# Big Data Technologies: Course Information

**Instructor**: Joseph Rosen
**Email**: jrosen@iit.edu (preferred contact method)
**Lecture:** Thursday, 6:45– 9:36 PM, online
**Telephone**: 312-860-0860 (m) (by appointment or for emergencies only)
**Office Hours**: Tuesday, 5:15 - 6:15 PM by appointment. Other times available by appointment.

## Course Description

Big data is the area of informatics focusing on datasets whose size is beyond the ability of typical database and other software tools to capture, store, analyze and manage. This course provides a rapid immersion into the area of big data and the technologies which have recently emerged to manage it.  We start with an introduction to the characteristics of big data and an overview of the associated technology landscape and continue with an in depth exploration of Hadoop, the leading open source framework for big data processing. Here the focus is on the most important Hadoop components such as Hive, Pig, stream processing and Spark as well as architectural patterns for applying these components. We continue with an exploration of the range of specialized (NoSQL) database systems architected to address the challenges of managing large volumes of data. Overall the objective is to develop a sense of how to make sound decisions in the adoption and use of these technologies. Prerequisites: CS 425 or equivalent.

## Required Texts

Tom White. 2015.  Hadoop: The Definitive Guide: Storage and Analysis at Internet Scale
 (4th ed.). O'Reilly Media, Inc (TW)

Pramod J. Sadalage and Martin Fowler. 2012. *NoSQL Distilled: A Brief Guide to the Emerging World of Polyglot Persistence*. Addison-Wesley.(PS)

## Free Online

Jimmy Lin and Chris Dyer. 2010. *Data-Intensive Text Processing with Mapreduce*. Morgan and Claypool Publishers. https://vgc.poly.edu/~juliana/courses/BigData2014/Textbooks/MapReduce-algorithms-Jan2013-draft.pdf

Anand Rajaraman and Jeffrey David Ullman. 2011. *Mining of Massive Datasets*. Cambridge University Press, New York, NY, USA. http://infolab.stanford.edu/~ullman/mmds/book0n.pdf

# Course Notes

Copies of each week's course lecture notes will be posted on the Blackboard system before class starts each week. Students should note that the lecture notes are meant to frame the lecture and class discussion, and alone will not provide the depth of knowledge required to successfully achieve the course objectives. Students are advised to take their own notes as necessary.

# Other Readings

Readings will be from published research or industry online materials and available via Blackboard "Articles" or "Free Books and Chapters."

# Equipment / Material Requirements

Students are required to have a 64-bit laptop/PC running Windows, MacOS, or Linux.

Students will also be required to set up an account on the Amazon Web Services cloud to install and configure a Hadoop environment. I will provide detailed instructions about this.

Students can try, on their own, to install the Hadoop sandbox on their laptop/PC. Often this will not work well, especially on Windows PCs due to lack of RAM. For ease of support and consistency, all assignments will be described and tested against the Azure environment.

# Grading Policy

| Short Quizzes | 15% |
| --- | --- |
| Assignments | 15% |
| Project / Paper Proposal | 5% |
| Project / Paper Draft | 5% |
| Project / Paper | 25% |
| Half Term Exam | 35% |

Each short quiz will be open notes and books (but closed Google Search) and consist of multiple choice and short answer questions and should take no more than 10-15 minutes to complete. The lowest short quiz score will be dropped when calculating the overall quiz score. Since you can miss a quiz with no penalty there are no make ups or exceptions (except as explicitly provided for by IIT policy or if negotiated ahead of time with me).

The half term exam will be open notes and books (again closed Google Search) and consist of multiple choice and short answer questions, some longer essay questions and take 90-120 minutes. Contact me well ahead of time for possible accommodations if you will be unable to attend on the day of the exam.

Assignments, project/paper proposals and project/paper drafts must be submitted by 6:45pm on their announced due dates (or whatever makes sense for your time zone). You have up to 3 grace days to apply across all of these without penalty. For example, one assignment can be up to 3 days late or three assignments can each be up to one day late. Beyond this grace period 5% will be deducted from your score for each day the assignment is late. But contact me well ahead of time for possible accommodations.

Without prior negotiated accommodations the final project or paper must be submitted by 11:59pm on its due date. Beyond this 10% will be deducted from your score for each day the assignment is late.

Assignments may be reviews of research papers or projects (with write ups) applying big data technology.

By the end of term you will have completed a paper exploring some topic in big data technology more deeply (ex. security, architecture, specific tools, etc.) or conducted an investigational project where you will have applied big data technology to a problem of interest to you, your community or an organization with which you are affiliated.

# Grade Distribution

$A = 100 - 90, \ B = 89.9999 - 75, \ C = 74.9999 - 60$ (Undergraduate only)

# Students with Disabilities

Reasonable accommodations will be made for students with documented disabilities. In order to receive accommodations, students must obtain a letter of accommodation from the Center for Disability Resources. The Center for Disability Resources (CDR) is located in 3424 S. State St., room 1C3-2 (on the first floor), telephone 312 567.5744 or disabilities@iit.edu.

