

CS205

HIGH PERFORMANCE COMPUTING FOR SCIENCE AND ENGINEERING

CLASS SYLLABUS

2022-12-09

Fabian Wermelinger
Harvard University

Objective:

With manufacturing processes reaching the limits in terms of transistor density on today's computing architectures, efficient modern code must exploit parallel execution to maintain scaling of available hardware resources. The use of computers in academia, industry and society is a fundamental tool for solving (scientific) problems while the "think parallel" mindset of code developers is still lagging behind. The aim of this course is to introduce the student to the fundamentals of parallel programming and its relationship on computer architectures. Various forms of parallelism are discussed and exploited through different programming models with focus on shared and distributed memory programming. The learned techniques are tried out by means of homework, lab sessions and a term project.

After completing this class, the student is capable of identifying parallelism in algorithms and exploit it with appropriate techniques. The student will further learn about methods to analyze the performance and parallel scaling of programs and know how to apply optimizations in software such that it can be deployed on large scale high performance computing (HPC) architectures.

Prerequisites:

The course assumes the student is comfortable with reading and writing code in the C, C++ or Fortran programming languages. Homework, lab and examples in class will be presented using C++. Familiarity with Linux command line tools, ssh, git and code editing tools is assumed.

Prerequisites include CS50, CS61, CS107, CS161 or AC207. This class is not suited for first-year students. CS205 will not teach basics of programming.

Textbooks:

The class does not follow a specific textbook. The following textbooks are suitable for additional reference:

- *"Introduction to High Performance Scientific Computing"*, V. Eijkhout, [free pdf 3rd edition 2020](#)
- *"Parallel Programming for Science and Engineering"*, V. Eijkhout, [free pdf 2nd edition 2020](#)
- *"An Introduction to Parallel Programming"*, P. Pacheco, Morgan Kaufmann 2011
- *"Introduction to High Performance Computing for Scientists and Engineers"*, G. Hager and G. Wellein, CRC Press 2011
- *"Computer Organization and Design"*, D. Patterson and J. Hennessy, Morgan Kaufmann 2018 (RISC-V edition)
- *"Computer Architecture"*, J. Hennessy and D. Patterson, Morgan Kaufmann 2019
- *"Programming Massively Parallel Processors"*, D. Kirk and W. Hwu, Morgan Kaufmann 2017

Course Format:

The course contains six main components:

1. **Lectures:** Deliver the main content of the class. Attendance is mandatory.
2. **Readings:** Accommodate lecture material. The reading assignments are discussed in class. Questions to individual students may be asked.
3. **Quizzes:** In-class quizzes intended to assess the learning progress of the student.
4. **Labs:** Lab sessions offer practice on topics addressed in class and help support homework assignments.
5. **Homeworks:** Homework assignments deepen the lecture material and include coding exercises (skeleton codes are provided in C++).
6. **Projects:** The class is accompanied by a project (teams of 3-4 students) to practice the methods learned in class on a real application. Topics are proposed by the teams and may involve research problems of individual team members.

Grading:

Homework:

35%

Project:

30%

Quizzes:

20%

Labs:

10%

Communal Contributions:

Via the class communication platforms.

The class does not have standard midterm or final exams. The project work involves presentations.

Collaboration and Class Policies:

You are welcome to discuss the course material and homework with others in order to better understand it, but the work you turn in must be your own (with exception of the project where collaborative work is permitted). Any work that is not your own, without properly citing the original author(s), is considered plagiarism. Failure to follow the academic integrity and dishonesty guidelines outlined in the [Harvard Student Handbook](#) will have an adverse effect on your final grade. This includes the removal of copyright notices in code. You may not submit the same or similar work to this course that you have submitted or will submit to another without permission.

The course related material will be distributed through git repositories on the <https://code.harvard.edu/> platform. Membership in the CS205 organization on that platform is required to access the material.

