

CAS CS 551 Streaming and Event-driven Systems

Course Syllabus Spring 2023

Instructor Name: Vasiliki Kalavri

Instructor Office Hours: Tue-Thu 5-6:30pm (CCDS 713, 7th floor).

Lectures: Tue-Thu 3:30 pm-4:45 pm, [MCS B29](#)

Discussions: Mon 9:05 am-9:55 am, [CAS 116](#)

Teaching Fellow: Emmanouil (Manos) Kritharakis

TA Office Hours: Tue 11:00 am-1:00pm, Open space 6th floor CCDS

IMPORTANT: Refrain from using email to reach the course staff. To contact the instructor or TA, send a **private Piazza post**.

1. Overview

1.1 Course Description

Modern data-driven applications increasingly require continuous, low-latency processing of large-scale, rapid data events such as clicks, search queries, online interactions, financial transactions, traffic records, and sensor measurements. Distributed stream processing has become highly relevant to industry and academia due to its capabilities to both improve established data processing tasks and to facilitate novel applications with real-time requirements. In this course, you will learn how to design, implement, and evaluate scalable and reliable stream processing and event-driven applications.

Specifically, we will cover the following topics:

- Publish/Subscribe systems
- Architecture of distributed stream processing systems
- Dataflow programming
- Fault-tolerance and processing guarantees
- Streaming state management
- Windowing semantics and optimizations
- Complex event processing
- Microservice architectures
- Serverless functions and their relationship to stream processing

1.2 Course Objectives

At the end of the course, successful students will have gained skills and hands-on experience on the following methods and technology:

- Design and implementation of dataflow stream processing applications

- Message queues, log-based message brokers, and publish/subscribe systems
- Ability to comprehensively compare the features, architecture, and processing guarantees of modern streaming systems
- Implementing, deploying, and evaluating event-based applications with Apache Flink and Apache Kafka.
- Operations for scalable and reliable stream processing, including logging, monitoring, debugging, and upgrading streaming applications
- A solid understanding of the challenges and trade-offs one needs to consider when designing and deploying streaming applications

Further, students will be exposed to recent developments in stream processing research through paper assignments and presentations. The collaborative semester-long project will prepare them for the practical aspects of their future careers and expose them to project management tools and software engineering best practices.

1.3 Prerequisites

CAS CS 112 and CAS CS 210; CAS CS 451 and CAS CS 460 or consent of instructor.

To be successful in this course, students will need to have strong programming skills, a solid understanding of Computer Systems fundamentals (CS 210) and some prior experience with object-oriented programming / Java (CS 211). Familiarity with Distributed Systems (CS 451) and Database Systems (CS 460) is highly recommended.

2. Instructional Format, Course Pedagogy, and Approach to Learning

2.1 Courseware

- We will use the **course website** to maintain an up-to-date class schedule:
<https://vasia.github.io/cs551/index.html>
- We will use **Piazza** for announcements, questions, discussions, and all other communication:
<https://piazza.com/bu/spring2023/cascs551>
- We will use **Gitlab** for the hands-on sessions, discussions, and the semester projects:
<https://cs551-gitlab.bu.edu/>
 - Sign up for an account using your BU email.
 - Once approved, you will be able to login and create projects.
- We will use **Gradescope** for assignment submissions:
<https://www.gradescope.com/courses/495250>

2.2 Lectures

Lectures will be held during the assigned time slots. Section 4 of the syllabus provides the topic and assigned readings for each lecture. You are expected to complete the readings before the day of the lecture and actively participate in class discussions. Lecture slides will be made available on the class website prior to the lectures or shortly after.

2.3 Discussions

Students are expected to attend the weekly discussion section they have been assigned to. The Teaching Fellows will lead the discussion sessions. The objectives are: to present material on the required tools such as Apache Flink, Apache Kafka, and Stateful Functions, that reinforce the concepts covered in the lectures, and answer questions (or provide clarifications) regarding the assignments and projects. The Teaching Fellow will post information to Piazza as necessary. In addition to the discussions, the Teaching Fellow will hold weekly Office Hours.

2.3.1 Software requirements

During the discussion sessions, you will solve a set of programming exercises using Apache Flink and Apache Kafka. Use your own laptop or desktop computer and make sure to set up your environment correctly as described below.

You can develop and execute Flink applications on Linux, macOS, and Windows. However, UNIX-based setups have complete tooling support and are generally preferred by Flink developers. **All assignments assume a UNIX-based setup.** If you are a Windows user, you are advised to use Windows subsystem for Linux (WSL), Cygwin, or a Linux virtual machine to run Flink in a UNIX environment.

To setup and run Flink, you additionally need:

- A Java 8 or 11 installation. To develop Flink applications and use its DataStream API in Java or Scala you will need a Java JDK. A Java JRE is not sufficient!
- Apache Maven 3.x. Flink provides Maven archetypes to bootstrap new projects.
- An IDE for Java and/or Scala development. Common choices are IntelliJ IDEA, Eclipse, or Netbeans with appropriate plugins installed. We recommend IntelliJ IDEA.

