

# CS 644: Introduction to Big Data

Daqing Yun  
New Jersey Institute of Technology



# **Who, where, and when**

- Daqing Yun, [dy83@njit.edu](mailto:dy83@njit.edu)
- CKB 313
- R, 2:30 – 5:20 PM EST
- See Moodle page for more detailed syllabus

# The 1st Class Attendance Check

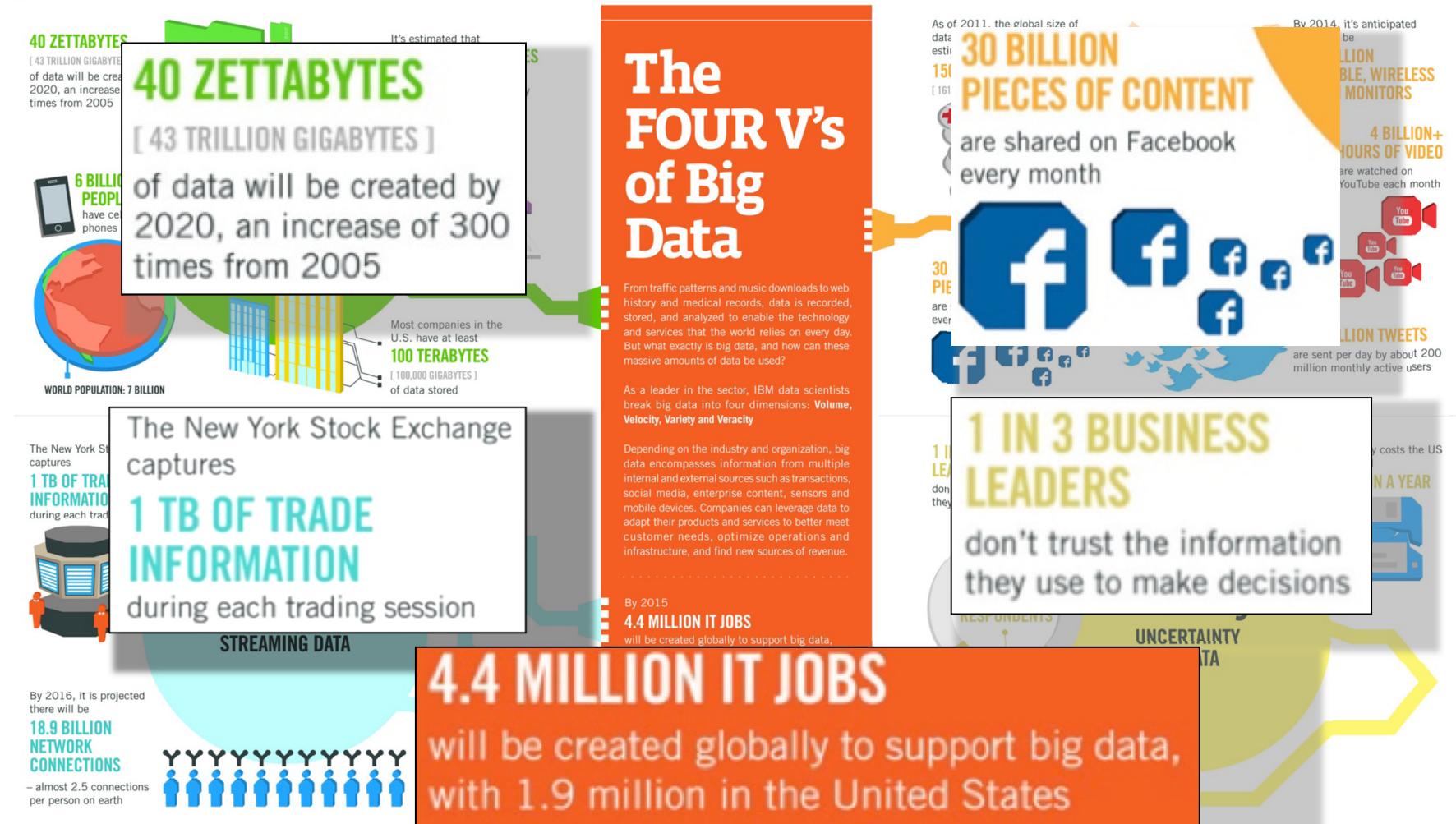
- Name
- Program (MS, Ph.D., etc.)
- Year
- Why do you take this course?
- What is the largest data size you've ever personally handled and in what context (application domain, data type, storage format, etc.)?

