



# Course Overview

CptS 415



# What is Big Data?

CptS 415

# What is Big Data?

- No official definition!
- A collection of datasets so large that
  - Difficult to process using on-hand database management tools
  - Difficult to process with traditional data processing applications
- The challenge include capture, curation, storage, search, sharing, transfer, analysis and visualization

# Big Data - Example

*Software application for storing and processing of images and videos:*

- store images / videos in a database*
- user interface for loading, tagging and searching the images / videos*



# Big Data - Example

*Software application for storing and processing of all WSU students' images and videos:*

- server computers*
- application and database running on different servers*



# Big Data - Example

*Software application for storing and processing of all (global) university students' images and videos:*

- multiple servers across different regions*
- distributed application and database architecture*



"Data center" by seeweb is licensed under CC BY-SA 2.0



# Big Data Applications



Search Google or type a URL



# Course Details

CptS 415



# Instructor Introduction

- B. Tech. - Indian Institute of Technology (India)
- M.S. - Vanderbilt University (Nashville, TN)
- 20+ yrs industry experience
- First Job:
  - Programmer / Data Analyst
- Recent Roles:
  - Software Engineer, SLAC
  - Founder - EpiData, Inc.
- Personal Interests:
  - Snowboarding, Golf, Hiking and Travel



# Objective and Scope

- Objective:

Understand core concepts of Big Data techniques, and develop ability to implement real-world Big Data solution

- Scope:

- Relational databases and SQL
- Document data store models
- NoSQL databases including Graph databases
- Hadoop and Spark for processing big data

# About the Course

- It is not a
  - Programming language course
  - Independent database or data mining course
- This course is
  - Provide design principles for Big Data challenges
  - Overview/survey of state-of-the-art Big Data techniques, tools and principles.
  - Provide pointers to Big Data research projects, papers, tools, and commercial/open-source projects
- This course is unique in
  - A complete overview of major Big Data techniques
  - Algorithm design techniques for Big Data
  - Academic & Industrial practice

# Course Format

- Seminar-style course:
  - No Exam!!
  - 6 homework assignments
  - 1 course projects with multiple milestones
- Suggested Textbook:
  - Database System: The Complete Book
  - Big Data Fundamentals: Concepts, Drivers and Techniques
  - Mining of Massive Datasets
- Online Tutorials & Papers
  - Research papers or chapters related to the topics
  - Checkout the resources listed in the syllabus

# Grading

Categories	Percent of Overall Grade
Participation	10%
Homework / Assignments	50%
Project	45%
<b>Total</b>	105%

# Grading (cont.)

- Participation Activities:
  - Panopto Videos
  - Survey / Poll responses
- Homework / Assignments:
  - Theoretical and design questions
  - Simple implementations
- Project:
  - Big Data project - from design to implementation
  - Presentation and report



# Grading (cont.)

Grade	Percent	Grade	Percent
A —	93 - 100	C	73 - 76.99
A-	90 - 92.99	C-	70 - 72.99
B+	87 - 89.99	D+	67 - 69.99
B	83 - 86.99	D	63 - 66.99
B-	80 - 82.99	D-	60 - 62.99
C+	77 - 79.99	F	0 - 59.99

# Academic Integrity

**Academic integrity** will be strongly enforced in this course.

**Cheating** includes, but is not limited to, plagiarism and unauthorized collaboration as defined in the Standards of Conduct for Students, WAC 504-26-010(3).

**Copyright** protects intellectual property and work of individuals, including instructors.

**Academic Integrity violations** examples (but is not limited to)

- Copying/taking a picture of another student's code/work
- Letting another student copy/take a picture of your code/work
- Sending your code/work to another student (i.e. digitally or in print)
- Receiving another's student code/work (i.e. digitally or in print)

# Weekly Plan

	Sunday	Monday - Friday	Sunday
Instructor	<ul style="list-style-type: none"><li>- Lecture Videos and Slides Released</li><li>- Assignments Released</li></ul>		
Students		<ul style="list-style-type: none"><li>- Watch Lecture Videos</li><li>- Read Textbook Chapters</li><li>- Complete Participation Activity</li></ul>	Submit Assignments / Project Milestones

