; on error, it returns an error number, and the contents of `*thread` are undefined.
- **Arguments:**
 - **thread:** identifier for the new thread returned by the subroutine.
 - **attr:** used to set thread attributes, such as joinable, detached, scheduling and stack size.

- **start_routine**: the C routine that the thread will execute once it is created.
- **arg**: argument passed to **start_routine**. It must be passed by address as a pointer cast of type void.

4.4.3.2 Termination

The thread returns from its startup routine when its “life” ends. The thread makes a call to the `pthread_exit` subroutine.

[Doc.](#) 

pthread API: pthread_exit

```
1 void pthread_exit(void *retval)
```

- **Return value**: this function does not return to the caller.
- **Arguments**:
 - **retval**: function terminates the calling thread and returns a value via **retval**.

The thread is canceled by another thread via the `pthread_cancel` routine.

[Doc.](#) 

pthread API: pthread_cancel

```
1 int pthread_cancel(pthread_t thread)
```

- **Return value**: on success, `pthread_cancel()` returns 0; on error, it returns a nonzero error number.
- **Arguments**:
 - **thread**: the `pthread_cancel()` function sends a cancellation request to the thread **thread**.

4.4.3.3 Joining

The join function **blocks the calling thread until the specified thread exits.**

Doc. 

pthread API: pthread_join

```
1 int pthread_join(pthread_t thread, void **retval)
```

- **Return value:** on success, `pthread_join()` returns 0; on error, it returns an error number.
- **Arguments:**
 - **thread:** the `pthread_join()` function waits for the thread specified by `thread` to terminate. If that thread has already terminated, then `pthread_join()` returns immediately. The thread specified by `thread` must be joinable.
If multiple threads simultaneously try to join with the same thread, the results are undefined. If the thread calling `pthread_join()` is canceled, then the target thread will remain joinable (i.e., it will not be detached).
 - **retval:** if `retval` is not NULL, then `pthread_join()` copies the exit status of the target thread (i.e., the value that the target thread supplied to `pthread_exit()`) into the location pointed to by `retval`. If the target thread was canceled, then `PTHREAD_CANCELED` is placed in the location pointed to by `retval`.

4.4.3.4 Detaching

The detach function **marks a thread as detached**. When a thread is detached, its **resources are automatically released back to the system when the thread terminates**, without the need for another thread to join with it.

❓ Why would I need to detach a thread and not join it?

Good question. The answer depends on what we have to do.

- **Fire and forget tasks.** When we start a thread to perform a task that doesn't require further interaction or result processing, releasing it ensures that the resources are automatically cleaned up when the task is complete.
- **Resource management.** Detaching avoids the need for another thread to call `pthread_join()`, which can save system resources and reduce the complexity of our code. It's especially useful in a highly concurrent application with many short-lived threads.
- **Avoid deadlocks.** When we have potential circular dependencies or complex synchronization between threads, detaching threads can help avoid deadlocks by eliminating the need for one thread to wait on another.
- **Long-running background tasks.** For tasks that should run independently in the background and not block the main program flow, detaching makes sense. We make sure they clean up after themselves without having to explicitly manage their lifecycle.

Doc. 

pthread API: pthread_detach

```
1 int pthread_detach(pthread_t thread)
```

- **Return value:** on success, `pthread_detach()` returns 0; on error, it returns an error number.
- **Arguments:**
 - **thread:** the `pthread_detach()` function marks the thread identified by `thread` as detached. When a detached thread terminates, its resources are automatically released back to the system without the need for another thread to join with the terminated thread. Attempting to detach an already detached thread results in unspecified behavior.

4.4.3.5 Joining through Barriers

The *barrier init* function **initializes a barrier object**, and the *barrier wait* function **blocks a thread until the specified number of threads have called it**. A **barrier object** is, in parallel computing, a synchronization tool that ensures that multiple threads reach a certain point of execution before any of them continue. It's like a **checkpoint that everyone must reach before continuing**, ensuring coordinated progress in a parallel algorithm.⁵

Doc. 

pthread API: pthread_barrier_init

```
1 int pthread_barrier_init(  
2     pthread_barrier_t * barrier,  
3     pthread_barrierattr_t * attr,  
4     unsigned int count  
5 )
```

pthread API: pthread_barrier_wait

```
1 int pthread_barrier_wait(pthread_barrier_t * barrier)
```

- **Return value:** on success, function return 0; on error, they return an error number.
- **Arguments:** the main and most important argument is `count`, which specifies the **number of threads to wait** for.

