

# Advanced Topics in Data Science CSCI 6720G

With Focus on Big Data Integration, Data  
Curation, Data Quality and Data Cleaning

Jarek Szlichta

Data Science Lab

Course Website:

<http://data.science.uoit.ca/teaching/data-science/>

# Agenda for Today

- Big Data
- Course Structure
  - Presentation
  - One pager
  - Class project
  - Mark breakdown
- Course Outline

# The Era of Big Data

- Unprecedented growth in data being generated and its potential uses/value [TPCTC-MC]
  - Tweets, social networks (statuses, check-ins, shared content), blogs, click streams, various logs, ...
  - *Facebook: > 1B active users, > 8B messages/day*
  - *Twitter: > 140M active users, > 340M tweets/day*
- Everyone is interested
  - Trade press and popular press: *“Big Data!”*
  - Enterprises, Web companies, online businesses, governments, public health researchers, social scientists, ...
  - Untapped value and countless new opportunities to understand, optimize, and/or compete

[TPCTC-MC] - Mike Carrey's keynote at TPCTC

# Social Networks - Facebook

- Daily user activity
  - 2.5B content items shared
  - 2.7B 'Likes'
  - 350M photos uploaded
- Data volume statistics
  - 100+ PB in a single HDFS cluster
  - 500+ TB of new data generated per day
  - 70K queries per day

# Data Integration



# Facebook Country

## The Republic of Facebook

If Facebook were a country....



It would be home to 1 in 7 of the world's entire population

**Sources**

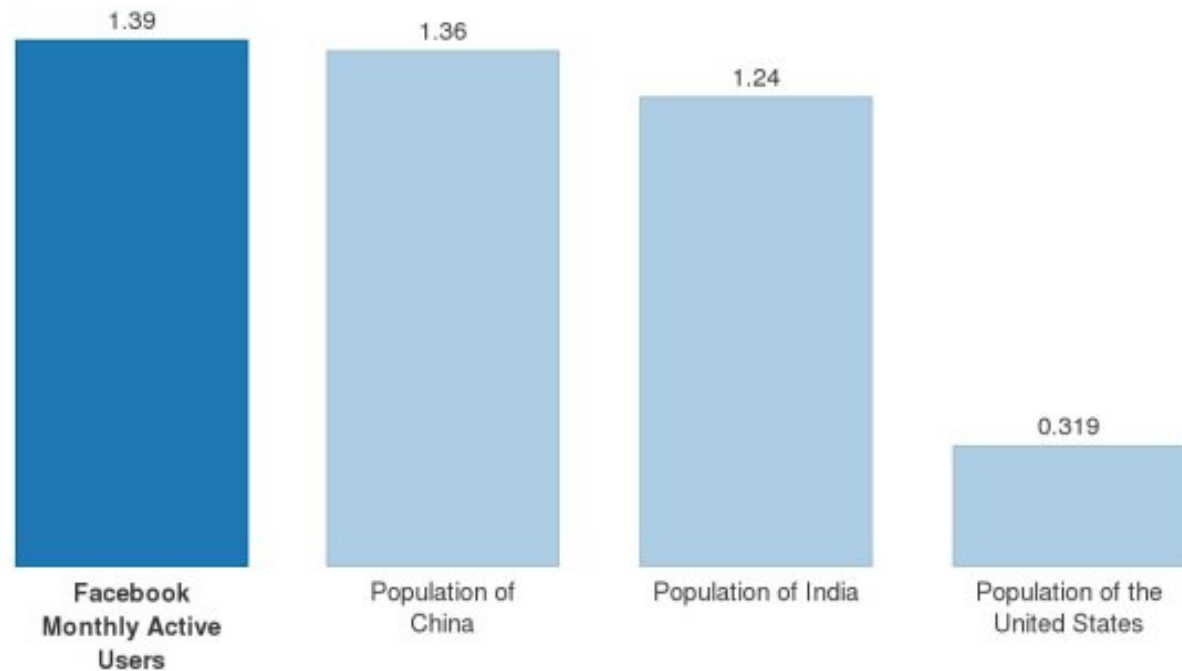
[www.newsroom.facebook.com/Key-Facts](http://www.newsroom.facebook.com/Key-Facts)  
[www.en.wikipedia.org/wiki/World-Population](http://www.en.wikipedia.org/wiki/World-Population)

