

Parallel Programming on Embedded Multicore System ESP32

Universitatea Politehnica Timișoara



Marian Belean
Franz Joseph Pal

WS2019/2020
November 5, 2019

Abstract

The following documentation will focus on the principles of parallel programming in general and the mathematical background. In addition to the different parallel programming architectures, the various models for their implementation are also discussed. Moreover, in this thesis the prerequisites for mathematical calculation models, which are suitable for Parallel Programming, are elaborated.

For a practical example, the ESP32 microcontroller was chosen, an embedded multicore system. After a brief introduction to the hardware itself, further details of the project structure and the development of the application will be presented. Therefore, a short example will be explained to focus on the basics of parallel programming.

Finally, the aim of the project and the documentation is an automatic benchmark setup and a webfrontend result overview for visualization purposes, which will be discussed in more detail in the conclusion.

Declaration

I hereby certify that I have done the final thesis on my own, that I have completely and accurately stated all the aids I have used and identified everything individually, which was taken from the work of others unchanged or with modifications.

The topic of the submitted work was jointly with Mr. / Mrs. (...) (Bachelor and Master Thesis No. (...)).

Timișoara, the November 5, 2019

Signature:

Non-disclosure notice

This work contains confidential information. In spite of the anonymous presentation of the researched organisations, readers might conclude their identity. Therefore copying, quoting or publishing is not allowed without my explicit authorisation. Furthermore, disclosure of the information to anyone other than the examination board or lecturers is not authorized.

Contents

1	Introduction	1
2	Overview	2
2.1	Problem definition	2
2.2	Objective of the documentation	2
3	Parallel Programming in General	3
3.1	Basic Concept	4
3.1.1	Amdahl's Law	4
3.1.2	Gustafson's Law's	6
3.1.3	Principles of Parallel Computing	7
3.2	Definition of parallel mathematical computations	9
3.3	Parallel Computer Architecture	10
3.3.1	Flynn's Taxonomy of Parallel Architectures	10
3.3.2	Thread Level Parallelism	12
3.4	Parallel Programming Models	12
3.4.1	Classification of Parallel Programming Models	13
4	Project documentation	14
4.1	Concept development	14
4.2	Project structure	15
4.3	Simple mathematical computation examples for Parallel Programming	16
4.4	Class diagramm	17
4.4.1	C++ Backend benchmark	17
4.4.2	Vuejs Frontend	17
4.5	Benchmark setup	18
5	Conclusion	19
A	Additional documents	20

Chapter 1

Introduction

Multicore systems are becoming increasingly popular as part of digitization and Industry 4.0¹ (German/EU) [2] [or 1] - also known as smart manufacturing in the USA [see 20, p1] - and are playing an important role in data processing and process automation [see 23, p294] [or 19, p1]. On the other hand, in addition to efficiency in energy consumption, performance in terms of computation time [see 23, p294] is required in almost every application field of multicore systems.

In fact, multicore hardware is not only expecially for smart manufacturing. Nowerdays in almost every smart application like smart phones², wearables³ or home automation we can find multicore embedded hardware platforms, which garantued high performance [see 4, p7], network connectivity, security and reliability [see 24, p5]. This field of application is also known as Internet of Things (IoT)⁴.

Especially for embedded systems mathematical models as well as numerical solutions, which can be executed both simply and parallel, are suitable. The question arises to what extent parallel execution of different sub-tasks to calculate a problem [see 22, p4] increases the desired cost factor in terms of energy consumption [see 10, chapter 3] and computational efficiency [see 22, p4 Figure 3].

¹add.: <https://www.epicor.com/en-ae/resource-center/articles/what-is-industry-4-0/>

²e.g. ARM based processors for mobile phones like <https://www.arm.com/solutions/mobile-computing/smartphones>

³add. information on ARM based solutions and the current trend in wearables: <https://www.arm.com/solutions/wearables>

⁴for additional information about Internet of Things, please see [14]

Chapter 2

Overview

2.1 Problem definition

Compared to single-core execution of tasks, multi-core embedded hardware platforms like the ESP32¹ provide the ability to develop advanced parallel computing software applications to reduce execution time and power consumption.

On the one hand, a major problem is choosing the right hardware platform to meet the cost and size factor, and moreover, whether a single-core or multi-core calculation is required. Therefore, context switching time, power consumption and total execution time must be included in the evaluation.