# Course Syllabus

For guidance, some order of the following may change a bit.

| Module | Description | Details |
|--------|-------------|---------|
| 1a | <ul><li>Course introduction</li><li>What is big data?</li></ul> | <ul><li>What is Big Data?</li><li>Characteristics of Big Data</li><li>Understanding Big Data with Examples</li><li>Big Data Processing Pipeline</li><li>Big Data Trends</li><li>The Risks of Big Data</li></ul> |
| 1b | <ul><li>Distributed systems concepts</li><li>Fault tolerance in distributed systems</li></ul> | |
| 2a | <ul><li>Hadoop Basics</li></ul> | <ul><li>Describe the case for Hadoop</li><li>Identify the Hadoop Ecosystem architecture</li></ul> |

| | | |
|---|---|---|
| | | • Data Management - HDFS, YARN<br>• Data Processing: MapReduce, Spark<br>• Data Access - Pig, Hive, HBase, Storm, Solr<br>• Data Governance & Integration - Falcon, Flume, Sqoop, Kafka, Atlas<br>• Security - Kerberos, Knox<br>• Operations - Ambari, Zookeeper, Oozie |
| 2b | Hadoop Distributed File System | • Concepts<br>• Operations<br>• Command Line Interface<br>• Java Interface<br>• Fault tolerance |
| 3 | YARN<br>MapReduce I | • YARN Concepts<br>• YARN Architecture |
| 4 | MapReduce II | • Introduction to Developing Hadoop Applications<br>• Illustrate the MapReduce model conceptually<br>• Discuss how MapReduce works at a high level<br>• Define how data flows in MapReduce<br>• Job Execution Framework<br>• Describe how jobs execute in YARN<br>• Describe how to manage jobs in YARN<br>• Write a MapReduce Program<br>• Design and implement the Mapper class, Reducer class and driver<br>• Build and execute the code then examine the output<br>• Describe data set for programming problem<br>• Use the MapReduce API<br>• API overview<br>• Mapper input processing and Reducer output processing data flow<br>• Explore the Mapper, Reducer and Job class API<br>• Managing, monitoring, and testing MapReduce jobs<br>• Work with counters<br>• Display job history and logs<br>• Write unit tests for MapReduce programs<br>• Characterizing and improving MapReduce job performance<br>• Enhance performance in your MapReduce jobs<br>• Working with different data sources in MapReduce<br>• Fault tolerance |
| 5a | Hive I | • Hive Basics |
| 5b | Hive II | • Hive Architecture<br>• Hive Query Language<br>• Loading and exporting data<br>• Use cases of Hive<br>• Steps in the data pipeline |

| | | |
|---|---|---|
| | | • Create and Load Data<br>• Create databases, internal tables, external tables, and partitioned tables<br>• Learn about data types<br>• Load data into tables and databases<br>• Query and Manipulate Data<br>• Query, sort, and filter data<br>• Hive Operators and Functions<br>• Hive Storage Formats |
| 6a | Pig | • Pig Basics<br>• Pig Latin<br>• Data Processing Operators<br>• Scripting with Pig |
| **Project or Paper Proposal Due** | | |
| 6b | Spark I (Spark Basics) | • Spark basics |
| 7 | Spark II (Spark Advanced) | • Using the Spark shell for interactive data analysis<br>• The features of Spark's Resilient Distributed Datasets<br>• Writing and Deploying Spark Applications<br>• Common Patterns in Spark Data Processing<br>• Fault tolerance |
| 8 | Spark III (Spark SQL) | • Spark SQL and the SQL Context<br>• SchemaRDD basics<br>• DataFrames basics<br>• Transforming and Querying DataFrames<br>• Loading and Saving Data<br>• Use with Hive |
| **9** | **Mid Term Exam** | • |
| 10a | Big Data Reference Architecture | • Data Warehouse and Data Mart Models<br>• Lambda Architecture<br>• Critiques and Alternatives |
| 10b | Kafka | • Basics<br>• Producers—Writing Messages<br>• Consumers—Reading Messages<br>• Architecture<br>• Fault Tolerance |
| 11 | Spark IV (Spark Streaming) | • Basics<br>• Transformations<br>• Inputs<br>• Outputs<br>• Fault Tolerance |
| 12a | NoSQL Database Landscape | • Motivations for NoSQL Databases<br>• ACID and BASE<br>    o ACID: Atomicity, Consistency, Isolation, and Durability |

| | | |
|---|---|---|
| | | <ul><li>○ BASE: Basically Available, Soft State, Eventually Consistent</li><li>○ Types of Eventual Consistency</li></ul><ul><li>Four Types of NoSQL Databases<ul><li>○ Key-Value Pair Databases</li><li>○ Document Databases</li><li>○ Column Family Databases</li><li>○ Graph Databases</li></ul></li></ul> |
| 12b | CAP Theorem | <ul><li>Basics</li><li>Critique</li><li>Other ways to characterize NoSQL databases</li></ul> |
| 13a | Key Value Databases | <ul><li>Essential Features of Key-Value Databases</li><li>Key-Value Database Terminology</li><li>Key-Value Architecture Terms</li><li>Keys: More Than Meaningless Identifiers</li><li>Values: Storing Just About Any Data You Want</li><li>Designing for Key-Value Databases</li><li>Limitations of Key-Value Databases</li><li>Design Patterns for Key-Value Databases</li></ul> |
| 13b<br><br>**Project or Paper Draft Due** | Wide Column Databases | <ul><li>Introduction to Wide Column Databases</li><li>Wide Column Database Terminology</li><li>Designing for Wide Column Databases</li><li>Usage Patterns</li><li>MapReduce Integration</li></ul> |
| 14 | Document Databases (MongoDB) | <ul><li>JSON</li><li>What Is a Document</li></ul> |
| 15 | Document Databases III (MongoDB) | <ul><li>Managing Multiple Documents in Collections</li><li>Avoid Explicit Schema Definitions</li><li>Basic Operations on Document Databases</li><li>MongoDB overview</li><li>Aggregation Framework</li><li>Designing for Document Databases</li><li>One-to-Many Relations</li><li>Many-to-Many Relations</li><li>The Need for Joins</li><li>Modeling Common Relations</li><li>One-to-Many Relations in Document Databases</li><li>Many-to-Many Relations in Document Databases</li><li>Modeling Hierarchies in Document Databases</li></ul> |
| **Final Project or Paper Due** | | |