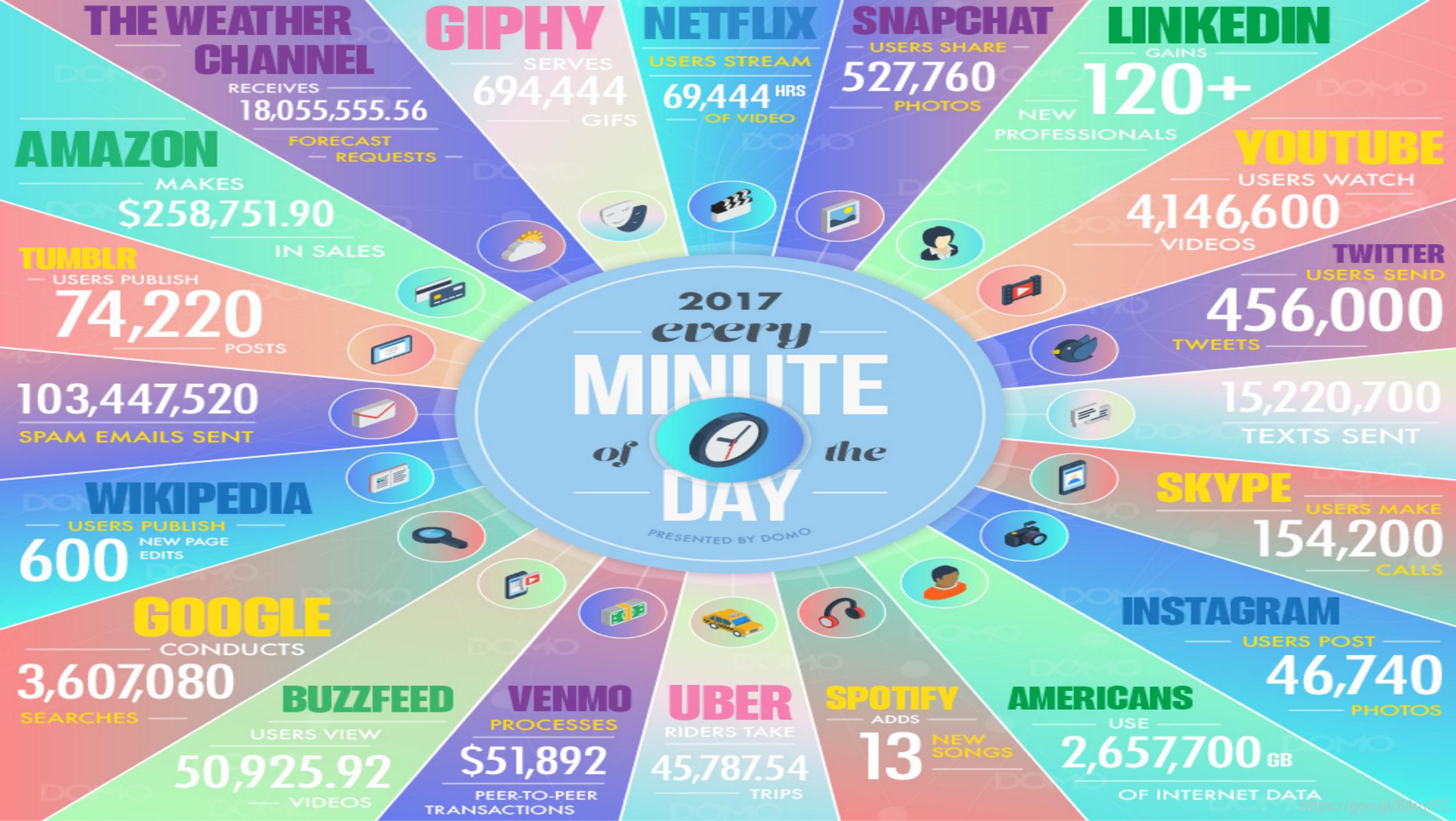
# About this course

- Recent Developments and Future Trends on Big Data Computing
- Overview of Big Data Analytics
- Platforms, Storages, Databases, and Algorithms for Big Data Analytics
- Advanced Topics:
  - Big-Data Visualization
  - Big-Data Movement
  - Big-Data Security
  - Big-Data Workflows

