

## 1. COURSE

CS3P1. Parallel and Distributed Computing (Mandatory)

## 2. GENERAL INFORMATION

**2.1 Course** : CS3P1. Parallel and Distributed Computing

**2.2 Semester** : 8<sup>vo</sup> Semestre.

**2.3 Credits** : 4

**2.4 Horas** : 2 HT; 4 HP;

**2.5 Duration of the period** : 16 weeks

**2.6 Type of course** : Mandatory

**2.7 Learning modality** : Blended

**2.8 Prerequisites** :

- CS212. Analysis and Design of Algorithms. (5<sup>th</sup> Sem)
- CS231. Networking and Communication. (7<sup>th</sup> Sem)
- CS212. Analysis and Design of Algorithms. (5<sup>th</sup> Sem)
- CS231. Networking and Communication. (7<sup>th</sup> Sem)

## 3. PROFESSORS

Meetings after coordination with the professor

## 4. INTRODUCTION TO THE COURSE

The last decade has brought explosive growth in computing with multiprocessors, including Multi-core processors and distributed data centers. As a result, computing parallel and distributed has become a widely elective subject to be one of the main components in the mesh studies in computer science undergraduate. Both parallel and distributed computing the simultaneous execution of multiple processes, whose operations have the potential to intercalate in a complex way. Parallel and distributed computing builds on foundations in many areas, including understanding the fundamental concepts of systems, such as: concurrency and parallel execution, consistency in state / memory manipulation, and latency. The communication and coordination between processes has its foundations in the passage of messages and models of shared memory of computing and algorithmic concepts like atomicity, consensus and conditional waiting. Achieving acceleration in practice requires an understanding of parallel algorithms, strategies for decomposition problem, systems architecture, implementation strategies and analysis of performance. Distributed systems highlight the problems of security and tolerance to Failures, emphasize the maintenance of the replicated state and introduce additional problems in the field of computer networks.

## 5. GOALS

- That the student is able to create parallel applications of medium complexity by efficiently leveraging machines with multiple cores.
- That the student is able to compare sequential and parallel applications.
- That the student is able to convert, when the situation warrants, sequential applications to parallel efficiently

## 6. COMPETENCES

- 1) Analyze a complex computing problem and to apply principles of computing and other relevant disciplines to identify solutions. (Usage)

6) Apply computer science theory and software development fundamentals to produce computing-based solutions. (**Usage**)

## 7. TOPICS

Unit 1: Parallelism Fundamentals (18)	
Competences Expected:	
Topics	Learning Outcomes
<ul style="list-style-type: none"> <li>• Multiple simultaneous computations</li> <li>• Goals of parallelism (e.g., throughput) versus concurrency (e.g., controlling access to shared resources)</li> <li>• Parallelism, communication, and coordination               <ul style="list-style-type: none"> <li>– Parallelism, communication, and coordination</li> <li>– Need for synchronization</li> </ul> </li> <li>• Programming errors not found in sequential programming               <ul style="list-style-type: none"> <li>– Data races (simultaneous read/write or write/write of shared state)</li> <li>– Higher-level races (interleavings violating program intention, undesired non-determinism)</li> <li>– Lack of liveness/progress (deadlock, starvation)</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>• Distinguish using computational resources for a faster answer from managing efficient access to a shared resource [Familiarity]</li> <li>• Distinguish multiple sufficient programming constructs for synchronization that may be interimplementable but have complementary advantages [Familiarity]</li> <li>• Distinguish data races from higher level races [Familiarity]</li> </ul>
<b>Readings :</b> [Pac11], [Mat14], [quinnz], [Geo10]	

Unit 2: Parallel Architecture (12)	
Competences Expected:	
Topics	Learning Outcomes
<ul style="list-style-type: none"> <li>• Multicore processors</li> <li>• Shared vs distributed memory</li> <li>• Symmetric multiprocessing (SMP)</li> <li>• SIMD, vector processing</li> <li>• GPU, co-processing</li> <li>• Flynn’s taxonomy</li> <li>• Instruction level support for parallel programming <ul style="list-style-type: none"> <li>– Atomic instructions such as Compare and Set</li> </ul> </li> <li>• Memory issues <ul style="list-style-type: none"> <li>– Multiprocessor caches and cache coherence</li> <li>– Non-uniform memory access (NUMA)</li> </ul> </li> <li>• Topologies <ul style="list-style-type: none"> <li>– Interconnects</li> <li>– Clusters</li> <li>– Resource sharing (e.g., buses and interconnects)</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>• Explain the differences between shared and distributed memory [Assessment]</li> <li>• Describe the SMP architecture and note its key features [Assessment]</li> <li>• Characterize the kinds of tasks that are a natural match for SIMD machines [Usage]</li> <li>• Describe the advantages and limitations of GPUs vs CPUs [Usage]</li> <li>• Explain the features of each classification in Flynn’s taxonomy [Usage]</li> <li>• Describe the challenges in maintaining cache coherence [Familiarity]</li> <li>• Describe the key performance challenges in different memory and distributed system topologies [Familiarity]</li> </ul>
<b>Readings :</b> [Pac11], [KH13], [SK10], [Geo10]	