⁵For example, imagine multiple threads working on different parts of a matrix. A barrier can ensure that all threads finish their part of the computation before moving on to the next phase, such as combining results or performing subsequent operations.

4.4.3.6 Mutexes

Mutex variables are the basic **method of protecting shared data when multiple writes occur**. Only **one thread can lock a mutex variable at a time**. If multiple threads attempt to lock a mutex, only one thread will succeed. **Threads that fail to acquire the mutex are blocked**. There is also the `trylock` function, which returns immediately if the mutex is currently locked (by any thread, including the current thread). Note that a **lock function has the potential to create a deadlock situation**.

There are also **three types of mutex** that can be set using the `settype` function:

- **Normal Mutex** (`PTHREAD_MUTEX_NORMAL`). A **normal mutex does not check for errors such as deadlock**. If a thread tries to lock a mutex it already owns, the thread will deadlock.
- **Error Check Mutex** (`PTHREAD_MUTEX_ERRORCHECK`). Provides **error checking**. If a thread tries to lock a mutex it already owns, `lock` function will return an error instead of deadlocking.
- **Recursive Mutex** (`PTHREAD_MUTEX_RECURSIVE`). Allows the same thread to lock the mutex multiple times without deadlocking. Each lock must have a corresponding unlock.

Doc. 

pthread API: `pthread_mutex_lock`

```
1 int pthread_mutex_lock(pthread_mutex_t *mutex)
```

pthread API: `pthread_mutex_trylock`

```
1 int pthread_mutex_trylock(pthread_mutex_t *mutex)
```

pthread API: `pthread_mutex_unlock`

```
1 int pthread_mutex_unlock(pthread_mutex_t *mutex)
```

4.4.3.7 Condition variables

Mutexes implement synchronization by serializing data accesses. Condition variables allow threads to synchronize explicitly by signaling the meeting of a condition. Without condition variables, the programmer would need to poll to check if the condition is met.

5 OpenMP v5.2

5.1 Introduction

OpenMP is a scalable model that gives parallel programmers a simple and flexible interface for developing portable parallel applications in C/C++ and Fortran.

- **Header file** to include in the main file: `omp.h`.
- To **compile** and use it, it is necessary to add the **flag** `-fopenmp` to the `gcc` (or `c++`) options.

✓ Benefits

- Standard across a variety of shared memory architectures and platforms.
- Supports three famous languages: Fortran, C and C++.
- Scalable from embedded systems to the supercomputer.
- The directives are intuitive, and with a limited set of directives we can implement parallel algorithms.
- Incremental parallelization of serial program.
- Coarse-grained and fine-grained parallelism. See [here](#) here an interesting difference between coarse-grained and fine-grained architecture:



❓ How it works?

OpenMP is based on the **fork-join paradigm**. A “master” thread forks a specified number of “slave” threads. Tasks are divided among the “slaves”, and each “slave” runs concurrently as the runtime allocates threads to different processors.

1. **Thread #0 born.** OpenMP programs start with a single thread;
2. **Fork.** At the start of a parallel region, the *master* creates a team of parallel worker threads (*slaves*).
3. **OpenMP code block.** Statements in the parallel block are executed in parallel by each thread.
4. **Join.** At the end of the parallel region, all threads synchronize (implicit barrier, see definition on page 49) and join the master thread. [5]