# Characterization of big data

- Four V's

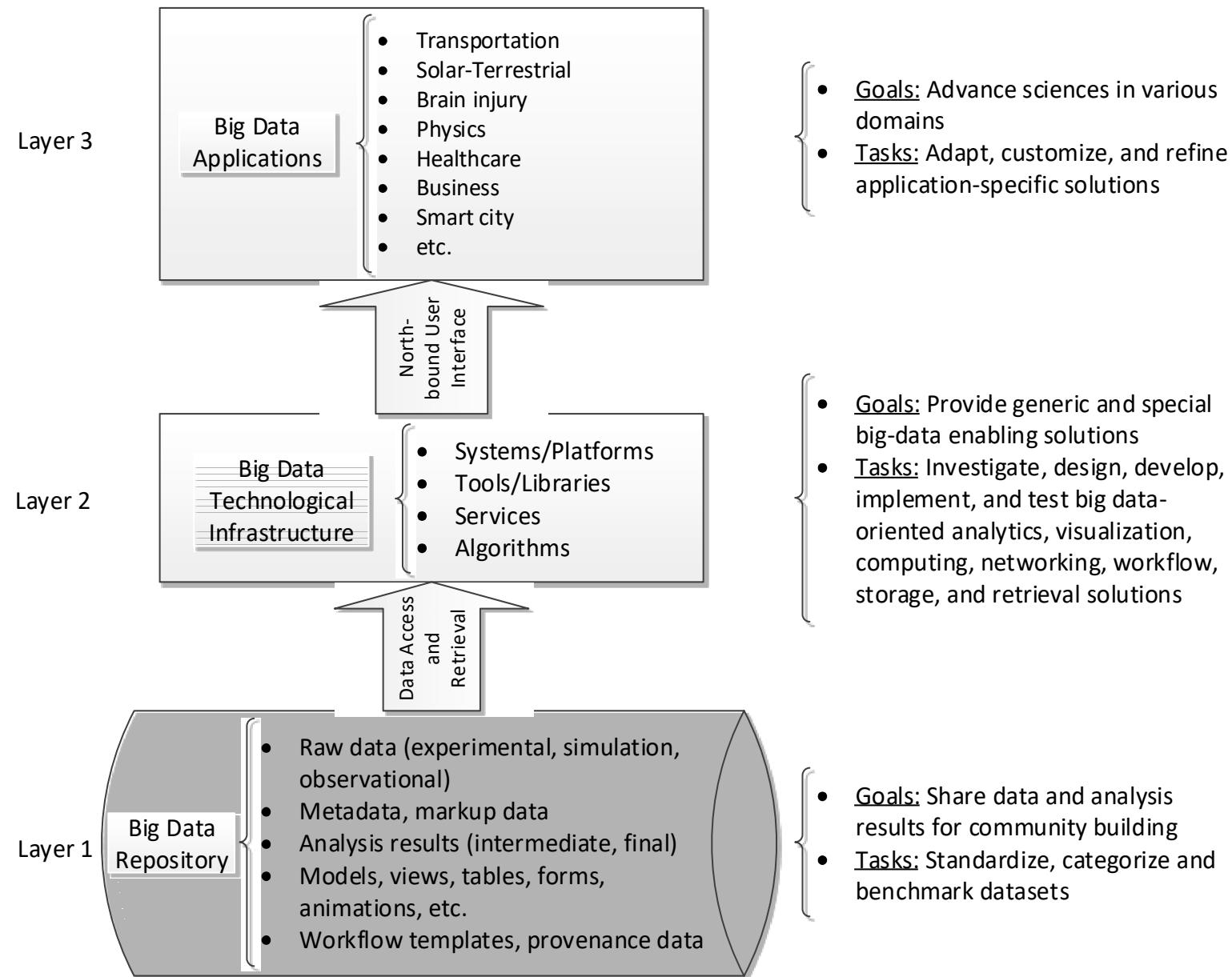




# Center for Big Data

- Director: Chase Wu, [chase.wu@njit.edu](mailto:chase.wu@njit.edu)
- <https://centers.njit.edu/bigdata>
- Location: GITC 4111
- Mission Statement:
  - Synergize the strong expertise in various disciplines across the NJIT campus
  - Build a unified platform that embodies a rich set of big data enabling technologies and services with optimized performance to facilitate research collaboration and scientific discovery
  - Investigate, develop, and apply cutting-edge technologies to address unprecedented challenges in big data with high Volume, high Velocity, high Variety, and high Veracity, in order to create high Value

# A Three-layer Structure of the CBD



# **A Three-layer Structure of the CBD**

- Layer 1: Big Data Repository**

- Store, manage, and provide a wide variety of data such as raw data (experimental, simulation, observational, and user-generated content), metadata, markup data, analysis results (intermediate and final) in various forms including models, views, tables, images, and videos, and workflow templates with provenance data.
  - Build a dedicated one-stop portal to share research data and analysis results for community building.