In order to develop an optimal solution, the hardware platform must be included in addition to the mathematical model of the problem itself. So in this case, suitable prerequisites and characteristics can be worked out in order to make an evaluation of "Parallel Computation Tasks on Embedded Multicore Systems" possible.

2.2 Objective of the documentation

The main goal of this documentation is to focus on the current parallel programming techniques, depending on the execution time in general and the required mathematical model. For this purpose, an application which can compute different sections of the Mandelbrot fractal [see 13, p11] will be developed to compare single core and multi-core calculations. Before the practical implementation, an investigation based on parallel architectures and programming models will be conducted.

The elaboration is divided into three different chapters: In the first chapter, the results of the general research are presented [see Chapter 3]. After that, the second chapter is pointing out the practical implementation of the developed application on the ESP32 [see Chapter 4]. In the end, the results including the webforntend and the automatic benchmark setup [see Chapter 5] will be discussed.

¹add. information: <https://www.espressif.com/en/products/hardware/socs>

Chapter 3

Parallel Programming in General

Since the 1970s [22], the decade in which the microprocessor era started, the overall performance of a processor has increased [10]. This goal was achieved by several points, including “*sophisticated process technology, innovative architecture or micro-architecture*” [see 10, Chapter 1, p2]. In fact, increasing the clock speed of a single core processor, like Moore’s Law predicted [22], was usually reached by increasing the number of transistors on the chip [22]. However, this go along side with the increase in complexity [see 22, Pollack’s rule], which mean, that doubling the logic of a processor result in a performance boost of only 40% [see 22, Chapter 2].

Another huge problem chip manufacturers have to deal with is leakage power [see 10, Chapter 2, p3], because the “*transistor leakage current increases as the chip size shrinks*” [see 22, p2] [see Chart 3.1]. An increase of leakage current of the transistors also result in a increase of the die’s temperature [10] along side the total power consumption as well.

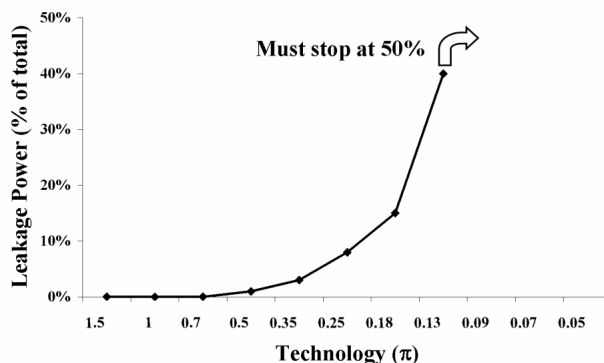


Figure 3.1: Leakage Power (% of total) vs. process technology [10]

Furthermore, a increase of the processor clock frequency to speed up the performance is only available to a suffisticating limit of 4GHz [22]. After this frequency threshold, also known as reaching the power wall, the “*power dissipation*” [see 22, p2] increases again.

Facing these types of problems such as “*chip fabrication costs, fault tolerance, power efficiency, heat dissipation*” [see 22, p3] along side with increasing processor performance, the only possible solution chip manufacturers and companies could offer was parallelism.

3.1 Basic Concept

Parallelism for processing is not something new. But due to the fact that real thread level parallelism [see Chapter 3.3.2] was only available after dual or multi-core processors were invented in 2005 [11], the topic itself and efficient software implementations are still treated in scientific work [3] [18].

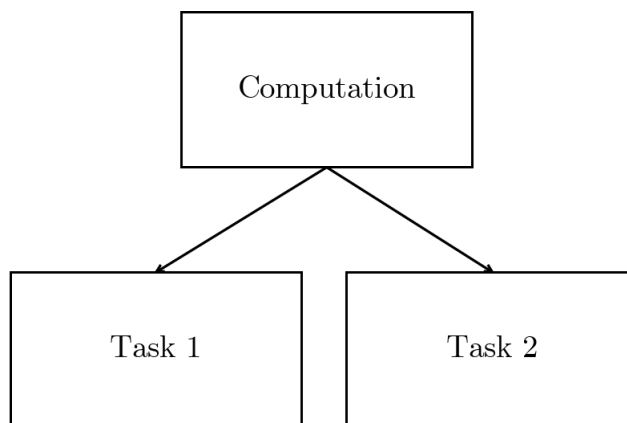


Figure 3.2: The basic concept of a simple concurrency computation.