[www.blogsession.co.uk](http://www.blogsession.co.uk)

# Facebook Country

## How Big Is Facebook?

Facebook has more active users than China has people (figures in billions)



Source: Facebook, CIA World Factbook

The Huffington Post



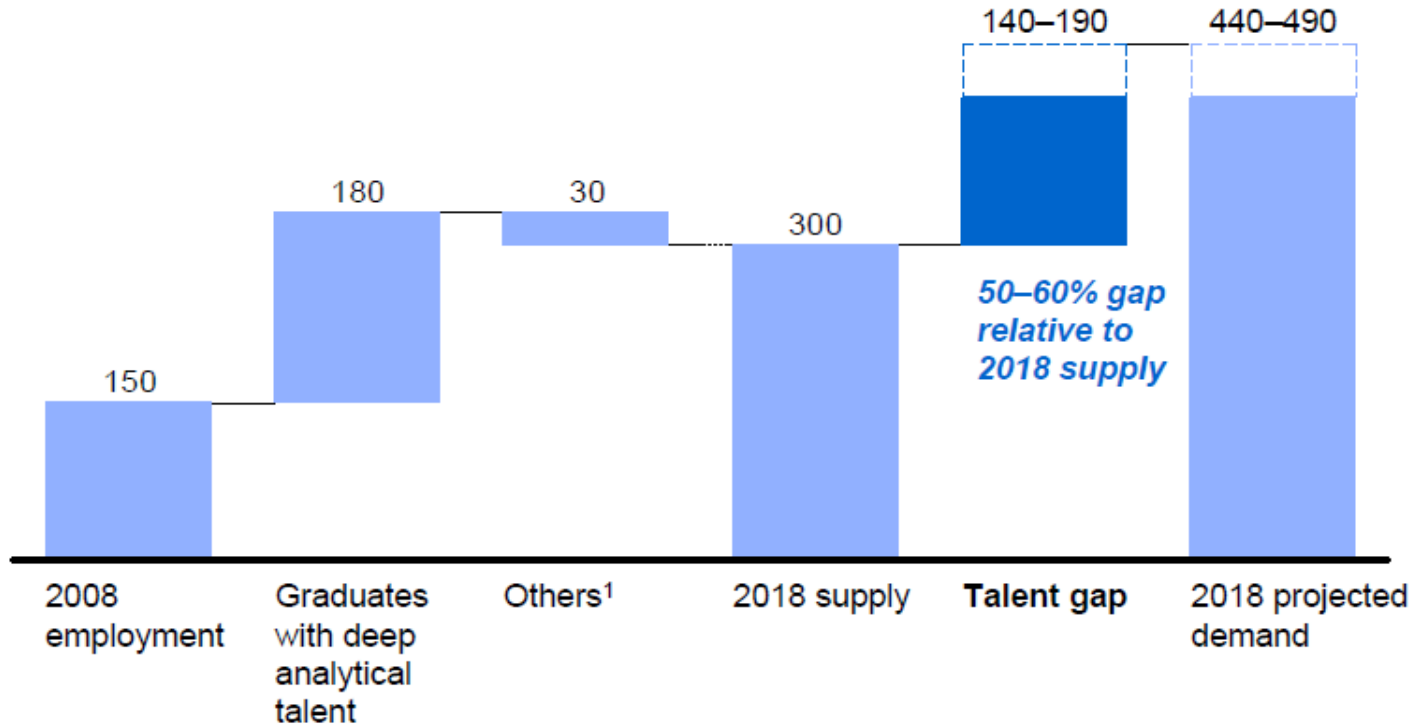
Data contains value and knowledge



# Good news: Demand for Big Data

**Demand for deep analytical talent in the United States could be 50 to 60 percent greater than its projected supply by 2018**