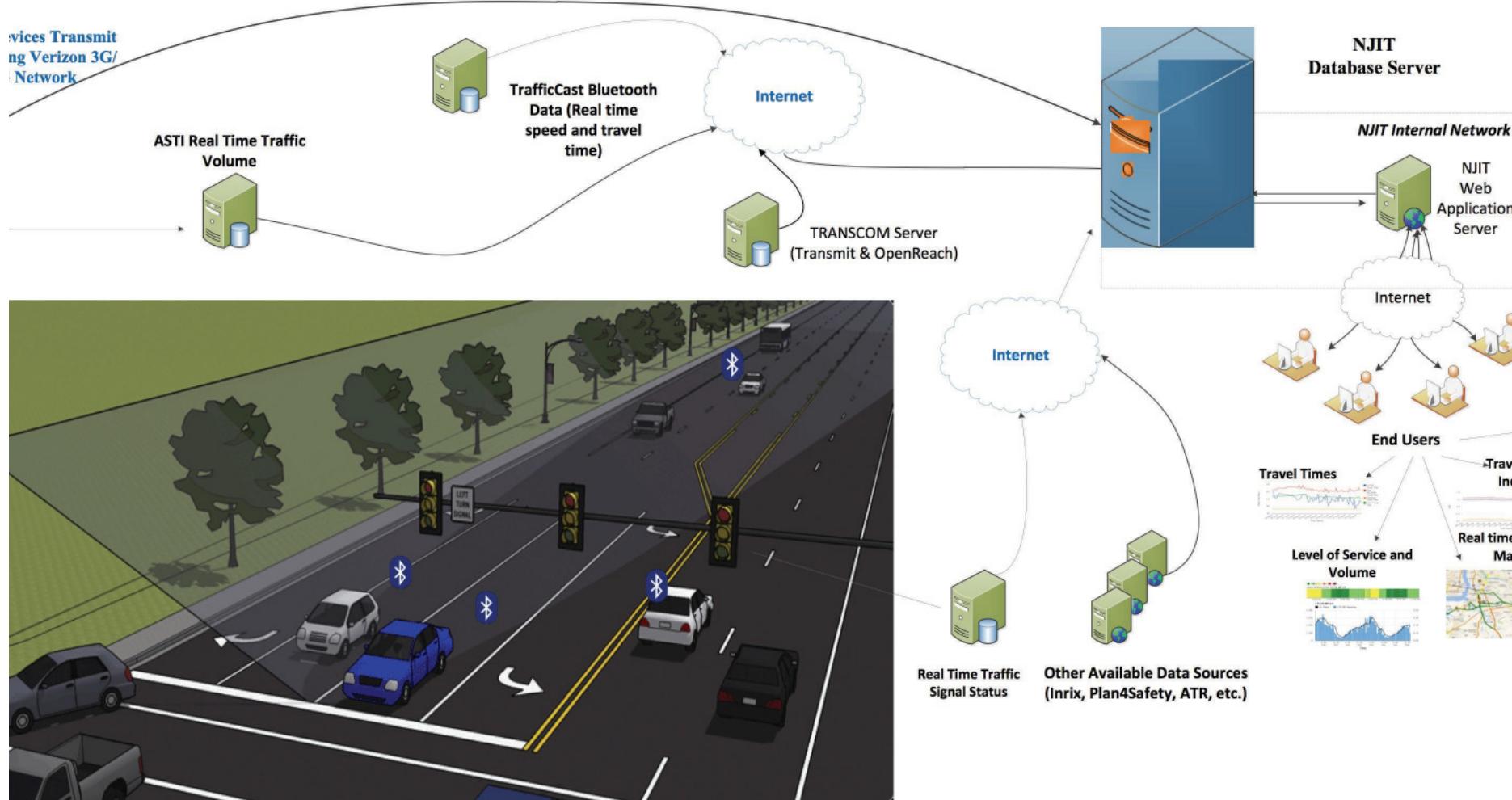
- Layer 2: Big Data Technological Infrastructure**

- Provide generic and domain-specific big data enabling solutions for data management, movement, and analytics.
  - Host and maintain a set of practical technical resources in the form of systems/platforms, tools/libraries, services, and algorithms in various areas including database management, data mining, machine learning, and parallel and distributed computing, which are needed to compose big data solutions in different application domains.

