

EC527: High Performance Programming with Multicore and GPUs -- Spring, 2022

Instructor: Martin Herbordt, PHO 333

Office Hours: T3-5, W3-5, other days/times by appointment

Phone: x3-9850 Email: herbordt@bu.edu Web page: <http://learn.bu.edu>

TFs: Sahan Bandara (GTF), Ben Li, Shaivya Gupta (UTFs)

Lab TA Hours: -- TBD Emails: sahanb@bu.edu, liben002@bu.edu, sgupta89@bu.edu

Mission Statement: "Programming for performance using the capabilities of modern processors"

Course Description (catalog): Considers theory and practice of hardware-aware programming. Key theme is obtaining a significant fraction of a program's potential performance through knowledge of the underlying computing platform and how the platform interacts with programs. Studies architecture of, and programming methods for, contemporary high-performance processors. These include complex processor cores, multicore CPUs, and graphics processors (GPUs). Labs include use and evaluation of programming methods on these processors through applications such as matrix algebra and the Fast Fourier Transform.

Prerequisites: Computer Organization (EC413 or equivalent), programming in C, and academic maturity sufficient, e.g., (i) to learn new programming tools from professional documentation, (ii) to design basic experiments, and (iii) perform simple data analysis, including using spreadsheets, plotting data, etc.

Course Motivation:

For several decades programmers found the von Neumann (vN) model to be an adequate worldview for obtaining most of the potential performance from target systems. This model is familiar: instructions are executed serially in a single stream and data are stored in a single image of memory. Instruction executions have uniform performance, as do all memory accesses. Except for specialized processors--DSPs, Supercomputers, MPPs--good vN programming plus a good compiler meant taking advantage of most of a computer's capability.

Current directions in processor architecture—complex memory hierarchies with many layers of cache, superscalar and **deeply pipelined CPUs**, multicore, and accelerators such as AVX extensions and GPUs—have made the vN approach to obtaining performance obsolete. For many applications performance is secondary; but for applications requiring performance, a deeper level of machine understanding is required. The programmer must be aware of the underlying hardware at all stages of software development from algorithm selection and numerical analysis, through coding, to interaction with system tools such as compilers and libraries, and finally debug, tuning, optimization, and maintenance.

Texts and Organization (supplemented with additional articles, lecture notes, and tutorials):

For complete (tentative) readings see "Readings" document.

Part 0 – Methods

- Timing and Timers – See Lab 0 documentation
- Performance Models – Patterson & Hennessy 5e, Chapter 6.1, 6.2, 6.10, 6.11

Part 1 – Single core

This part of the course is based on sections of courses taught at CMU and ETH.

- "How to Write Fast Code," Markus Pueschel, Lecture Notes from CMU and ETH
- "How to Write Fast Numerical Code: A Small Introduction," S. Chellappa, et al.; CMU Tech Report
- *Computer Systems: A Programmer's Perspective*, Bryant & O'Hallaron, Chapter 5 and parts of Chapter 6
- H&P 4th Edition -- Appendix F: Vector processors

Part 2 – Multicore

This part is a condensed and applied version of material from EC713 Parallel Computer Architecture.

- Computer Architecture (Chapter 4): Hennessy & Patterson 4e
- Parallel Computer Architecture (Chapter 5): D. Culler, et al.
- Parallel Programming (Chapters 2 and 3): D. Culler, et al.
- Threads Primer: A Guide to Multithreaded Programming (Chapters 2-5): Lewis & Berg
- OpenMP Lecture Notes: SCV at BU

Part 3 – GPUs

This part is based on a course taught at UIUC by Wen-mei Hwu.

- CUDA Reference Manual(s): NVIDIA
- Programming Massively Parallel Processors: Kirk & Hwu

Course Mechanics

- **Style:** The primary mission of this course is how to create high performance code. Another mission to be a practicum associated with the computer organization and architecture curriculum. For both we explore contemporary high-end processors in some depth and then practice using that knowledge to obtain high resource utilization with real programs. The emphasis is therefore on programming, with lectures in support of the labs. Lectures will also introduce appropriate theory when necessary, especially with respect to performance evaluation.
- **Grading:** Exams: 35%
 Programming/Homework Assignments: 40%
 Final Project: 25%

Please note that these percentages are tentative. Also, that the impact of an assignment/exam grade on the final grade depends on the variance in addition to the percentage.