Even though Apache Flink is a distributed data processing system, you will typically develop and run initial tests on your local machine. This makes development easier and simplifies cluster deployment, as you can run the exact same code in a cluster environment without making any changes.

2.4 Classroom recordings

Class sessions might be recorded for the benefit of registered students who are unable to attend live sessions (either in person or remotely) due to illness. Recorded sessions will be made available to registered students ONLY via their password-protected BU account. Students may not share such sessions with anyone not registered in the course and may certainly not repost them in a public platform. Students have the right to opt-out of being part of the class recording. Please consult the following site for further details: <https://digital.bu.edu/lfa-classroom-recordings>.

2.5 Course Materials

There is no required textbook for this class. Slides, lecture notes, and other publicly available resources will be published on the course website and on Piazza. A list of readings is provided in the course website: <https://vasia.github.io/cs551/readings.html>. You should be able to access

all readings when connected to the campus network. Please contact the instructor if any of the listed readings is unavailable or inaccessible.

3. Assignments and Grading Criteria

3.1 Semester Project

This class is highly collaborative and research-oriented. During the first week, you will be provided with a list of semester projects and you will be asked to select your top-3 preferences. You will then be assigned to a project team with 3-5 students. During the semester, the team will be working together to deliver:

- A **design document** outlining (1) the project goals, (2) an implementation and evaluation plan, (3) the task distribution among team members.
- A **midterm project demo**. Demos will be presented during class time.
- A **final demo** to be presented during the last week of class.
- The project's **gitlab repository**, including code, tests, automation and plotting scripts, and documentation.

3.2 Paper assignments and guest lectures

During the semester, we will read and discuss various technical papers. For each paper, you will be asked to submit:

1. a short summary, describing the core ideas of the paper
2. a list of questions to be discussed during lecture time. **These deliverables are individual.**

We will also host 2 guest lectures. You are expected to participate in the guest lectures by asking questions.

3.3 Grading Scheme

Your final grade will be determined as follows:

1. Participation & effort (**20%**):
 - In-class participation.
 - Discussion participation.
 - Piazza contributions.
 - Git activity (project + discussions).
 - Group Meeting minutes.
 - Office Hours participation.
2. Paper readings & guest lectures (**20%**):
 - Paper summaries, questions (10%)
 - Paper discussion participation (5%)
 - Guest lecture participation (5%)

3. Semester project (**60%**) (in teams of 3-5 students):

- Design document (maximum 3 pages) 5%.
- Midterm demo 20%.
- Final demo and deliverables: 35%.

The final deliverables include (1) the full code implementing the project tasks as defined in the project design document, (2) auxiliary code for data pre-processing, deployment, and testing, (3) complete supporting documentation.

Individual contributions to collaborative assignments will be assessed by taking into account the following:

- The quality of individual task deliverables outlined in the project design document.
- The individual's ability to answer questions about the project during demo presentations and office hours.
- The individual's performance during the paper and demo presentations.
- The individual's contribution to the project's gitlab repository (git history).

There is no final exam at the end of the course.

4. Class and University Policies

4.1 Homework submission

All assignments and the project deliverables will be submitted via the course Gitlab. All deliverables are **due by 11:59pm on the day of the respective deadline**.

4.2 Attendance

Students are expected to attend each class session unless they have a valid reason for being absent. Acceptable excused absences include observing religious holidays and illness. In such cases, students are advised to contact the instructor as soon as possible, so that reasonable accommodations can be provided. Please review the **BU attendance policy** and the **BU Policy on Religious Observance** for more information.

4.3 Late work policy

Students who submit homework late will only be eligible for up to **50% of the original score**.

4.4 Academic conduct

Academic standards and the code of academic conduct are taken very seriously at our university. Please take the time to review the CAS Academic Conduct Code: <http://www.bu.edu/academics/resources/academic-conduct-code/> and the GRS Academic Conduct Code: <http://www.bu.edu/cas/students/graduate/grs-forms-policies-procedures/academic-discipline-procedures/>. Please review the sections regarding plagiarism and cheating carefully. Copies of the CAS Academic Conduct Code are also available in room CAS 105. A

student suspected to violate this code will be reported to the Academic Conduct Committee, and if found culpable, the student will receive a grade of "F" for the course

All assignments must be completed individually, unless instructed otherwise. Discussion with fellow students via Piazza or in-person are encouraged, but presenting the work of another person as your own is expressly forbidden. This includes "borrowing", "stealing", copying programs/solutions or parts of them from others. Note that we may use an automated plagiarism checker. Cheating will not be tolerated under any circumstances.

Any resources, including material from other students (current or past), that are used, beyond the text or that provided by the TF or professor must be clearly acknowledged and attributed. Using such material may at the discretion of the TF or professor result in a lower grade. However, if such material is used and not acknowledged and 12 attributed, it will automatically be considered as possible academic misconduct.

5. Accommodations

If you are a student with a disability or believe you might have a disability that requires accommodations, please contact the Office for Disability Services (ODS) at (617) 353-3658 or access@bu.edu to coordinate any reasonable accommodation requests. ODS is located at 25 Buick Street on the 3rd floor.