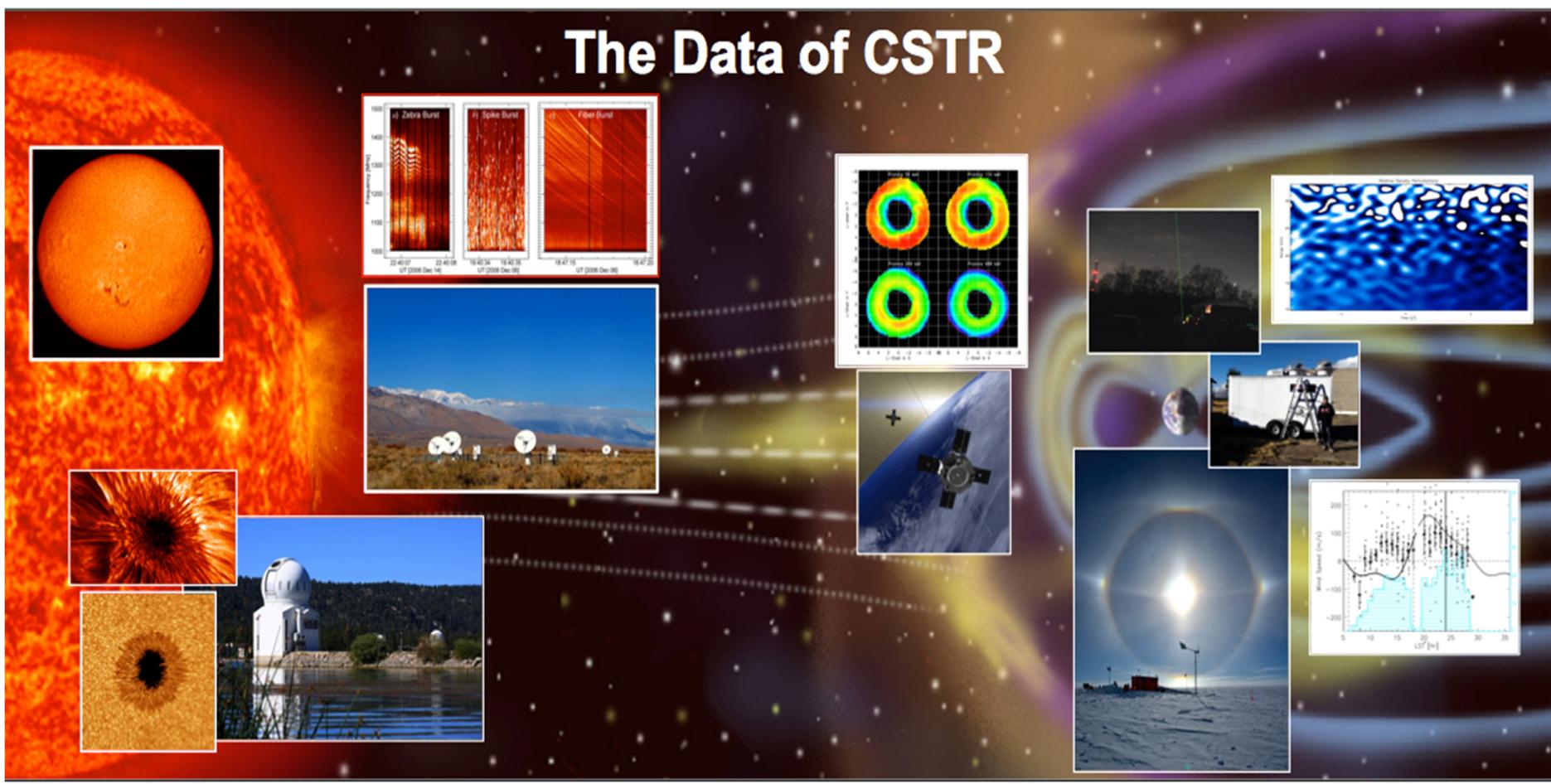
- Layer 3: Big Data Applications**

- Present a common portal to big data applications spanning across a wide spectrum of research fields, including transportation, solar-terrestrial, brain injury, physics, healthcare, business, smart city
  - Provide researchers powerful and customized big data solutions to advance the frontier of sciences in various application domains.

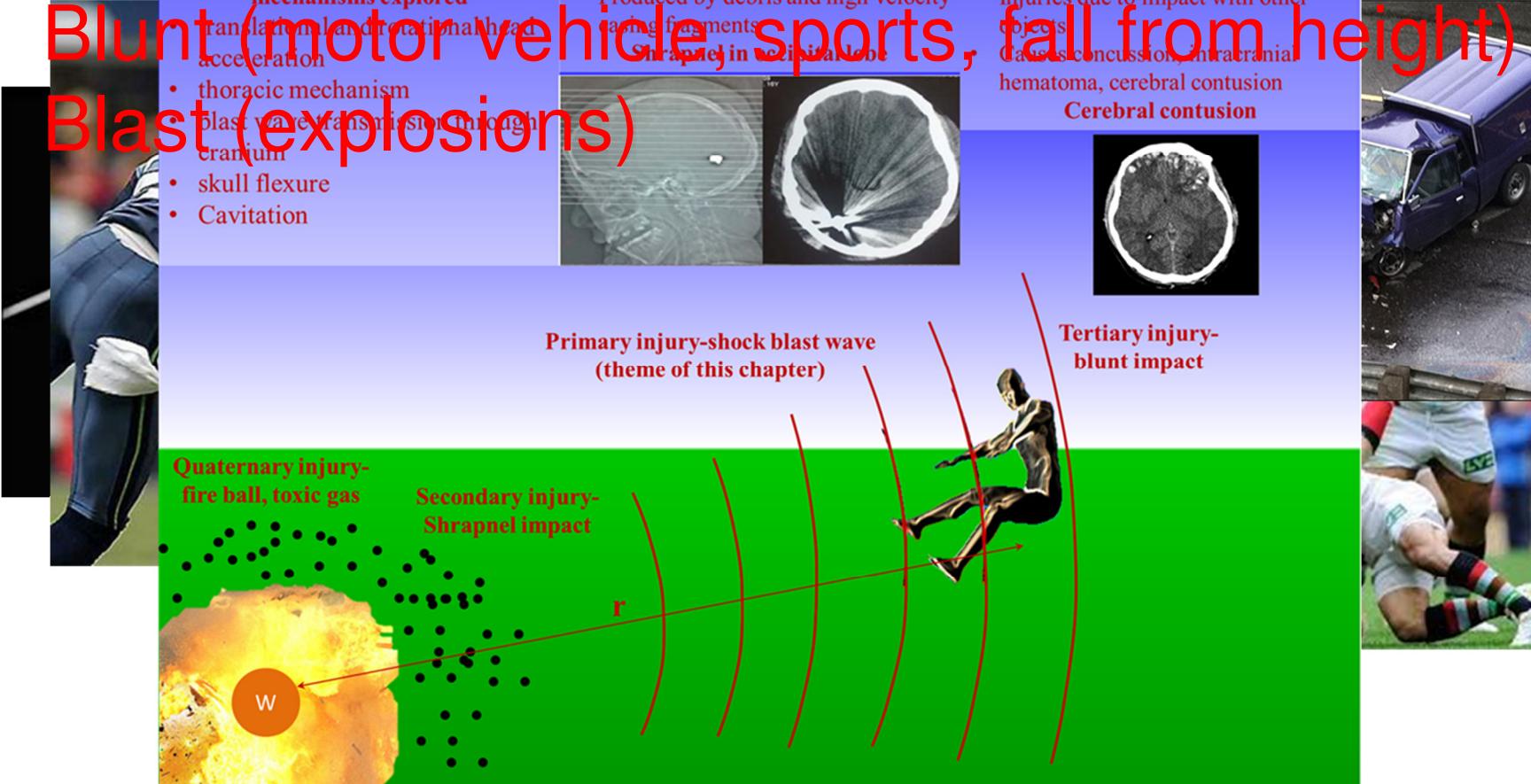
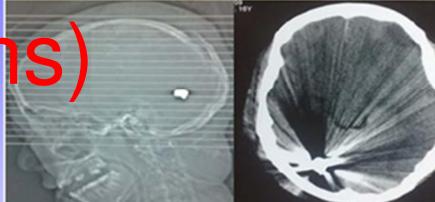
# Transportation