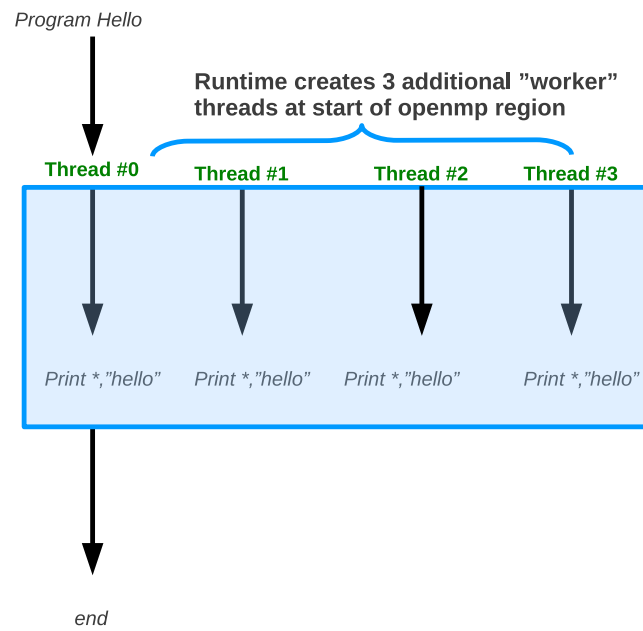


Figure 13: Example of an OpenMP program called `Hello`. At runtime, the *master* thread forks 3 additional *slave* threads to print `hello`. The example is very trivial, but here is a graphical representation of the workflow. An interesting thing to note is that **OpenMP creates a sort of region where each thread executes all the instructions in the OpenMP block**. This is important to understand. [5]

5.2 Basic syntax

We can manage OpenMP work flow using the directive syntax. We remember that the reference guide of OpenMP is available on their website:

[Reference Guide](#)



A **directive** is a combination of the base-language mechanism and a *directive-specification* (the *directive-name* followed by *optional clauses*). A construct consists of a directive and, often, additional base language code. In C++ directives are formed from either pragmas or attributes.

OpenMP: pragma omp

```
1 #pragma omp directive-specification
```

The **number of OpenMP threads can be set** using:

- At compilation time: using the environment variable OMP_NUM_THREADS
- At runtime: using the function

OpenMP: omp_set_num_threads

```
1 void omp_set_num_threads(int num_threads)
```

Other useful function to get information about threads:

- The **number of threads in the current team**:

OpenMP: omp_get_num_threads

```
1 int omp_get_num_threads()
```

The binding region for an `omp_get_num_threads` region is the innermost enclosing `parallel` region. If called from the sequential part of a program, this routine returns 1.

- The **upper bound on the number of threads** that could be used to form a new team if a `parallel` construct without a `num_threads` clause were encountered after execution returns from this routine.

OpenMP: `omp_get_max_threads`

```
1 int omp_get_max_threads()
```

- The **thread number of the calling thread**, within the current team.

OpenMP: `omp_get_thread_num`

```
1 int omp_get_thread_num()
```

OpenMP programs execute serially until they reach a `parallel` directive. As we have explained at page 52, the thread that was executing the code spawns a group of “slave” threads and becomes the “master” (thread ID 0). The code in the structured block is replicated, each thread executes a copy. At the end of the block there is an implied barrier, only the “master” thread continues.

OpenMP: `pragma omp parallel`

```
1 #pragma omp parallel optional-clauses
```

The parallel directive has **optional** clauses, the most commonly used are:

- Specify the **number of threads to spawn**:

```
1 #pragma omp parallel num_threads(int)
```

- Conditional parallelization with:

```
1 #pragma omp parallel if (condition)
```

Example 1: parallel if condition

```
1 #include <stdio.h>
2 #include <omp.h>
3
4 void test(int val)
5 {
6     #pragma omp parallel if (val != 0)
7     if (omp_in_parallel()) {
8         #pragma omp single
9         printf_s(
10             "val = %d, parallelized with %d threads\n",
11             val, omp_get_num_threads()
12         );
13 }
```

```
13     } else {  
14         printf_s("val = %d, serialized\n", val);  
15     }  
16 }  
17  
18 int main( )  
19 {  
20     omp_set_num_threads(2);  
21     test(0);  
22     test(2);  
23 }
```

The output will be:

```
1 val = 0, serialized  
2 val = 2, parallelized with 2 threads
```

- Data scope clauses (explained in the following pages).

The **number of threads in a parallel region is determined by the following factors**, in order of priority (high to low):

1. Evaluation of the `if` clause;
2. Value of the `num_threads` clause;
3. Use of the `omp_set_num_threads()` library function;
4. Setting of the `OMP_NUM_THREADS` environment variable;
5. Implementation default, e.g., the number of CPUs on a node.