- **Weekly Programming Assignments:** Until the beginning of the project there will be weekly assignments, 9 in all (0 - 8). All involve programming, mostly exploring small amounts of code in great depth. Some assignments include pencil-and-paper problems in addition to programming.
- **Lab Groups:** Programming assignments have been designed to be done by single students. But by popular demand, you may work in groups of at most two students.
- **Late Policy:** Assignments must be submitted on time, usually on Fridays at midnight. There is a 20% per day penalty. You will get a total of 5 “free” late days to handle special (but common) occurrences such as illness and interviews. Otherwise the only acceptable excuses will be “uncommon” events such as long term disability or a family crisis.
- **Academic (Dis)Honesty versus Collaboration:** You are encouraged to work together to learn the material and to discuss approaches to solving problems. However, *you must come up with and write up the programs and other solutions on your own, or, if you are working with another student, only from you and your partner.*
- **Exams:** There will likely be two mid-term exams. One will be after the multicore part of the course, the second towards the end of the semester. There may also be a final exam.
- **Academic Conduct Code:** Read the academic conduct code. Please note that the penalties are severe.
- **Final Project:** The purpose is to add depth and to practice the concepts learned in an extended case study. Results will be written up conference paper style and presented to the class. For the project only, you may work in teams of up to three students. Please note – the larger the group, the larger the expectations. Group projects work the best when there are distinct parts done by the different group members.

Course Objectives

Review

- Computer architecture including memory hierarchy and basic pipelined CPUs.

Learn about

- Various contemporary high-end processors, in particular, recent CPU cores and memory hierarchy, multicore cache, and GPUs

- Methods of performance evaluation
- Methods of hardware-aware code development
- How to program complex hardware to obtain high utilization

Gain experience with developing efficient programs, including

- using extended instruction sets (AVX) with implicit code and intrinsics
- synchronization
- methods of parallel programming, including PThreads and OpenMP
- cache-aware optimizations
- CUDA for GPU programming (and possibly OpenCL or OpenACC)

From the Course Requisition Form

Basic Goals

1. Students should learn enough about processor architecture and programming to write fast code (code with high utilization of available resources) on contemporary processors.
2. The knowledge and experience should enable students to extend this capability to new processors and to large and varied applications.

Detailed Goals

Students should have a good understanding of theory and practice of

1. Measuring and analyzing performance
2. Developing fast code using methods such as decomposition, mapping, load balancing, blocking, basic block optimizations, and many others
3. Using advanced capabilities such as SIMD vector extensions and use of intrinsics
4. Parallel processing with small-scale (multicore) shared memory processors with both PThreads and OpenMP
5. Programming GPUs, including dealing with the standard inhibitors to getting good performance

Students should also develop a deeper understanding of one of the three technologies (CPU, multicore, GPU) with an extended project. This will consist of examining a more complex numerical or data processing problem.

Course Outcomes

1. Sufficient knowledge of various processor architectures to be able to write high-performance programs
2. Basic knowledge of principles and practice of performance evaluation and writing high-performance programs with the goal of applying this knowledge to other processors.
3. Ability to use AVX instructions set extensions
4. Ability to write parallel programs using PThreads and OpenMP
5. Ability to write GPU programs in CUDA
6. Ability to formulate and design programs at a high level accounting for a target architecture

READINGS & SCHEDULE

Part I – Classes 1-7

L00: Class 1 (1) – Intro, motivation, overview

- “How to Write Fast Numerical Code: A Small Introduction,” S. Chellappa, et al.; CMU Tech Report
- “Programmability: Design Costs and Payoffs Using AMD GPU Streaming Languages and Traditional Multi-Core Libraries,” Weber, SAAHPC, Extended Abstract and Lecture Slides
- “How to Write Fast Code,” Markus Pueschel, Lecture Slides, Class 1
- “Debunking the 100x GPU versus CPU Myth: An evaluation of throughput computing on CPU and GPU,” V.W. Lee, et al. ISCA10.

L01a, L01b: Class 2-3 (2) – Performance models, review memory hierarchy, Memory-aware optimizations

- P&H 5e Chapter 6.1, 6.2, 6.10, 6.11 on performance models
- “How to Write Fast Numerical Code: A Small Introduction,” S. Chellappa, et al.; CMU Tech Report
- Basic memory hierarchy from any standard computer architecture textbook
- *Computer Systems: A Programmer’s Perspective*, Bryant & O’Hallaron, Chapter 6 (parts).