# Solar Terrestrial Research

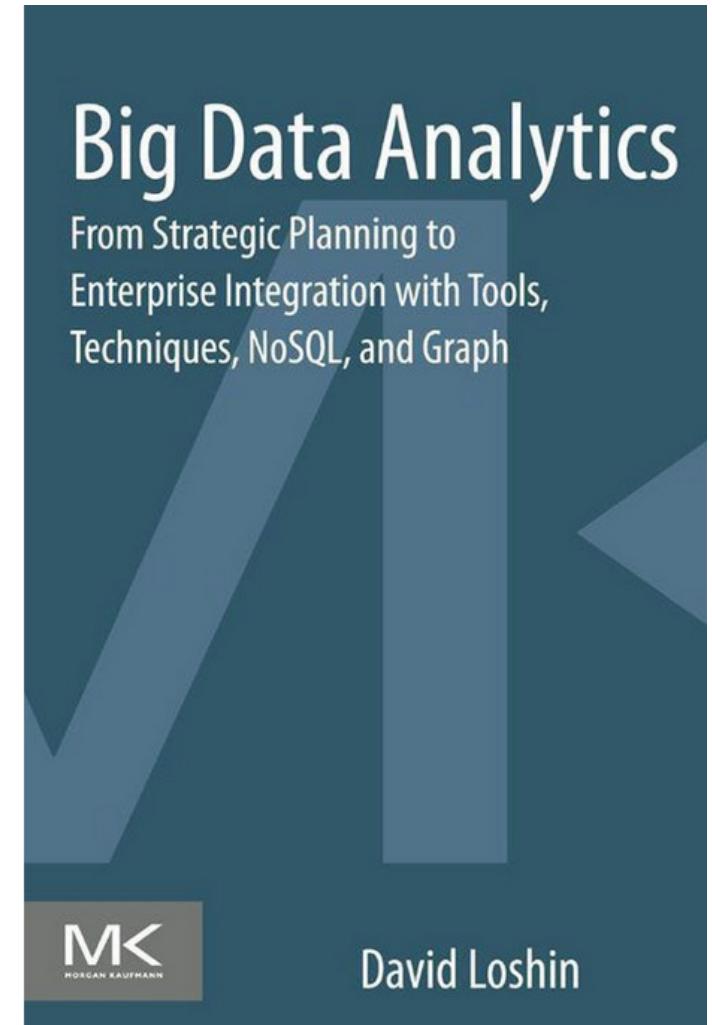


# Classification of Traumatic Brain Injury

Ballistic (bullet)	Blunt Injury-most prevalent	Blast (military)
<ul style="list-style-type: none"><li>• Ballistics (Bullet, shrapnel)</li><li>• Blunt (motor vehicle, sports, fall from height)</li><li>• Blast (explosions)</li></ul>  <p>Primary shock blast wave Current mechanisms explored Secondary injury-Shrapnel impact Produced by debris and high velocity shrapnel fragments Shrapnel in occipital lobe</p> <p>Primary injury-shock blast wave (theme of this chapter) Secondary injury-Shrapnel impact Tertiary injury-blunt impact Injuries due to impact with other objects Causes concussion, intracranial hematoma, cerebral contusion Cerebral contusion</p> <p>Quaternary injury- fire ball, toxic gas</p>		 