In general, parallelism for programming means to split up a task or a computation into several sub tasks [e.g. Figure 3.2] or results, to decrease the execution time. Depending on the problem itself, these separated tasks can be independent or connected [see Chapter 3.1.3]. If we want to talk about the general concept of parallelism, we have to take a closer look to some mathematical laws, which try to describe the availability to parallel task execution and their limits.

3.1.1 Amdahl's Law

The first one, which quantifies parallelism, is called Amdahl's Law [6]. During the publication period of Amdahl's paper [5], critics claimed "*that the organization of a single computer has reached its limits and that truly significant advances can be made only by interconnection of a multiplicity of computers*" [see 6, p80]. Of course this can be transferred on single and multi-core processors or even on multi threading [see 15, Chapter 1.3, p2], but in fact, like Amdahl claimed too, addressing hardware [6], and nowadays switching context time was not considered in this case.

Amdahl's Law wants "*to provide an upper limit on speedup*" [see 6, p81] in general to point out that there is an overhead [6], which can not be implemented in parallel, but at the same time, "*apart from the sequential fraction, the remaining computations are perfectly parallelizable*" [see 6, p81].

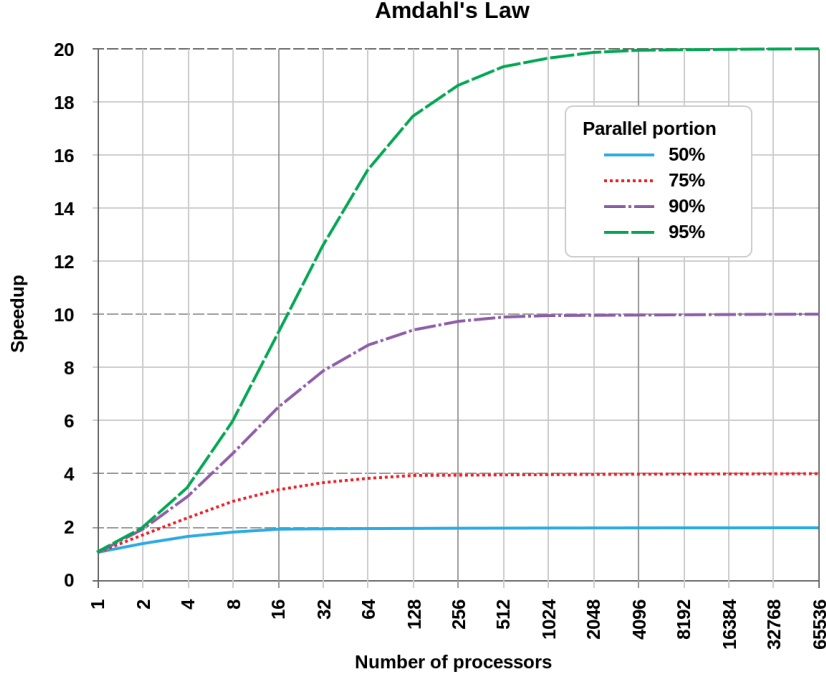


Figure 3.3: The limited speed-up of a program, which can be parallelized, depending on the number of parallel executions [9] [similar to 15, p4].

So regarding to Amdahl's Law, there has to be an upper limit of parallelism, due to the fact, that some sequential fractions still exists. In order to be able to describe this relationship, the time required to perform a calculation at once is related to the total time in parallel execution:

“Let t_1 be the time taken by one processor solving a computational problem and t_p be the time taken by p processors solving the same problem. Finally let us denote the supposed inherently sequential fraction of instructions by f . Then, according to Amdahl, $t_p = t_1(f + (1-f)/p)$ and the speedup obtainable by p processors can be expressed as” [see 6, p81]:

$$s = \frac{t_1}{t_p} = \frac{1}{f + (1-f)/p} \quad (3.1)$$

For example, a program which contains 90% of parallelizable code [see Figure 3.3] reaches his speed up limit at around 512 cores [formula 3.2]; in this case we substitute $f = 0.1$. After that number of processor cores, an significant speed up increase is not noticeable anymore .

$$\lim_{p \rightarrow \infty} \left(\frac{t_1}{t_p} \right) = \lim_{p \rightarrow \infty} \left(\frac{1}{0.1 + (1-0.1)/p} \right) = 10 \quad (3.2)$$

Many other authors tried to claim that this upper limit of speed up, both in theory and practice, is not the final end. To proof that, only to mention a few , they took into account the “energy per instruction (EPI)” [Annavaram et al. in 6, Chapter 3, p81], a case study depending on “asymmetric (or heterogeneous) chip multiprocessor architectures” [Kumar et al. in 6, Chapter 3, p81] or even considering “disk arrays to reduce input-output requirements” [Patterson et al. in 6, Chapter 3, p81].

