# Course Introduction

## CMPT 732, Fall 2018

### Greg Baker

https://coursys.sfu.ca/2018fa-cmpt-732-g1/pages/

## Us

Instructor: Greg Baker, University Lecturer in CS. Fourth time teaching this course.

TAs:

- Sethuraman Annamalai
- Vishal Shukla
- Raman Singh

All of the TAs are previous-cohort Big Data students.

## Welcome



Greg's orientation summary:

- SFU/CS/Big Data provide more support than you may be used to. Use it.
- The co-op staff (Paula, Eunice, Laura) will guide you into the job search process. Listen to them.
- Katie is a good first-contact for all other administrative/academic things.
- The TAs (and I) have experience with the program and can provide guidance on academics/technology.

## This Course

It's "Programming for Big Data I". So it's about:

1. programming,
2. big data.

## What is Big Data?

A quote you have probably seen before:

> "Big data is like teenage sex: everyone talks about it, nobody really knows how to do it, everyone thinks everyone else is doing it, so everyone claims they are doing it..." [Dan Ariely]

It's a buzzword. But a useful one.

## How big is "Big Data"?

Answer 1: Big enough that traditional techniques can't handle it.

The "traditional techniques" strawman is usually (1) highly structured data (2) in a relational database (3) on a single database server.

PostgreSQL can handle single tables up to 32 TB [*] (but maybe not as quickly as you'd like). Big data must be bigger than that?

Answer 2: Big enough that one computer can't store/process it.

"One computer" can mean dozens of cores, TBs of RAM and many TB of storage. So bigger than that?

## "Big data" isn't always big.

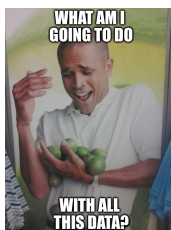Many describe with "The Four V's" (or 5 Vs' or 7 V's...).

- Volume: the amount of data.
- Velocity: the data arrives quickly and constantly.
- Variety: many differently-structured (or less-structured) input data sets.
- Veracity: some data might be incorrect, or of unknown correctness.

Honestly, the term *big data* is often used to mean "modern data processing, 2013–2018 edition".

Even if most people don't work with truly-big data most of the time, it's nice to have the tools to do it when necessary.

Sometimes it's nice to know your computation can scale if a megabyte of data becomes many gigabytes.

Or maybe "can't be processed on one computer" should be "can't be processed in a time I'm willing to wait on one computer".

An "overnight" calculation that doesn't complete until noon isn't very useful.

Greg's functional definition: people say they're doing "Big Data" when they have so much data, it's annoying.



## Clusters

If our data is going to be too big for one computer, we presumably need many. Each one can store some of the data and do some of the processing and they can work together to generate "final" results.

This is a *cluster*.

Computer cluster

Actually managing work on a cluster sucks. You have all of the problems from an OS course (concurrency, interprocess communication, scheduling, …) except *magnified* by being a distributed system (some computers fail, network latency, …).

The MPI tools are often used to help with this, but are still very manual.

Do you want to worry about all that? Me neither. Obvious solution: let somebody else do it.

## Hadoop

We will be using (mostly) the Apache Hadoop ecosystem for storing and processing data on a cluster. This includes:

- YARN: managing jobs in the cluster.
- HDFS: Hadoop Distributed File System, for storing data on the cluster's nodes.
- MapReduce: system to do parallel computation on YARN.
- Spark: another system to do computation on YARN (or elsewhere).
  ⋮

## Our Environment

We have a cluster for use in this course: 6 nodes, each 16 cores and 110 GB memory. We will explore in the assignments.

Not a big cluster, but big enough we can exercise the Hadoop tools, and do calculations that aren't realistic for one computer.

In many ways, a two or three node cluster is enough to learn with: you have many of the same issues as 1000 nodes.

## Things you will do

- Lectures (2 hour/week) and labs (4 hours/week).
- Assignments: 10 of them. Complete in labs + later.
- Project: a more open-ended big data analysis.

## Lecture and Labs

Lecture: ASB 9204W on Tuesday 13:30–15:20.

Labs in ASB 10920:

- Tues/Thurs 9:00–10:50.
- Wed/Fri 11:30–13:20.
- Wed/Fri 13:00–15:50.

Greg and one TA will be in each lab.

There are computers in the lab rooms, or whatever laptop you have should be workable. Go to the lab section you're registered for.

No Tuesday/Wednesday labs this week.

## Course Topics

My current plan for the assignments is a good outline:

1. YARN, HDFS, MapReduce
2. MapReduce vs Spark
3. Spark RDDs
4. RDDs vs DataFrames
5. Spark DataFrames
6. DataFrames and ML
7. NoSQL and Cassandra
8. Spark + Cassandra
9. Spark Streaming
10. Hadoop reliability

The "express computation on a cluster" tools we'll see: MapReduce, Spark RDDs, Spark DataFrames.

There is a progression from "lower-level and more explicit" to "higher level and more declarative".

Programming languages:

MapReduce will use **Java**: we will use it relatively little. Suggestion: don't bother setting up an IDE and go command-line-only for the few times you need it.

We'll use Spark with **Python**. Most of your programming in this course will be in Python.

# Expectations

- The assignments and project are your chance to learn things, not a thing you have to do for marks.
- They are individual work: copying and understanding your friend's solutions isn't "your" work. [Academic Honesty]
- You are expected to be in the labs during your lab times, and working on the assignments.
- Everybody be nice to each other. [Code Of Conduct]