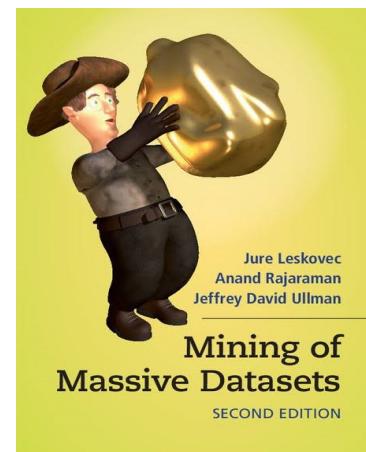
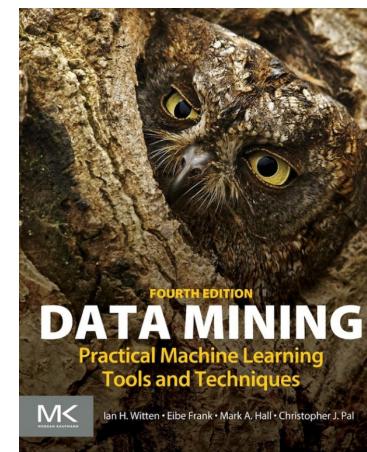
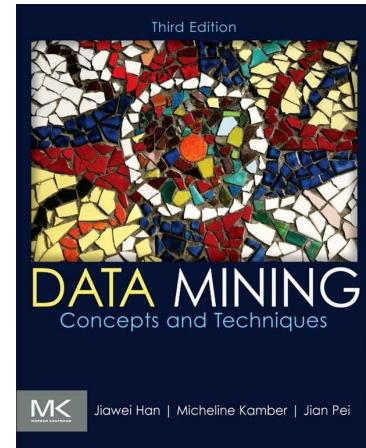
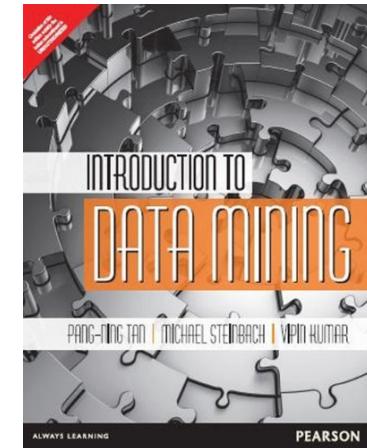
# Reading materials

- Chapter 1: Market and Business Drivers for Big Data Analysis
- Chapter 2: Business Problems Suited to Big Data Analytics
- Chapter 3: Achieving Organizational Alignment for Big Data Analytics
- Chapter 4: Developing a Strategy for Integrating Big Data Analytics into the Enterprise
- Chapter 5: Data Governance for Big Data Analytics: Considerations for Data Policies and Processes
- Chapter 6: Introduction to High-Performance Appliances for Big Data Management
- Chapter 7: Big Data Tools and Techniques
- Chapter 8: Developing Big Data Applications
- Chapter 9: NoSQL Data Management for Big Data
- Chapter 10: Using Graph Analytics for Big Data
- Chapter 11: Developing the Big Data Roadmap