3.1.2 Gustafson's Law's

Due to the fact that Gustafson's Law's is based on “*the same concepts as the bases of Amdahl's law, it is more a variant, rather than a refutation*” [see 6, p81]. But in fact, it is another mathematical consideration, which offers, in comparison to Amdahl's Law, no upper speed up limit regarding parallelism. Related to Gustafson, the time a single core processor needs solving the same computational problem on the sequential would be ft_p , and on the parallelizable part $(1-f)pt_p$. Therefore, the total amount of achievable speed up by p processors can thus be calculated

$$s = \frac{t_1}{t_p} = \frac{ft_p + (1-f)pt_p}{t_p} = f + (1-f)p \quad (3.3)$$

using the formula 3.3 above. In this case, “*f is the same “inherently sequential” fraction of instructions as in the case of Amdahl's law*” [see 6, p81]. In addition to that, he doesn't take the “*sequential input-output requirements proportional to input and output sizes into account*” [see 6, p81].

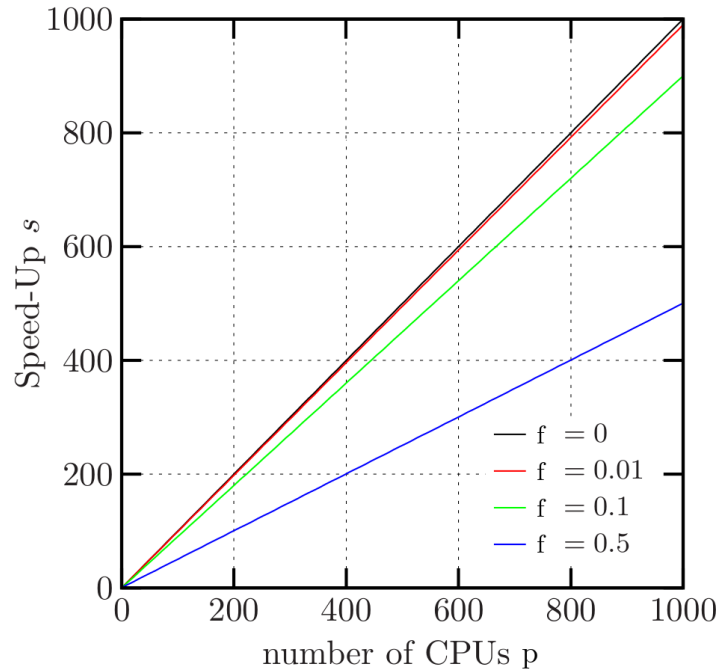


Figure 3.4: Gustafson's Law. In contrast to Amdahl we now have no upper limit to the speed up. f only determines the slope of the speed-up [12].

Regarding to [6], neither Amdahl's or Gustafson's Law are suitable to quantify parallelism in theory and practice, because both don't take into account, that “*sequential fractions of computations have negligible effect on speedup if the growth rate of the parallelizable fraction is higher than that of the sequential fraction*” [see 6, Chapter 7, p88].

Furthermore, [6] point out, that no simple formula governing parallelism exists. Both laws are more of an attempt to describe experimental results, and therefore understood rather as a draft rule of thumb as a law.

3.1.3 Principles of Parallel Computing

Despite trying to quantify the ability to parallelism task processing, it is also important to keep in mind the basic aim of parallelism to mention a “*design for concurrency*” [see 16, p4]. First of all, reducing the execution time is one of the most important objective in concurrency to make applications more **efficient**. Computations, and even tasks, are getting more and more complicated [7], not only in the application field of scientific researches. Today’s software applications require sophisticated hardware like multi-core processors, to offer a suitable user experience, for example in gaming (e.g VR), augmented reality in automation (e.g. AR) or for IoT [see Chapter 1] devices.

