

# CAS CS 591 A1: Data Systems Architectures

Boston University

Spring 2019

## Class Syllabus

**Course Description:** Data is everywhere. As scientists, users, and citizens we are both generating and exploiting large, ever-growing, diverse sets of data. For several applications – ranging from scientific discovery to business analysis, governance, and every-day activities – we are directly using and indirectly affecting hundreds of data systems! The big challenge is to turn data into useful knowledge, and to do so quickly, in order to increase the impact of the new insights. Achieving these goals comes with a number of technical challenges. How to exploit the continuously evolving hardware (storage, computation, network)? How to collect all incoming data efficiently? How to query dynamic collections of data that keep accumulating incoming data? How to parallelize query processing from one core to a few (scale-up), and then to thousands (scale out)? What are the needs of evolving workloads (hybrid transactional/analytical processing, graph analytics, Internet-of-Things, micro-payments, monitoring)? In this course, we will discuss how to design data systems that can address these challenges. We will see in detail the two driving forces behind innovation in data systems: hardware and workloads, and we will discuss recent and future trends of both. We will use examples from several data management areas including relational systems, distributed database systems, key value stores, newSQL and NoSQL systems, data systems for machine learning (and machine learning for data systems), interactive analytics, and data management as a service. In a quickly moving industry and research landscape, such skills are essential.

**Prerequisites:** The class requires familiarity with database systems at the level CS460/660, and with algorithms, data structures, computer systems, and system programming at the level of CS210 or CS350. Please see the instructor if you are not sure about the level of your preparation.

**Instructor:** Manos Athanassoulis ([mathan@bu.edu](mailto:mathan@bu.edu))

office hours: Tu/Th 2-3pm

office: MCS 279

**Teaching Assistants:** Subhadeep Sarkar ([ssarkar1@bu.edu](mailto:ssarkar1@bu.edu))

### Meeting Times and Places

lectures: Tu/Th, 12:30-1:45 pm, MCS B23

**Course Website:** <http://manos.athanassoulis.net/classes/CS591/>

**Textbooks (not required):** There is no textbook that covers cutting edge research, however, the data management community has produced top quality textbooks that can serve as reference to provide background material. The class is based on recent research papers which will be available to you through the BU network.

- R. Ramakrishnan and J. Gehrke. *Database Management Systems*. Third Edition. McGraw-Hill 2002.

An excellent collection of classic papers in the database field is the following:

- [Readings in Database Systems](#). P. Bailis, J. Hellerstein, M. Stonebraker, editors.

Other good background material is:

- Architecture of a Database System. By J. Hellerstein, M. Stonebraker and J. Hamilton. Foundations and Trends in Databases, 2007
- The Design and Implementation of Modern Column-store Database Systems. By D. Abadi, P. Boncz, S. Harizopoulos, S. Idreos, S. Madden. Foundations and Trends in Databases, 2013
- Massively Parallel Databases and MapReduce Systems. By Shivnath Babu and Herodotos Herodotou. Foundations and Trends in Databases, 2013

### Course Requirements

*Paper Presentation:* After the initial 4 classes (2 weeks) all students will take turns presenting papers. In each class we will discuss one or two main paper(s) (and there will be a few background papers), and each student will present once in the semester, either alone or as a group of two students). The student(s) presenting will be responsible to outline strong and weak points of the paper and propose at least one idea for improving the approach presented in the paper. **All students** will read the presented paper.

*Paper Reviews:* All students should read all papers and provide a review. *Reading the paper and writing a review is very important to help the students prepare for the class presentation and discussion.* Every student is expected to deliver a **long review for 5 papers** and a **short review for all other papers**. Throughout the semester each student can skip 3 papers from reviewing them. Every review for a given paper has to be handed-in **before the class**, having the class starting time as a hard deadline.

A **short review** consists of two paragraphs (up to half a page), (a) the first outlining what is the problem and why is it important, and (b) the second presenting the key ingredients of the approach taken.

A **long review** consists of a few paragraphs answering the following questions: (i) what is the problem, (ii) why it is important, (iii) why is it hard, (iv) why older approaches are not enough, (v) what is key idea and why it works (a list of at least three key points), (vi) what might be missing and how can we improve this idea (a list of at least three key points), (vii) an evaluation as to whether the paper supports its claims, and (viii) possible next steps of the work presented in the paper. The ideal size of the review is about 1 page, single column, 10pt font, 1 inch margin (and it may only exceed 1 page if the student wants to elaborate on how to improve the ideas on the paper). Every paper should have at least two reviewers.

*Project:* Finally, this class requires a semester-long project and a final report in the style of a

conference paper. The project will be either implementation-heavy or research-oriented. Students will work in teams of 3-4 and after the first two weeks each team will have been associated with a specific project. Students can propose their own research project upon approval by the instructor.