Unit 3: Parallel Decomposition (18)	
Competences Expected:	
Topics	Learning Outcomes
<ul style="list-style-type: none"> <li>• Need for communication and coordination/synchronization</li> <li>• Independence and partitioning</li> <li>• Basic knowledge of parallel decomposition concept</li> <li>• Task-based decomposition <ul style="list-style-type: none"> <li>– Implementation strategies such as threads</li> </ul> </li> <li>• Data-parallel decomposition <ul style="list-style-type: none"> <li>– Strategies such as SIMD and MapReduce</li> </ul> </li> <li>• Actors and reactive processes (e.g., request handlers)</li> </ul>	<ul style="list-style-type: none"> <li>• Explain why synchronization is necessary in a specific parallel program [Usage]</li> <li>• Identify opportunities to partition a serial program into independent parallel modules [Familiarity]</li> <li>• Write a correct and scalable parallel algorithm [Usage]</li> <li>• Parallelize an algorithm by applying task-based decomposition [Usage]</li> <li>• Parallelize an algorithm by applying data-parallel decomposition [Usage]</li> <li>• Write a program using actors and/or reactive processes [Usage]</li> </ul>
<b>Readings :</b> [Pac11], [Mat14], [Qui03], [Geo10]	

Unit 4: Communication and Coordination (18)	
Competences Expected:	
Topics	Learning Outcomes
<ul style="list-style-type: none"> <li>• Shared Memory</li> <li>• Consistency, and its role in programming language guarantees for data-race-free programs</li> <li>• Message passing <ul style="list-style-type: none"> <li>– Point-to-point versus multicast (or event-based) messages</li> <li>– Blocking versus non-blocking styles for sending and receiving messages</li> <li>– Message buffering (cross-reference PF/Fundamental Data Structures/Queues)</li> </ul> </li> <li>• Atomicity <ul style="list-style-type: none"> <li>– Specifying and testing atomicity and safety requirements</li> <li>– Granularity of atomic accesses and updates, and the use of constructs such as critical sections or transactions to describe them</li> <li>– Mutual Exclusion using locks, semaphores, monitors, or related constructs <ul style="list-style-type: none"> <li>* Potential for liveness failures and deadlock (causes, conditions, prevention)</li> </ul> </li> <li>– Composition <ul style="list-style-type: none"> <li>* Composing larger granularity atomic actions using synchronization</li> <li>* Transactions, including optimistic and conservative approaches</li> </ul> </li> </ul> </li> <li>• Consensus <ul style="list-style-type: none"> <li>– (Cyclic) barriers, counters, or related constructs</li> </ul> </li> <li>• Conditional actions <ul style="list-style-type: none"> <li>– Conditional waiting (e.g., using condition variables)</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>• Use mutual exclusion to avoid a given race condition [Usage]</li> <li>• Give an example of an ordering of accesses among concurrent activities (eg, program with a data race) that is not sequentially consistent [Familiarity]</li> <li>• Give an example of a scenario in which blocking message sends can deadlock [Usage]</li> <li>• Explain when and why multicast or event-based messaging can be preferable to alternatives [Familiarity]</li> <li>• Write a program that correctly terminates when all of a set of concurrent tasks have completed [Usage]</li> <li>• Give an example of a scenario in which an attempted optimistic update may never complete [Familiarity]</li> <li>• Use semaphores or condition variables to block threads until a necessary precondition holds [Usage]</li> </ul>
<b>Readings :</b> [Pac11], [Mat14], [Qui03], [Geo10]	