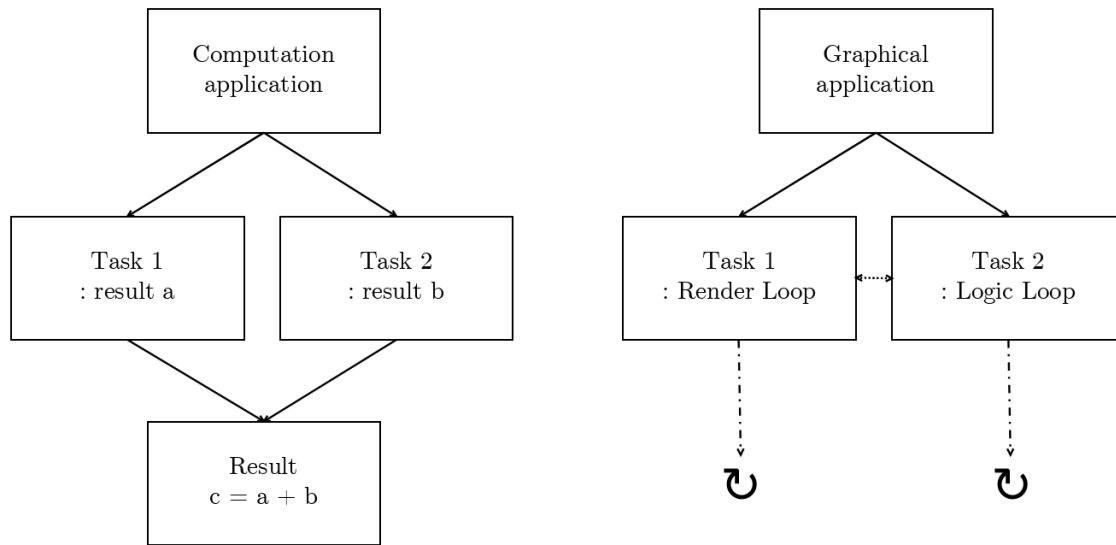


Figure 3.5: Comparison of independent and dependent tasks in a general concurrency application.

In case of software parallelism, which is the main objective of this chapter, we have to distinguish between tasks, which are independent and tasks, which are connected (called dependent). It depends largely on the application itself: e.g. for example, a numerical calculation [see Chapter 3.2] can be easily subdivided into sub tasks to speed up the computation, whereas a graphical application needs a logical loop and a render loop, both independently in different tasks [see Fig. 3.5]; a data exchange usually takes place via a shared memory [more on this in Chapter 3.3]. As already suggested in Chapter 3.2, our work will be limited to the division of numerical or generic mathematical calculations.

In addition to that, multi-core software applications offers also the opportunity, to provide low power systems, because **power consumption** results from performance and clock frequency [see 22, Fig. 3, p4], along side instructions per core (IPC) [10] [further inf. in Chapter 3].

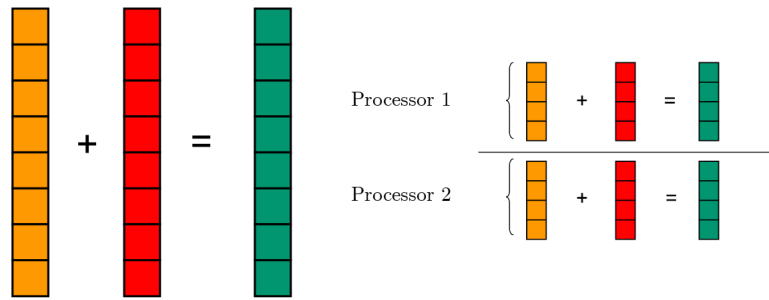


Figure 3.6: Adding to arrays of integers for speed up in concurrency [see 7, p3].

Furthermore, to guarantee **flexibility** for both, the software application and hardware platform, multi-core systems and concurrency is the way to go to fit customers and companies claims. This can easily be proofed due to the fact that today's hard- and software manufacturers always try to provide generalized solutions.

It can not be denied that as the degree of parallelism increases, the complexity of the application including the hardware realization also increases [22] [16]. A basic principle and design pattern of parallel computing is therefore to ensure maximum efficiency through parallelism while reducing complexity. A usual rule can be transferred on this topic: ensure **simplicity**.

Regarding to Figure 3.5 and Chapter 3.2, the best way to implement and computation concurrent is to guarantee that the sub tasks can work independent from each other. This has a major effect on the performance and complexity of the software implementation: “massively parallel vs. embarrassingly serial” [see 16, p15], alongside the ability to take less attention on control issues such as task orders, synchronization, (shared) memory access or even task communication [8]. In fact, complete concurrency is not possible, due to the fact that splitting a computation into sub tasks requires at least on remaining worker task to collect and combine the sub task results. Anyway, **Independence** of sub tasks enables almost the best efficiency.