L02a, L02b: Classes 4-5 (2) -- CPU-aware optimizations, compiler interactions

- *Computer Systems: A Programmer’s Perspective*, Bryant & O’Hallaron, Chapter 5.
- “How to Write Fast Numerical Code: A Small Introduction,” S. Chellappa, et al.; CMU Tech Report

L03a, L03b: Class 6-7 (2) – Vector Processing, SIMD, Data Parallel Programming, Vector extensions

- Slides from David Patterson
- “How to Write Fast Numerical Code: A Small Introduction,” S. Chellappa, et al.; CMU Tech Report
- H&P 5e Appendix G, Vector Processors
- “How to Write Fast Code,” Markus Pueschel, Lecture Slides, Classes 13-14
- “Introduction to SSE Programming,” Alex Fr, The Code Project

Part II – Classes 8-14

L04: Classes 8,10 (2) – Parallel Programming

- “Type Architectures,” L. Snyder, *Annual Review of Computer Science*, 1986
- *Optimizing Compilers for Modern Architectures: A Dependence-based Approach*, Allen & Kennedy, Chapter 2
- *Parallel Computer Architecture -- Draft*, Culler, et al.: Chapter 2, “Parallel Programming”
- *Parallel Computer Architecture -- Draft*, Culler, et al., Chapter 3, Sections 3.1, 3.2.2, 3.4.1, “Programming for Performance”
- *Computer Organization and Design*, Patterson and Hennessy: Chapter 6, “Parallel Processing”

L05: Class 9 – Programming with Threads, Review Concurrency and Synchronization

- Slides from H. Casanova, U Delaware, and U Hawaii
- “POSIX Threads Programming,” B. Barney, LLNL

L06, L07: Class 11 – Review Concurrency and Synchronization

- “Operating Systems Concepts” by Peterson and Silberschatz, Chapter 9

L08: Classes 12,14 (1.5) – Cache for Shared Memory

- *Parallel Computer Architecture*, by Culler, et al. Chapter 5, Sections 5.1-5.4

L09: Class 13 – OpenMP

- Slides from BU SCV and other sources

L10: Class 14 (0.5) – Synchronization implementation for shared memory processors

- *Parallel Computer Architecture*, by Culler, et al. Chapter, Section 5.5

Part III – Classes 15,17-21

L11-L16: Classes 15,17-21 (6) – NVIDIA GPUs, CUDA, Performance optimization, case studies

- Slides from Kirk and Hwu
- *Programming Massively Parallel Processors*, by Kirk and Hwu
- Various case studies, esp. from Molecular Dynamics

Class 16 – Mid-Term 1

Class 22,25 – Group meetings with instructor

Class 23 – Review

Class 24 – Mid-Term 2

Class 26-28 – Project Presentations

Exam dates are tentative, Lab dates are fixed
Labs will be posted at least a week before they are due

Wk	Cl	Date	Lecture Topic	HW	LAB
1	1	24-Jan	Intro, administration, goals & expectations Motivation		Lab 0: Timers, CPE, roofline, etc.
	2	26-Jan	Performance models, Roofline model, Review memory hierarchy		
2	3	31-Jan	Memory-aware optimizations	Lab 0 due	Lab 1: Memory
	4	2-Feb	Intel processor assembly language CPU-aware optimizations	Lab 1 due 2/4	
3	5	7-Feb	CPU-aware optimizations		Lab 2: Pipelines
	6	9-Feb	Intro to dataparallel, SIMD, and Vector processing	Lab 2 due 2/10	
4	7	14-Feb	Programming using AVX		Lab 3: Vectors
	8	16-Feb	Intro to parallel programming, part 1	Lab 3 due 2/18	
5	9	22-Feb	Thread-based programming and Pthreads		Lab 4: Threads tutorial
	10	23-Feb	Intro to parallel programming, part 2	Lab 4 due 2/25	
6	11	28-Feb	Review Concurrency		Lab 5: Threads programming
	12	2-Mar	Multicore Cache -- DMA, state machines, snooping, invalidate/update, protocols	Lab 5 due 3/4	
Spring Break					
7	13	14-Mar	OpenMP		Lab 6: OpenMP
	14	16-Mar	Multicore synchronization implementation	Lab 6 due 3/18	
8	15	21-Mar	GPU -- First pass	Mid-term week	
	16	23-Mar	Mid-term		
9	17	28-Mar	GPU -- Thread Organization		Lab 7: GPU threads (SOR)
	18	30-Mar	GPU -- Memory Organization		
10	19	4-Apr	GPU -- catch up	Lab 7 due 4/8	
	20	6-Apr	GPU -- In Depth, processor		Lab 8: GPU memory (MMM)
11	21	11-Apr	GPU -- In Depth, memory, part 1		
	22	13-Apr	GPU -- In Depth, memory, part 2, miscellaneous	Lab 8 due 4/15	
12		18-Apr	Patriots Day		
	23	20-Apr	Project Meetings		
13	24	25-Apr	Project Meetings		
	25	27-Apr	Project Meetings		
14	26	2-May	Project presentations		
	27	4-May	Project presentations		
15		10-May	Project Writeups Due -- 12:00 Noon		