Supply and demand of deep analytical talent by 2018  
Thousand people



<sup>1</sup> Other supply drivers include attrition (-), immigration (+), and reemploying previously unemployed deep analytical talent (+).  
SOURCE: US Bureau of Labor Statistics; US Census; Dun & Bradstreet; company interviews; McKinsey Global Institute analysis

# Course Structure

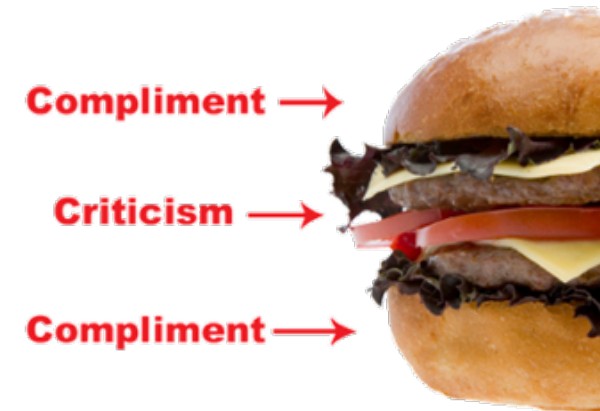
- Presentation & Discussion
  - Each week 3 presentations
  - Recent research papers on
    - Big data integration (data integration, data analytics, data curation, data cleaning, machine learning)
- Class project
  - 1 person per team..
  - Related to the course topic

# Presentation

- Everyone should read the paper before the lecture!
- Presentation time: 30 minutes
- Prepare slides
- In general your slides should:
  - be easy to read,
  - avoid too much text!
  - use a lot **figures, examples** to introduce concepts for ease understanding!
  - not be too long. Usually you will not need more than one slide per minute. As a rule of thumb you will need 1.5-2 minutes per slide
  - have a clear and easy to follow structure
  - practice your talk before presenting at class!

# Discussion

- After each presentation
  - Led by the instructor
- Prepare your questions
- Give feedback to the presentation
  - Feedback Burger
    - Basis: Positive details
    - Meat: What could be improved
    - Top: Conclusion / Motivation



# One Pager

- **One** page summary of one paper per 3 weeks
  - Typed, 10pt font
- It should:
  - summarize the problem(s) addressed/solved by the research paper (1-2 sentences that clearly describe the problem: "The problem is ... .")
  - briefly sketch the main ideas on which the solution of the problem is based
    - briefly describe the research methodology of the paper (1-2 sentences)
  - identify 3 strong points and 3 weak points of the paper
  - summarize any assumptions the solution in the paper is based upon (any restrictions; stated assumptions and non-stated assumptions)
  - raise three non-trivial questions about the paper (including future work)
- Is marked

# Class Project

- Independent exploration of specific problem
  - within the course topics
  - modeling, design, algorithms, simulation, and analysis
  - Propose solution
- Evaluation of project:
  - timeliness, development and presentation of idea (i.e., proposal, progress, and final project report, and final in class presentation).
- No more than 12 pages (LNCS template)
- Project agenda
  - Proposal – submit online
  - Progress report – submit online
  - Final presentation – in-class
  - Final report – conference style paper

# Mark Breakdown

- 15% One pagers
- 30% Presentations and discussion leading
- 20% Participation and interactions (discussion, readings, ideas etc.)
- 35% Course project (proposal, progress reports, presentation, final report)

# Short Term Agenda

- Next week presentation – my research and discussions
- The week after we will start with your presentations
  - **Everyone** has to read papers before the class!
- Course Website: [data.science.uoit.ca](http://data.science.uoit.ca)
- Questions
  - [jarek@uoit.ca](mailto:jarek@uoit.ca)



# Paper Assignments

- Top tier big data conferences: VLDB, SIGMOD, ICDE, EDBT, CIKM

# Sample Paper Assignments

- Fei Chiang, Renée J. Miller: A unified model for data and constraint repair. ICDE 2011: 446-457
- Lukasz Golab, Howard J. Karloff, Flip Korn, Avishek Saha, Divesh Srivastava: Sequential Dependencies. PVLDB 2(1): 574-585 (2009)