Wk	Tuesday	Thursday	Labs	Events
1(4)	<p>Lecture 1: 2023-01-24</p> <ul style="list-style-type: none"> • Class introduction/organization • Moore's Law • Transistor density and power limit • Parallel computing • Flynn's taxonomy • Overview of parallelism treated in class: DLP, ILP, TLP, shared memory and distributed memory 	<p>Lecture 2: 2023-01-26</p> <ul style="list-style-type: none"> • Computer architecture • von Neumann architecture • Memory pyramid • Linux process anatomy • Introduction to compute cluster: access, job submission • Reading: Leiserson paper 	<p>Sign-up: Select one of the offered lab session days according to your schedule</p>	<p>Note: The "Reading" assignments are relevant for the lecture and due on the day of the lecture! Handouts are typeset in green and deadlines in red. All deadlines are due 11:59 pm.</p> <ol style="list-style-type: none"> 1. Lab section preferences submitted on my.harvard (2023-01-27)
2(5)	<p>Lecture 3: 2023-01-31</p> <ul style="list-style-type: none"> • Cache memories: why are they there, how they work • Cache lines and the 3 C's • What is temporal and spatial locality • Cache associativity: fully, n-way, direct mapped • Memory access patterns (differences row-major / column-major) 	<p>Lecture 4: 2023-02-02</p> <ul style="list-style-type: none"> • Shared memory introduction • Examples of concurrency and concurrent memory access • Why is shared memory programming hard: what is a race condition and why/how does it happen • Quiz 1 	<p>Lab 1: Accessing cluster, SLURM, Linux, compiler and C++ tutorials.</p>	<ol style="list-style-type: none"> 1. HW1 release (2023-01-31)
3(6)	<p>Lecture 5: 2023-02-07</p> <ul style="list-style-type: none"> • Memory model for shared memory programming and its implications on compilers • Sequential consistency • Mutual exclusion / critical sections / locks • Overview of thread libraries 	<p>Lecture 6: 2023-02-09</p> <ul style="list-style-type: none"> • Introduction to OpenMP: why OpenMP and how to use it in new or existing codes • OpenMP: fork/join parallel regions • OpenMP: work sharing constructs • Reading: OpenMP specification 5.1 Chap. 1 (until 1.4 inclusive) 		<ol style="list-style-type: none"> 1. Lab 1 due (2023-02-10) 2. Project team formation due (2023-02-07)
4(7)	<p>Lecture 7: 2023-02-14</p> <ul style="list-style-type: none"> • OpenMP: data environment • OpenMP: synchronization constructs • OpenMP: library routines • OpenMP: environment variables 	<p>Lecture 8: 2023-02-16</p> <ul style="list-style-type: none"> • OpenMP: data environment • OpenMP: synchronization constructs • OpenMP: library routines • OpenMP: environment variables • Quiz 2 	<p>Lab 2: OpenMP locks, critical sections and atomic clauses.</p>	<ol style="list-style-type: none"> 1. HW1 due (2023-02-14) 2. HW2 release (2023-02-14)
5(8)	<p>Lecture 9: 2023-02-21</p> <ul style="list-style-type: none"> • UMA/NUMA memory architectures and processor affinity • What is cache coherency and why is it required in shared memory programming • Cache coherency protocols (focus on MESI) • False sharing 	<p>Lecture 10: 2023-02-23</p> <ul style="list-style-type: none"> • Performance analysis (single node) • Relationship of compute performance (flop) to memory bandwidth • Roofline model • Reading: Williams paper 	<p>Lab 3: False sharing and cache thrashing.</p>	<ol style="list-style-type: none"> 1. Lab 2 due (2023-02-24) 2. Project high-level description due (2023-02-21)

Wk	Tuesday	Thursday	Labs	Events
6(9)	Lecture 11: 2023-02-28 <ul style="list-style-type: none"> Introduction to distributed programming (recap Flynn's taxonomy) What is the Message Passing Interface (MPI) Simple parallel MPI program example 	Lecture 12: 2023-03-02 <ul style="list-style-type: none"> MPI: blocking point-to-point MPI: blocking collective <i>Reading:</i> MPI 4.0 Standard 3.1, 3.2, 3.4, 3.5 		1. Lab 3 due (2023-03-03)
7(10)	Lecture 13: 2023-03-07 <ul style="list-style-type: none"> MPI: non-blocking point-to-point MPI: non-blocking collective <i>Reading:</i> MPI 4.0 Standard 3.7 	Lecture 14: 2023-03-09 <ul style="list-style-type: none"> MPI: I/O file management MPI: I/O read and write routines Parallel I/O for data compression example <i>Quiz 3</i> 	Lab 4: MPI reductions and scans.	1. HW2 due (2023-03-07) 2. HW3 release (2023-03-07)
8(11)	Spring break: 2023-03-14	Spring break: 2023-03-16		
9(12)	Presentations for project proposals: 2023-03-21	Presentations for project proposals: 2023-03-23		1. Lab 4 due (2023-03-24) 2. Project proposals due
10(13)	Lecture 15: 2023-03-28 <ul style="list-style-type: none"> Parallel scaling analysis Strong scaling / Amdahl's law Weak scaling 	Lecture 16: 2023-03-30 <ul style="list-style-type: none"> Instruction set architecture (ISA) / RISC / CISC Processor pipelining (ILP) <i>Reading:</i> Hennessy and Patterson Turing lecture 	Lab 5: Linking your code with third party libraries. Examples for BLAS and LAPACK.	1. HW3 due (2023-03-28) 2. HW4 release (2023-03-28)
11(14)	Lecture 17: 2023-04-04 <ul style="list-style-type: none"> Assembly language (x86-64) Recap Flynn's taxonomy: SIMD Instruction set architecture extensions What is vectorization and why is it important Floating-point operations in x86-64 	Lecture 18: 2023-04-06 <ul style="list-style-type: none"> Memory alignment and relation to cache lines Manual vectorization Intel intrinsics 		1. Lab 5 due (2023-04-07)
12(15)	Presentations for project designs: 2023-04-11	Presentations for project designs: 2023-04-13		1. Project designs due
13(16)	Lecture 19: 2023-04-18 <ul style="list-style-type: none"> Intel intrinsics Compiler auto vectorization Examples for vectorization and performance impact (DLP in roofline) 	Lecture 20: 2023-04-20 <ul style="list-style-type: none"> SPMD programming model Intel ISPC compiler <i>Reading:</i> Pharr paper <i>Quiz 4</i> 	Lab 6: Understanding machine instructions by learning how to debug code.	1. HW4 due (2023-04-18) 2. HW5 release (2023-04-18)

Wk	Tuesday	Thursday	Labs	Events
14(17)	<p>Lecture 21: 2023-04-25</p> <ul style="list-style-type: none"> GPU computing: <ul style="list-style-type: none"> Streaming processors Main difference between CPU and GPU architectures SIMD and SIMD Streaming multiprocessor and Little's Law Introduction to CUDA CUDA warps and threads Class summary 	Reading period: 2023-04-27		<ol style="list-style-type: none"> HW5 due (2023-04-30) Lab 6 due (2023-04-28)
15(18)	Reading period: 2023-05-02	Exam period: 2023-05-04		<ol style="list-style-type: none"> Project deliverables due (TBD) Project final presentations due (TBD)
16(19)	Exam period: 2023-05-09	Exam period: 2023-05-11		