The overall grade will be based on the following policy:

- Class participation: 5%
- Long Paper reviews: 15%
- Short Paper reviews: 10%
- Paper presentation: 25%
- Mid-semester project progress report: 10%
- Project: 35%

**Topics:** Throughout the class we will cover data systems design principles from the following different angles.

1. What affects new data systems designs (data and applications, emerging hardware, and new workloads)
2. Traditional Data Systems for Modern Hardware
3. Distributed Database Systems
4. Scale-out Systems: from Map-Reduce to SQL-on-Hadoop
5. NoSQL, NewSQL and Key-Value stores

### Important Dates for all classes

February 4<sup>th</sup>, last day to add a class

February 26<sup>th</sup>, last day to drop (without a “W”)

### Important Dates for CS591

February 8<sup>th</sup>, last day to select a project

March 8<sup>th</sup>, submit the mid-way project report

April 28<sup>th</sup>, submit final project report

### Tentative Schedule

Week 1	<b>Lecture 1</b> (1/22)	<b>Introduction to Data Systems and CS591</b>
	<b>Lecture 2</b> (1/24)	<b>Data Systems Architectures Essentials – Part 1</b>
Week 2	<b>Lecture 3</b> (1/29)	<b>Data Systems Architectures Essentials – Part 1</b>
	<b>Lecture 4</b> (1/31)	<b>Class Project Overview</b>
Week 3	<b>Lecture 5</b> (2/5)	<b>Storage Layouts: Row-Stores vs. Column-Stores</b>
	<b>Lecture 6</b> (2/7)	<b>Storage Layouts: Adaptive &amp; Hybrid Layouts</b>
Week 4	<b>Lecture 7</b> (2/12)	<b>New Hardware: Data Systems for Flash &amp; SMR Disks</b>
	<b>Lecture 8</b> (2/14)	<b>New Hardware: Data Systems for Multi-Core</b>
Week 5	2/19 – No Class, Monday replacement day	
	<b>Lecture 9</b> (2/21)	<b>Indexing A: B+-Trees, Bitmaps, Hash-Index</b>
Week 6	<b>Lecture 10</b> (2/26)	<b>Indexing B: Access Path Selection</b>
	<b>Lecture 11</b> (2/28)	<b>Modern Storage Engines: HTAP Systems</b>
Week 7	<b>Lecture 12</b> (3/5)	<b>Modern Storage Engines: Log-Structured Merge Trees</b>

	<b>Research talk 1</b> (3/7)	<b>Widening the LSM design space</b> (visiting lecture)
Week 8	No Class – Spring Break	
Week 9	<b>Lecture 13</b> (3/19)	<b>Indexing C:</b> Data Skipping
	<b>Lecture 14</b> (3/21)	<b>Indexing D:</b> Adaptive Indexing
Week 10	<b>Research talk 2</b> (3/26)	<b>Fast &amp; Accurate Time Series Clustering</b> (visiting lecture)
	<b>Research talk 3</b> (3/28)	<b>Storage &amp; Indexing for Data Series</b> (visiting lecture)
Week 11	<b>Lecture 15</b> (4/2)	<b>In-Situ Data Processing:</b> Efficiently Accessing Raw Data Files
	<b>Lecture 16</b> (4/4)	<b>Scientific Databases:</b> Multi-dimensional Data Management
Week 12	<b>Lecture 17</b> (4/9)	<b>Distributed Databases:</b> Global Distributed Systems
	<b>Lecture 18</b> (4/11)	<b>Map/Reduce:</b> Data Management at Scale
Week 13	<b>Lecture 19</b> (4/16)	<b>Data Systems for ML:</b> Data Processing Primitives for ML
	<b>Lecture 20</b> (4/18)	<b>ML for Data Systems:</b> Automatic Data System Design
Week 14	<b>Lecture 21</b> (4/23)	<b>Indexing E:</b> Learned Indexes
	<b>Lecture 22</b> (4/25)	<b>Indexing F:</b> The Periodic Table of Access Methods
Week 15	<b>Lecture 23</b> (4/30)	<b>Project Presentations A</b>
	<b>Lecture 24</b> (5/2)	<b>Project Presentations B</b>

### Collaboration Policy

You are strongly encouraged to collaborate with one another in studying the lecture materials and preparing for reviews and presentations.

You may discuss ideas and approaches to the projects with others (provided that you acknowledge doing so in your solution), but such discussions should be kept at a high level, and should not involve actual details of the code or of other types of answers. **You must complete the actual solutions on your own.**

### Academic Misconduct

We will assume that you understand BU's Academic Conduct Code:

<http://www.bu.edu/academics/policies/academic-conduct-code>

Prohibited behaviors include:

- copying all or part of someone else's work, even if you subsequently modify it; this includes cases in which someone tells you what you should write for your solution
- viewing all or part of someone else's work
- showing all or part of your work to another student
- consulting solutions from past semesters, or those found online or in books
- posting your work where others can view it (e.g., online).

Incidents of academic misconduct will be reported to the Academic Conduct Committee (ACC). The ACC may suspend/expel students found guilty of misconduct. ***At a minimum, students who engage in misconduct will have their final grade reduced by one letter grade (e.g., from a B to a C).***