**Unit 5: Parallel Algorithms, Analysis, and Programming (18)****Competences Expected:**

Topics	Learning Outcomes
<ul style="list-style-type: none"><li>• Critical paths, work and span, and the relation to Amdahl's law</li><li>• Speed-up and scalability</li><li>• Naturally (embarrassingly) parallel algorithms</li><li>• Parallel algorithmic patterns (divide-and-conquer, map and reduce, master-workers, others)<ul style="list-style-type: none"><li>– Specific algorithms (e.g., parallel MergeSort)</li></ul></li><li>• Parallel graph algorithms (e.g., parallel shortest path, parallel spanning tree) (cross-reference AL/Algorithmic Strategies/Divide-and-conquer)</li><li>• Parallel matrix computations</li><li>• Producer-consumer and pipelined algorithms</li><li>• Examples of non-scalable parallel algorithms</li></ul>	<ul style="list-style-type: none"><li>• Define “critical path”, “work”, and “span” [Familiarity]</li><li>• Compute the work and span, and determine the critical path with respect to a parallel execution diagram [Usage]</li><li>• Define “speed-up” and explain the notion of an algorithm's scalability in this regard [Familiarity]</li><li>• Identify independent tasks in a program that may be parallelized [Usage]</li><li>• Characterize features of a workload that allow or prevent it from being naturally parallelized [Familiarity]</li><li>• Implement a parallel divide-and-conquer (and/or graph algorithm) and empirically measure its performance relative to its sequential analog [Usage]</li><li>• Decompose a problem (eg, counting the number of occurrences of some word in a document) via map and reduce operations [Usage]</li><li>• Provide an example of a problem that fits the producer-consumer paradigm [Usage]</li><li>• Give examples of problems where pipelining would be an effective means of parallelization [Usage]</li><li>• Implement a parallel matrix algorithm [Usage]</li><li>• Identify issues that arise in producer-consumer algorithms and mechanisms that may be used for addressing them [Usage]</li></ul>
<b>Readings :</b> [Mat14], [Qui03], [Geo10]	

Unit 6: Parallel Performance (18)	
Competences Expected:	
Topics	Learning Outcomes
<ul style="list-style-type: none"> <li>• Load balancing</li> <li>• Performance measurement</li> <li>• Scheduling and contention (cross-reference OS/Scheduling and Dispatch)</li> <li>• Evaluating communication overhead</li> <li>• Data management <ul style="list-style-type: none"> <li>– Non-uniform communication costs due to proximity (cross-reference SF/Proximity)</li> <li>– Cache effects (e.g., false sharing)</li> <li>– Maintaining spatial locality</li> </ul> </li> <li>• Power usage and management</li> </ul>	<ul style="list-style-type: none"> <li>• Detect and correct a load imbalance [Usage]</li> <li>• Calculate the implications of Amdahl's law for a particular parallel algorithm (cross-reference SF/Evaluation for Amdahl's Law) [Usage]</li> <li>• Describe how data distribution/layout can affect an algorithm's communication costs [Familiarity]</li> <li>• Detect and correct an instance of false sharing [Usage]</li> <li>• Explain the impact of scheduling on parallel performance [Familiarity]</li> <li>• Explain performance impacts of data locality [Familiarity]</li> <li>• Explain the impact and trade-off related to power usage on parallel performance [Familiarity]</li> </ul>
<b>Readings :</b> [Pac11], [Mat14], [KH13], [SK10], [Geo10]	

## 8. WORKPLAN

### 8.1 Methodology

Individual and team participation is encouraged to present their ideas, motivating them with additional points in the different stages of the course evaluation.

### 8.2 Theory Sessions

The theory sessions are held in master classes with activities including active learning and roleplay to allow students to internalize the concepts.

### 8.3 Practical Sessions

The practical sessions are held in class where a series of exercises and/or practical concepts are developed through problem solving, problem solving, specific exercises and/or in application contexts.

## 9. EVALUATION SYSTEM

\*\*\*\*\* EVALUATION MISSING \*\*\*\*\*

## 10. BASIC BIBLIOGRAPHY

- [Geo10] Gerhard Wellein Georg Hager. *Introduction to High Performance Computing for Scientists and Engineers (Chapman & Hall/CRC Computational Science)*. Ed. by CRC Press. 1st. 2010. ISBN: 978-1439811924.
- [KH13] David B. Kirk and Wen-mei W. Hwu. *Programming Massively Parallel Processors: A Hands-on Approach*. 2nd. Morgan Kaufmann, 2013. ISBN: 978-0-12-415992-1.
- [Mat14] Norm Matloff. *Programming on Parallel Machines*. University of California, Davis, 2014. URL: <http://heather.cs.ucdavis.edu>
- [Pac11] Peter S. Pacheco. *An Introduction to Parallel Programming*. 1st. Morgan Kaufmann, 2011. ISBN: 978-0-12-374260-5.
- [Qui03] Michael J. Quinn. *Parallel Programming in C with MPI and OpenMP*. 1st. McGraw-Hill Education Group, 2003. ISBN: 0071232656.
- [SK10] Jason Sanders and Edward Kandrot. *CUDA by Example: An Introduction to General-Purpose GPU Programming*. 1st. Addison-Wesley Professional, 2010. ISBN: 0131387685, 9780131387683.