# Course Schedule

Dates	Lesson Topic	Participation	Assignment
<b>Week 1</b> Aug. 21 - Aug. 27	• Course Overview, Big Data Overview	• Week 1 Videos ( <b>due Aug. 25</b> ) • Introductions ( <b>due Aug. 27</b> )	No Assignment
<b>Week 2</b> Aug. 28 - Sept. 3	• Relational DBMS	• Week 2 Videos ( <b>due Sept. 1</b> )	• Project Team Formation ( <b>Sept. 2 - Sept. 3</b> )
<b>Week 3</b> Sept. 4 - Sept. 10	• SQL, Relational Algebra	• Week 3 Videos ( <b>due Sept. 8</b> )	• Assignment 1 ( <b>due Sept. 10</b> )
<b>Week 4</b> Sept. 11 - Sept. 17	• XML and JSON	• Week 4 Videos ( <b>due Sept. 15</b> )	• Project Milestone 1 ( <b>due Sept. 17</b> ) • CATME Peer Eval. Survey ( <b>due Sept. 17</b> )
<b>Week 5</b> Sept. 18 - Sept. 24	• Graphs and RDF	• Week 5 Videos ( <b>due Sept. 22</b> )	• Assignment 2 ( <b>due Sept. 24</b> )
<b>Week 6</b> Sept. 25 - Oct. 1	• Distributed Systems and NoSQL Databases	• Week 6 Videos ( <b>due Sept. 29</b> )	• Project Milestone 2 ( <b>due Oct. 1</b> ) • CATME Peer Eval. Survey ( <b>due Oct. 1</b> )
<b>Week 7</b> Oct. 2 - Oct. 8	• Query Processing and Query Optimization	• Week 7 Videos ( <b>due Oct. 6</b> ) • Plus/Delta Survey ( <b>due Oct. 7</b> )	No Assignment
<b>Week 8</b> Oct. 9 - Oct. 15	• Graph Query Processing and Approximate Query Processing	• Week 8 Videos ( <b>due Oct. 13</b> )	• Assignment 3 ( <b>due Oct. 15</b> )
<b>Week 9</b> Oct. 16 - Oct. 22	• MapReduce	• Week 9 Videos ( <b>due Oct. 20</b> )	• Project Milestone 3 <b>due Oct. 22</b> • CATME Peer Eval. Survey ( <b>due Oct. 22</b> )
<b>Week 10</b> Oct. 23 - Oct. 29	• Hadoop	• Week 10 Videos ( <b>due Oct. 27</b> )	• Assignment 4 ( <b>due Oct. 29</b> )
<b>Week 11</b> Oct. 30 - Nov. 5	• Apache Spark	• Week 11 Videos ( <b>due Nov. 3</b> )	• Project Milestone 4 ( <b>due Nov. 5</b> ) • CATME Peer Eval. Survey ( <b>due Nov. 5</b> )
<b>Week 12</b> Nov. 6 - Nov. 12	• Big Data Theory & Practice	• Week 12 Videos ( <b>due Nov. 10</b> )	• Assignment 5 ( <b>due Nov. 12</b> )
<b>Week 13</b> Nov. 13 - Nov. 19	• Data Quality and Data Privacy	• Week 13 Videos ( <b>due Nov. 17</b> )	• Project Milestone 5 ( <b>due Nov. 19</b> ) • CATME Peer Eval. Survey ( <b>due Nov. 19</b> )
<b>Thanksgiving Break</b>	Nothing due this week		
<b>Week 14</b> Nov. 17 - Dec. 3	• NewSQL and In-Memory DBMS	• Week 14 Videos ( <b>due Dec. 1</b> )	• Assignment 6 ( <b>due Dec. 3</b> )
<b>Week 15</b> Dec. 4 - Dec. 10	Project Completion Week	• Course Evaluation (Course Feedback)	• Project Presentations and Demos ( <b>Dec. 8 - Dec. 9</b> )
<b>Week 16</b> Dec. 11 - Dec. 15	• Finals Week		• Project Report ( <b>due Dec. 11</b> ) • CATME Peer Eval. Survey ( <b>due Dec. 11</b> )



# Summary