Emphasises design [see 16, p5]: In summary, therefore, the following points should be seen as goals and thus as groundbreaking for the basic principles of parallelism:

1. Efficiency:

Concurrency in hard- and software to solve large problems in less amount of time to reduce execution time and power consumption.

2. Flexibility:

Environments will be more heterogeneous and the use in different application areas will be enabled.

3. Simplicity:

4. Independence:

These principles result in four different **design patterns** [see 16, p11 ff.] for parallel software applications:

- Finding concurrency:

This should be the main aim of all software applications today, set the case it makes sense.

- Algorithm structure:

To ensure proper efficient, the implementation should be based on usual parallel programming models [see Chapter 3.4].

- Supporting structures:

“Useful idioms rather than unique implementations” [see 16, p25] like Loop Parallelism, Fork/Join, Shared Data or Shared Queue [16].

- Implementation Mechanism:

Using programming languages which offer the opportunity to real parallel computing (e.g. C++, OpenMP & Pthreads, MPI, OpenCL) [16].

3.2 Definition of parallel mathematical computations

Mathematical examples [6]:

- ...[see 8, p8]
- ...[see 7, p4]
- ...[see 17, p398]

3.3 Parallel Computer Architecture

The terminology of computer architecture was invented in 1960s by the designers of IBM System to describe the structure of a computer. Computer architect's task is to write a suitable program code for the machine, keeping in mind every time this structure of computer, understanding all the factors like state-of-the-art technologies at each design level and changing those designs tradeoffs for their specific applications.

Parallel computing means the situation where tasks are separated into discrete parts that can be executed concurrently. Each part is diffused into a larger series of instructions which will be executed simultaneously on different CPU's or even in a pipeline [see Figure 3.7]. These kind of parallel systems have to deal with the simultaneous use of multiple computer resources that can include either a single computer with multiple CPU's, or a number of computers connected by a network creating a parallel processing cluster or combination of both.

The crux of parallel processing are CPU's. Based on the number of instruction and data streams that can be processed simultaneously, computing systems are classified into four major categories based on Flynn's Taxonomy.

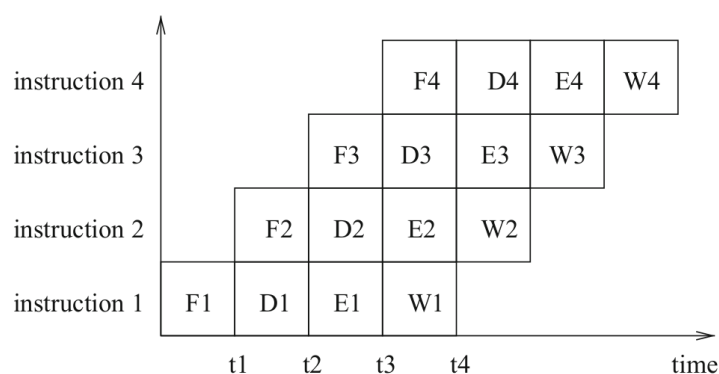


Figure 3.7: Overlapping execution of four independent instructions by pipelining. The execution of each instruction is split into four stages: *fetch (F)*, *decode (D)*, *execute (E)*, and *write back (W)* [see 21, Fig. 2.1, p11].

3.3.1 Flynn's Taxonomy of Parallel Architectures

Flynn's classification was first elaborated and proposed by Michael Flynn in 1966 and represents a scheme which is based on the notion of information stream. The term 'stream' defines a sequence or flow of either one of both existent types of information which flows and are operated into a processor: instructions or data.

Instruction stream defines the sequence of instructions performed by CPU, as in the same time the data stream defines the data traffic exchanged between the memory and CPU. His taxonomy left aside the machine's structure for classification of parallel computers and took over a whole new concept focusing on multiplicity of instructions and data streams observed by the CPU during execution.

The major four categories are the followings [comp. to Fig. 3.8]:

1. **SISD** (single-instruction, single-data) systems:

It designs an sequential computer which exploits no parallelism in either the instruction stream nor data stream. An SISD computing system is a uniprocessor machine capable of single stream executions.