References

- [1] Gene M Amdahl. Validity of the single processor approach to achieving large scale computing capabilities, reprinted from the afips conference proceedings, vol. 30 (atlantic city, nj, apr. 18–20), afips press, reston, va., 1967, pp. 483–485, when dr. amdahl was at international business machines corporation, sunnyvale, california. *IEEE Solid-State Circuits Society Newsletter*, 12(3):19–20, 2007.
- [2] Ferrandi Fabrizio. Parallel computing. Slides from the HPC-E master’s degree course on Politecnico di Milano, 2024.
- [3] John L Gustafson. Reevaluating amdahl’s law. *Communications of the ACM*, 31(5):532–533, 1988.
- [4] Laxman Dhulipala Guy E. Blelloch and Yihan Sun. Introduction to parallel algorithms. <https://www.cs.cmu.edu/~guyb/paralg/paralg/parallel.pdf>, 2024. [Accessed 22-10-2024].
- [5] Johnston Hans. OpenMP by Example. https://people.math.umass.edu/~johnston/PHI_WG_2014/OpenMPSlides_tamu_sc.pdf, 2024. [Accessed 23-10-2024].
- [6] Vipin Kumar, Ananth Grama, Anshul Gupta, and George Karypis. *Introduction to parallel computing*, volume 110. Benjamin/Cummings Redwood City, CA, 1994.
- [7] M. McCool, J. Reinders, and A. Robison. *Structured Parallel Programming: Patterns for Efficient Computation*. ITPro collection. Elsevier Science, 2012.
- [8] M. Nemirovsky and D. Tullsen. *Multithreading Architecture*. Synthesis Lectures on Computer Architecture. Springer International Publishing, 2022.
- [9] Wikipedia. Cache (computing) - Wikipedia. [https://en.wikipedia.org/wiki/Cache_\(computing\)](https://en.wikipedia.org/wiki/Cache_(computing)). [Accessed 20-10-2024].

Index

A

Amdahl's Law	19
Arbitrary Concurrent Write	6

B

Barrier Object	49
Block Multithreading	29

C

Cache	27
Cache Hit	27
Cache Miss	27
Coarse-Grain Multithreading	29
Common Concurrent Write	6
Composition Rules	38
Concurrent Read (CR)	6
Concurrent tasks in parallel algorithms	37
Concurrent Write (CW)	6

D

Data Race	44
Directed Acyclic Graph (DAG)	37

E

Error Check Mutex (PTHREAD_MUTEX_ERRORCHECK)	50
Exclusive Read (ER)	6
Exclusive Write (EW)	6

F

Fine-Grain Multithreading (FGMT)	29
Flynn's taxonomy	43

G

Gustafson's Law	20
-----------------	----

I

Implicit SPMD Program Compiler (ISPC)	30
Instruction-Level Parallelism (ILP)	23
Interleaved Multithreading	29

M

Machine Model	5
Matrix Multiply (MM) algorithm	15
Matrix-Vector Multiply (MVM) algorithm	9
Memory Access Latency	25
Memory Bandwidth	26
Multi-Core Processor (MCP)	24
Multithreaded Processor	27

N

Normal Mutex (PTHREAD_MUTEX_NORMAL) 50

O

OpenMP 51

P

Parallel Algorithm 35

Parallel Program 35

Parallel Random-Access Machine (parallel RAM or PRAM) 5

Parallel tasks in parallel algorithms 37

Parallelism metric for parallel algorithms 38

PCAM (Partitioning, Communication, Agglomeration, Mapping) 35

POSIX Threads 44

Prefetching 27

Prefix Sum 14

Priority Concurrent Write 6

Processor Stall 25

R

Race Condition 44

Random Access Machine (RAM) 5

Random Concurrent Write 6

Recursive Mutex (PTHREAD_MUTEX_RECURSIVE) 50

S

Shared Address Space 34

Shared Variables 34

Simultaneous Multithreading (SMT) 29

Single Instruction, Multiple Data (SIMD) 24

Single Program Multiple Data (SPMD) 11

Single Program, Multiple Data (SPMD) 30

Span metric for parallel algorithms 37

SPMD programming model 31

Superscalar Processor 23

Switch-On-Event Multithreading 29

T

Temporal Multithreading 29

Thread 43

W

Work metric for parallel algorithms 37