6. Detailed Schedule

The rest of the syllabus is tentative and might be updated during the semester. We will be keeping you informed of any changes made to the readings or assignment deadline via Piazza.

Make sure to become familiar with the **Official Semester Dates**. Some of the critical Semester Dates are:

- The Last Day to DROP Classes (without a 'W' grade) February 23, 2023.
- The Last Day to DROP Classes (with a 'W' grade) March 31, 2023.

Date	Topic	Readings	Assignment
1/19	Introduction to stream processing	[1]	
1/23	Disc #1: How to read a paper		
1/24	Publish/Subscribe systems	[2]	
1/26	Paper 1: Stream ingestion & indexing	[3]	Projects announced
1/30	Disc #2: Intro to Flink		
31/1	Dataflow stream processing systems	[4]	
2/2	Paper 2: Realtime data processing	[5]	Project selection due
2/6	Disc #3: Intro to Kafka		
2/7	Notions of time	[6]	
2/9	Paper 3: Tracking computation progress	[7]	
2/13	Disc #4: Writing Flink programs & DataStream API		
2/14	Windowing semantics	[8]	
2/16	Paper 4: Window aggregation	[9]	Design document due
2/21	Disc #5: Windows & event-time		
2/23	Streaming state management	[10]	
2/27	Disc #6: State management		
2/28	Paper 5: Consistent regions	[11]	
3/2	Distributed snapshots	[12]	

3/13	Disc #7: Flink + ML		
3/14	Guest Lecture: Confluent @ CCDS 701		
3/16	Paper 6: Exactly-once fault tolerance	[13]	
3/20	Disc #8: Checkpoints		
3/21	NO CLASS		Midterm presentation due
3/23	Midterm project presentations (Teams 1, 2, 3a, 4)		
3/27	Disc #9: Reconfiguration & upgrading		
3/28	Midterm project presentations (Teams 3b, 6, 8)		
3/30	Guest Lecture - Materialize @ CCDS TBD		
4/3	Disc #10: Metrics & monitoring		
4/4	Flow control & backpressure	[14]	
4/6	Elasticity & state migration	[15]	
4/10	Disc #11: Flink Stateful Functions		
4/11	Paper 7: State migration	[16]	
4/13	Stream query optimization	[17]	
4/18	NO CLASS		
4/19	Disc #12: Project hacking		
4/20	Stateful functions		
4/24	Disc #13: Project hacking		
4/25	Emerging topics in data stream processing		
4/26	Final demos due		
4/27	Demo presentation (teams 1-4)		
5/1	Disc #14: Project hacking		
5/2	Demo presentation (teams 5-8)		Final project repositories due

Readings

- [1] Streaming 101: <https://www.oreilly.com/radar/the-world-beyond-batch-streaming-101/>
- [2] The many faces of publish/subscribe: <https://dl.acm.org/doi/pdf/10.1145/857076.857078>
- [3] **Data Ingestion for the Connected World:**
https://people.csail.mit.edu/tatbul/publications/sstore_cidr17.pdf
- [4] Streaming 102: <https://www.oreilly.com/radar/the-world-beyond-batch-streaming-102/>
- [5] **Realtime Data Processing at Facebook**
<https://research.facebook.com/publications/realtime-data-processing-at-facebook/>
- [6] Flexible time management in data stream systems:
<https://dl.acm.org/doi/pdf/10.1145/1055558.1055596>
- [7] **Watermarks in Stream Processing Systems: Semantics and Comparative Analysis of Apache Flink and Google Cloud Dataflow** <https://www.osti.gov/servlets/purl/1823361>
- [8] SECRET: a model for analysis of the execution semantics of stream processing systems:
<https://dl.acm.org/doi/pdf/10.14778/1920841.1920874>
- [9] **Efficient Window Aggregation with General Stream Slicing:**
https://openproceedings.org/2019/conf/edbt/EDBT19_paper_171.pdf
- [10] State Management in Apache Flink: <https://dl.acm.org/doi/10.14778/3137765.3137777>
- [11] **Consistent regions: guaranteed tuple processing in IBM streams:**
<https://sariyuce.com/papers/vldb16.pdf>
- [12] Lightweight Asynchronous Snapshots for Distributed Dataflows:
<https://arxiv.org/pdf/1506.08603.pdf>
- [13] **MillWheel: fault-tolerant stream processing at internet scale:**
<https://research.google/pubs/pub41378/>
- [14] A Survey on the Evolution of Stream Processing Systems:
<https://arxiv.org/pdf/2008.00842.pdf>
- [15] Three steps is all you need: fast, accurate, automatic scaling decisions for distributed streaming dataflows: <https://www.usenix.org/system/files/osdi18-kalavri.pdf>
- [16] **Meces: Latency-efficient Rescaling via Prioritized State Migration for Stateful Distributed Stream Processing Systems**
<https://www.usenix.org/system/files/atc22-gu-rong.pdf>
- [17] A catalog of stream processing optimizations: <https://dl.acm.org/doi/10.1145/2528412>