2. **SIMD** (single-instruction, multiple-data) systems:

It designs a multiprocessor machine capable of executing a single instruction stream on multiple different data streams. Instructions can be executed sequentially, such as by pipe-lining, or in parallel by multiple functional units.

3. **MISD** (multiple instruction streams, single data stream) systems:

It designs a multiprocessor machine capable of executing different instructions streams on the same data stream.

4. **MIMD** (multiple instruction streams, multiple data streams) systems:

It designs a multiprocessor machine capable of executing multiple instructions streams on multiple data streams. This architectures include multi-core superscalar processors and distributed systems.

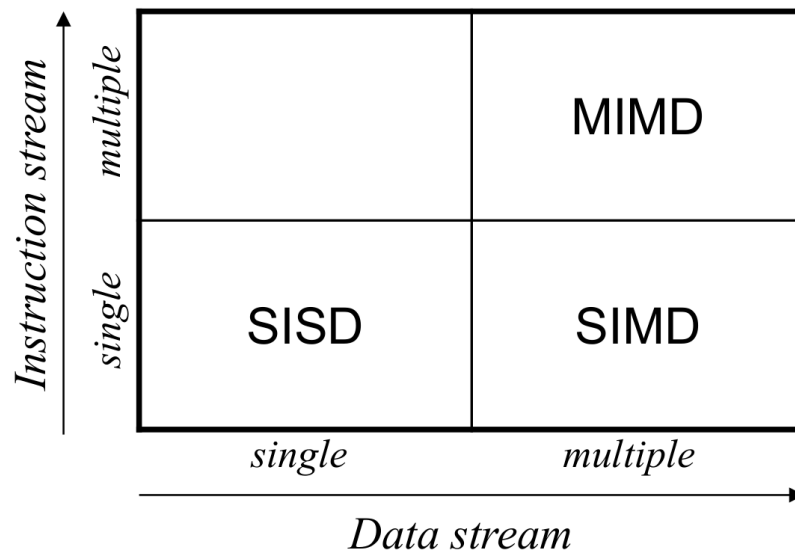


Figure 3.8: Visualization of Flynn's Taxonomy [see 8, p5].

3.3.2 Thread Level Parallelism

...[see 21, p24]

...[see 16, p14]

3.4 Parallel Programming Models

Steps to evaluate a proper parallel design [see 16, p6]:

1. Finding Concurrency
2. Algorithm Structure
3. Supporting Structures
4. Implementation Mechanisms

3.4.1 Classification of Parallel Programming Models

3.4.1.1 Process Interaction

...[see 8, p4]

3.4.1.2 Problem decomposition

...[see 21, p105 ff.]

Chapter 4

Project documentation

4.1 Concept development

...

4.2 Project structure

...

4.3 Simple mathematical computation examples for Parallel Programming

...

4.4 Class diagramm

...

4.4.1 C++ Backend benchmark

...

4.4.2 Vuejs Frontend

...

4.5 Benchmark setup

...

Chapter 5

Conclusion

...

Appendix A

Additional documents

Bibliography

- [1] Plattform Industrie 4.0. “Positionspapier, Leitbild 2030 für Industrie 4.0 - Digitale Ökosysteme global gestalten”. In: *Plattform Industrie 4.0* (Apr. 2019). URL: https://www.plattform-i40.de/PI40/Redaktion/DE/Downloads/Publikation/Positionspapier%20Leitbild.pdf?__blob=publicationFile&v=5.
- [2] Plattform Industrie 4.0. *Was is Industrie 4.0?* 2019. URL: <https://www.plattform-i40.de/PI40/Navigation/DE/Industrie40/WasIndustrie40/was-ist-industrie-40.html>.
- [3] Umut A. Acar and Guy E. Blelloch. *Algorithms: Parallel and Sequential*. Pittsburgh, USA: Carnegie Mellon University, Department of Computer Science, Feb. 2019.
- [4] Tosiron Adegbija et al. “Microprocessor Optimizations for the Internet of Things: A Survey”. In: *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* PP (June 2017), pp. 1–1. DOI: <http://dx.doi.org/10.1109/TCAD.2017.2717782>.
- [5] G.M. Amdahl. *Validity of the single processor approach to achieving large scale computing capabilities*. In: *Proceedings of Spring Joint Computer Conference*. New York: ACM, 1967, pp. 483–485.
- [6] Frank Devai. “The Refutation of Amdahl’s Law and Its Variants”. In: Sept. 2018, pp. 79–96. ISBN: 978-3-662-58038-7. DOI: http://dx.doi.org/10.1007/978-3-662-58039-4_5.
- [7] Gabriel Edgar. *An Introduction to Parallel Computing*. URL: http://www2.cs.uh.edu/~gabriel/courses/mpicourse_03_06/Introduction.pdf.
- [8] Prof. Robert van Engelen. *Parallel Programming Models*. URL: <https://www.cs.fsu.edu/~engelen/courses/HPC/Models.pdf>.
- [9] Daniels220 at English Wikipedia. *SVG Graph illustrating Amdahl’s law*. [Online; accessed October 29, 2019]. Apr. 2008. URL: <https://commons.wikimedia.org/w/index.php?curid=6678551>.
- [10] Pawel Gepner and Michal Kowalik. “Multi-Core Processors: New Way to Achieve High System Performance.” In: Jan. 2006, pp. 9–13. DOI: <http://dx.doi.org/10.1109/PARELEC.2006.54>.
- [11] Computer Hope. *Computer processor history*. 2019. URL: <https://www.computerhope.com/history/processor.htm>.

