# City University of Hong Kong Course Syllabus

# offered by Department of Computer Science with effect from Semester A 2015/16

## Part I Course Overview

Course Title:	Big Data Algorithms and Techniques
Course Code:	CS5488
Course Duration:	1 Semester
Credit Units:	3 Credits
Level:	P5
Medium of Instruction:	English
Medium of Assessment:	English
<b>Prerequisites</b> : (Course Code and Title)	CS3402 Database Systems
<b>Precursors</b> : (Course Code and Title)	Nil
<b>Equivalent Courses</b> : (Course Code and Title)	Nil
<b>Exclusive Courses</b> : (Course Code and Title)	Nil

## 1. Abstract

This course is aimed at equipping students with the ability to manage very large data sets (Big Data) using a cluster of commodity machines with the main focus on the Hadoop ecosystem. It has three specific objectives: (1) to familiarize students with software systems and techniques for implementing distributed data-parallel programs, (2) to provide insight into internal mechanisms of large-scale data analytical systems, and (3) to acquaint students with big data solutions deployed in real-world settings. Students will also have the opportunity to analyse and to compare real-world big data solutions in a class project case study.

### 2. Course Intended Learning Outcomes (CILOs)

(CILOs state what the student is expected to be able to do at the end of the course according to a given standard of performance.)

No.	CILOs	Weighting	Discov	very-enr	riched
		(if	curricu	ılum rel	lated
		applicable)	learnin	g outco	omes
			(please	e tick	where
			approp	riate)	
			A1	A2	A3
1.	Identify data parallelism to be exploited in large-scale data		<i>✓</i>	~	~
	processing problems.				
2.	Implement data parallel algorithms using techniques		$\checkmark$	$\checkmark$	
	action the course				
	covered in the course.				
3.	Understand the internal mechanisms of the Hadoop		~		
	framework				
4	Tune work.				
4.	Design scalable solutions to a real-world problem and		~	~	~
	sufficiently provide rationalizations to the design decisions.				
5	Analyse existing hig data solutions deployed in real world				
5.	Analyse existing big data solutions deployed in real-world		•	~	v
	settings through case studies.	1000/			
		100%			

### A1: Attitude

Develop an attitude of discovery/innovation/creativity, as demonstrated by students possessing a strong sense of curiosity, asking questions actively, challenging assumptions or engaging in inquiry together with teachers.

A2: Ability

Develop the ability/skill needed to discover/innovate/create, as demonstrated by students possessing critical thinking skills to assess ideas, acquiring research skills, synthesizing knowledge across disciplines or applying academic knowledge to self-life problems.

A3: Accomplishments

Demonstrate accomplishment of discovery/innovation/creativity through producing /constructing creative works/new artefacts, effective solutions to real-life problems or new processes.

#### 3. Teaching and Learning Activities (TLAs)

(TLAs designed to facilitate students' achievement of the CILOs.)

### Teaching pattern:

Suggested lecture/laboratory mix: 2 hrs. lecture; 1 hr. tutorial.

TLA	Brief Description			No			Hours/week
		1	2	3	4	5	(if applicable)
1. Looturo	Lectures will cover (1) different types of data-parallel	<	~	~			2 hours
Lecture	problems; (2) the APIs and tools for handling big data;						
	(3) the internal mechanisms for job scheduling, failure						
	handling, and task execution of the Hadoop framework;						
	(4) case studies on real-world big data algorithms and						
	solutions.						
2. Lab	Laboratory classes will provide the students with the		~	<	~		1 hour
	opportunity to (1) familiarize themselves with different						
	big data tools; (2) implement data-parallel algorithms; (3)						
	design experimental studies.						
3. Class	For the class project, the students will work on a use case				1	~	
Project	study on Hadoop-based solutions deployed in real-world						
5	settings.						

# 4. Assessment Tasks/Activities (ATs)

(ATs are designed to assess how well the students achieve the CILOs.)

Assessment Tasks/Activities	CILO No.					Weighting	Remarks
	1	2	3	4	5		
Continuous Assessment: 60%							
Class Project	$\checkmark$	$\checkmark$		~	~	40%	
In-class/online Discussions	~	~	~	~		5%	Participating in class and online discussions
Midterm Examination	$\checkmark$	$\checkmark$	$\checkmark$			15%	
Examination <sup>*</sup> : <u>40</u> % (duration: 2 h	ours)						
Final Examination	$\checkmark$	$\checkmark$	$\checkmark$	~		40%	
						100%	

<sup>^</sup> For a student to pass the course, at least 30% of the maximum mark for the examination must be obtained.

### 5. Assessment Rubrics

(Grading of student achievements is based on student performance in assessment tasks/activities with the following rubrics.)

Assessment Task	Criterion	Excellent $(A + A - A)$	Good (B+ B B-)	Adequate $(C + C - C)$	Marginal (D)	Failure (F)
1. Group Project	1.1 Ability to identify challenges in various types of data parallel problems	High	Significant	Moderate	Basic	Inadequate
	1.2 Ability to critique existing big data solutions	High	Significant	Moderate	Basic	Inadequate
2.In-class/Online Discussions	2.1 Ability to contribute to discussions on principle and concepts of scalable data processing	High	Significant	Moderate	Basic	Inadequate
<ol> <li>Midterm Exam</li> <li>Final Exam</li> </ol>	3.1, 4.1 Ability to demonstrate a good understanding of materials covered in the course	High	Significant	Moderate	Basic	Inadequate

Part III Other Information (more details can be provided separately in the teaching plan)

#### 1. Keyword Syllabus

(An indication of the key topics of the course.)

Big Data, Analytics, Hadoop, MapReduce, Distributed File System, Parallel Processing, Data-parallel Systems, RDBMS, NOSQL, Distributed Indexes, Key-value Stores, Query Languages, Data Manipulation Languages, Consistency, Reliability, Commodity Cluster, Failure Handling, Search Algorithms, Social Network Applications, In-memory Processing, Use Case Studies, Nutch, Log Processing, Cascading, TeraByte Sort, Graph Data

### 2. Reading List

#### 2.1 Compulsory Readings

(Compulsory readings can include books, book chapters, or journal/magazine articles. There are also collections of e-books, e-journals available from the CityU Library.)

1. Tom White, Hadoop: The Definitive Guide, Third Edition, O'Reilly Media

#### 2.2 Additional Readings

(Additional references for students to learn to expand their knowledge about the subject.)

1.	Holden Karau, Andy Konwinski, Patrick Wendell, Matei Zaharia, Learning Spark: Lightning-
	Fast Big Data Analysis, First Edition, O'Reilly Media