- [12] Harald Koestler, C. Moeller, and F Deserno. “Performance Results for Optical Flow on an Opteron Cluster Using a Parallel 2D/3D Multigrid Solver”. In: (Jan. 2006). URL: <https://www.researchgate.net/publication/236892123>.
- [13] Nigel Lesmoir-Gordon. “THE MANDELBROT SET, FRACTAL GEOMETRY AND BENOIT MANDELBROT - The Life and Work of a Maverick Mathematician”. In: *Medicographia* 34 (June 2012), p. 353. URL: <https://www.researchgate.net/publication/270285889>.
- [14] Sheik Dawood M. “Review on Applications of Internet of Things (IoT)”. In: (Dec. 2018). URL: <https://www.researchgate.net/publication/329672903>.
- [15] Khaled Mashfiq. *NONLINEAR EARTHQUAKE ENGINEERING SIMULATION USING PARALLEL COMPUTING SYSTEM*. Dec. 2012. DOI: <http://dx.doi.org/10.13140/RG.2.2.21215.87208>.
- [16] Tim Mattson, Beverly Sanders, and Berna Massingill. “Patterns for Parallel Programming”. In: (Sept. 2004).
- [17] Zhang N. “A Novel Parallel Scan for Multicore Processors and Its Application in Sparse Matrix-Vector Multiplication”. In: *IEEE Transactions on Parallel and Distributed Systems* 23.3 (Mar. 2012), pp. 397–404. DOI: <http://dx.doi.org/10.1109/TPDS.2011.174>.
- [18] Kota Sujatha et al. “Multicore Parallel Processing Concepts for Effective Sorting and Searching”. In: *IEEE* (2015). DOI: <http://dx.doi.org/10.1109/SPACES.2015.7058238>.
- [19] Karlsruhe Institute of Technology. “Multi-core processors for mobility and industry 4.0”. In: *PHYSORG* (2016). URL: <https://phys.org/news/2016-12-multi-core-processors-mobility-industry.html>.
- [20] Klaus-Dieter Thoben, Stefan Wiesner, and Thorsten Wuest. “”Industrie 4.0” and Smart Manufacturing – A Review of Research Issues and Application Examples”. In: *International Journal of Automation Technology* 11 (Jan. 2017), pp. 4–19. DOI: <http://dx.doi.org/10.20965/ijat.2017.p0004>.
- [21] Gudula Rünger Thomas Rauber. *Parallel Programming for Multicore and Cluster Systems*. Germany: Second Edition, Springer Verlag, 2013.
- [22] Balaji Venu. “Multi-core processors - An overview”. In: (Oct. 2011). URL: <https://www.researchgate.net/publication/51945986>.
- [23] Dragan Vuksanović, Jelena Vešić, and Davor Korčok. “Industry 4.0: the Future Concepts and New Visions of Factory of the Future Development”. In: Jan. 2016, pp. 293–298. DOI: <http://dx.doi.org/10.15308/Sinteza-2016-293-298>.
- [24] Yousaf Zikria et al. “Internet of Things (IoT) Operating Systems Management: Opportunities, Challenges, and Solution”. In: *Sensors* 8 (Apr. 2019), pp. 1–10. DOI: <http://dx.doi.org/10.3390/s19081793>.