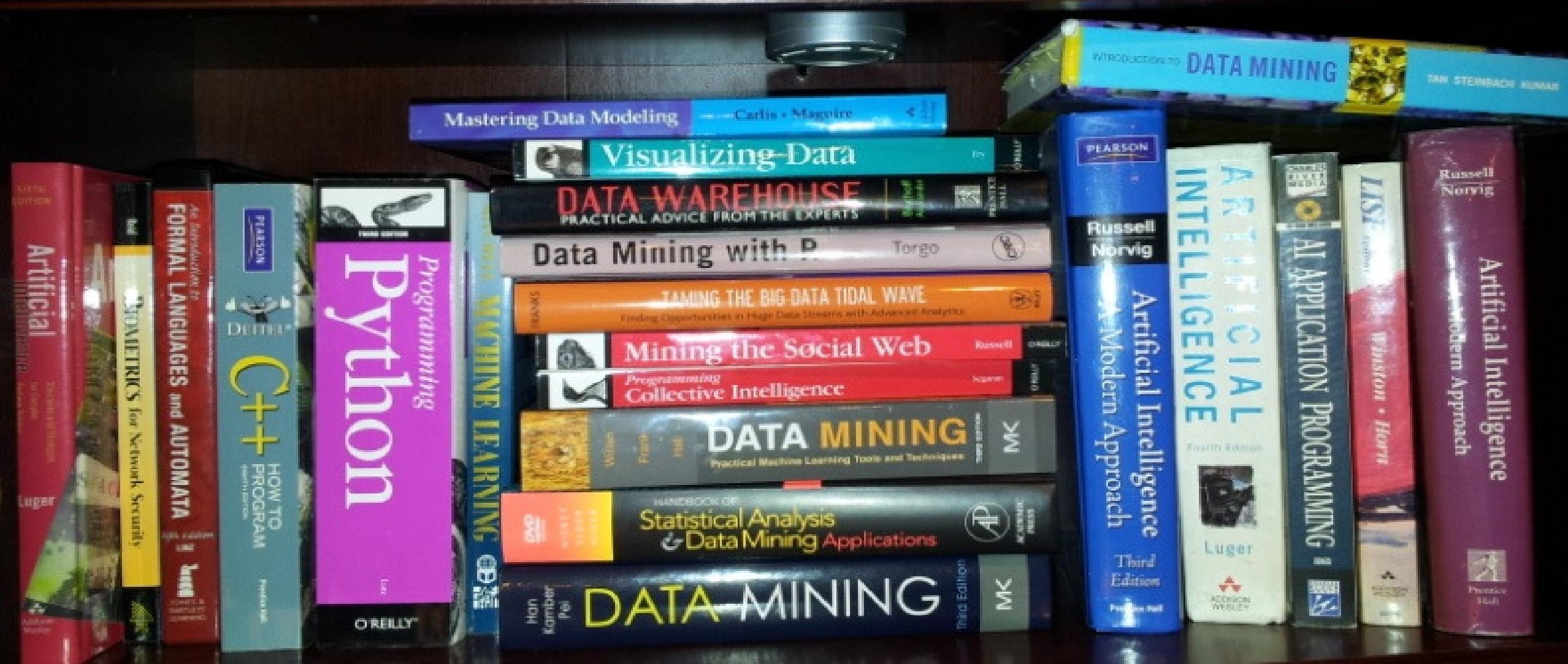


# Reading materials

- Introduction to Data Mining, by P. Tan, M. Steinbach, and V. Kumar, ISBN: 978-93-325-1865-0
- Data Mining: Concept and Techniques (3rd Ed.), by J. Han, M. Kamber, and J. Pei, ISBN-13: 978-9380931913
- Data Mining: Practical Machine Learning Tools and Techniques (4th Ed.), by I. Witten, E. Frank, M. Hall, and C. Pal, ISBN: 978-0-12-804291-5
- Mining of Massive Datasets, by J. Leskovec, A. Rajaraman, and J. Ullman



# If you are really interested...



# Exascale Computing and Big Data

- By Daniel A. Reed and Jack Dongarra, July 2015,  
Communications of the ACM
- <https://vimeo.com/129742718>



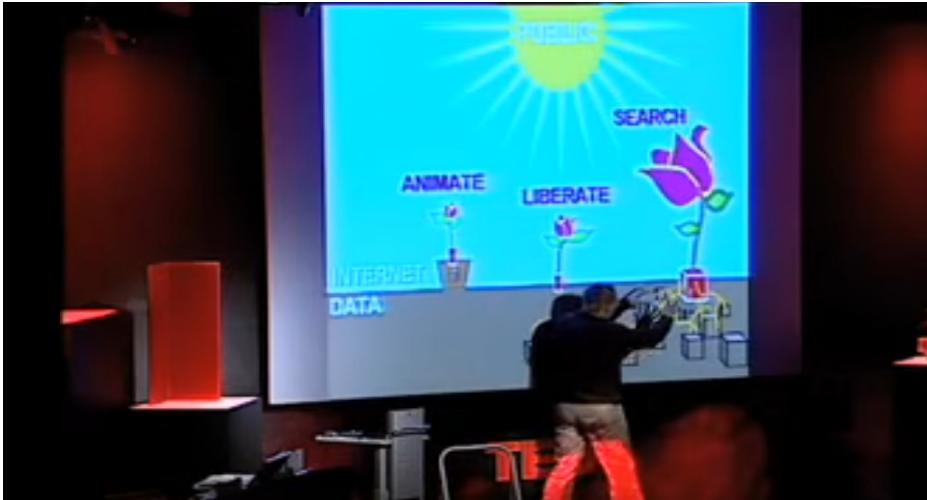
# Future of Data Mining



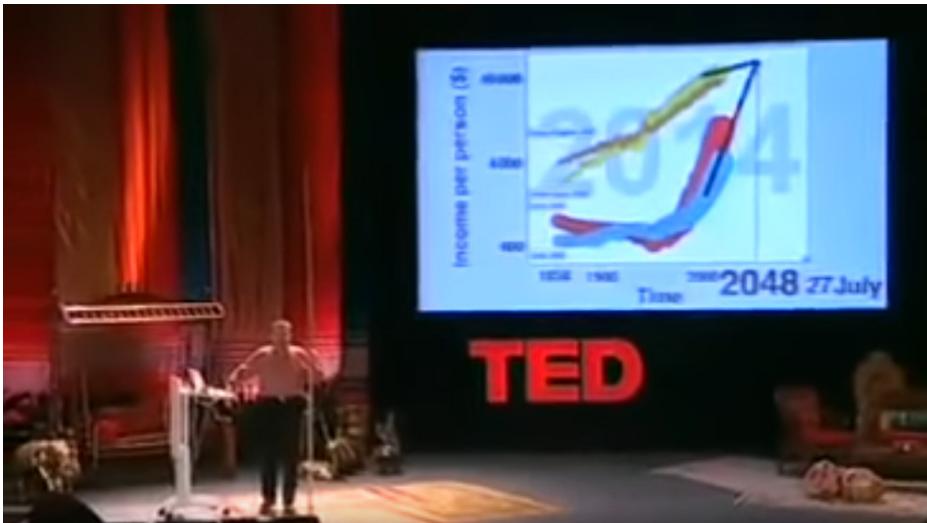
<https://www.youtube.com/watch?v=DS310JMdu2s>

By Bill Weir, C. Michael Kim, David Miller, Justin Bare & Mark Monroy | [This Could Be Big](#) – Tue, Dec 4, 2012

# The Magic of Data

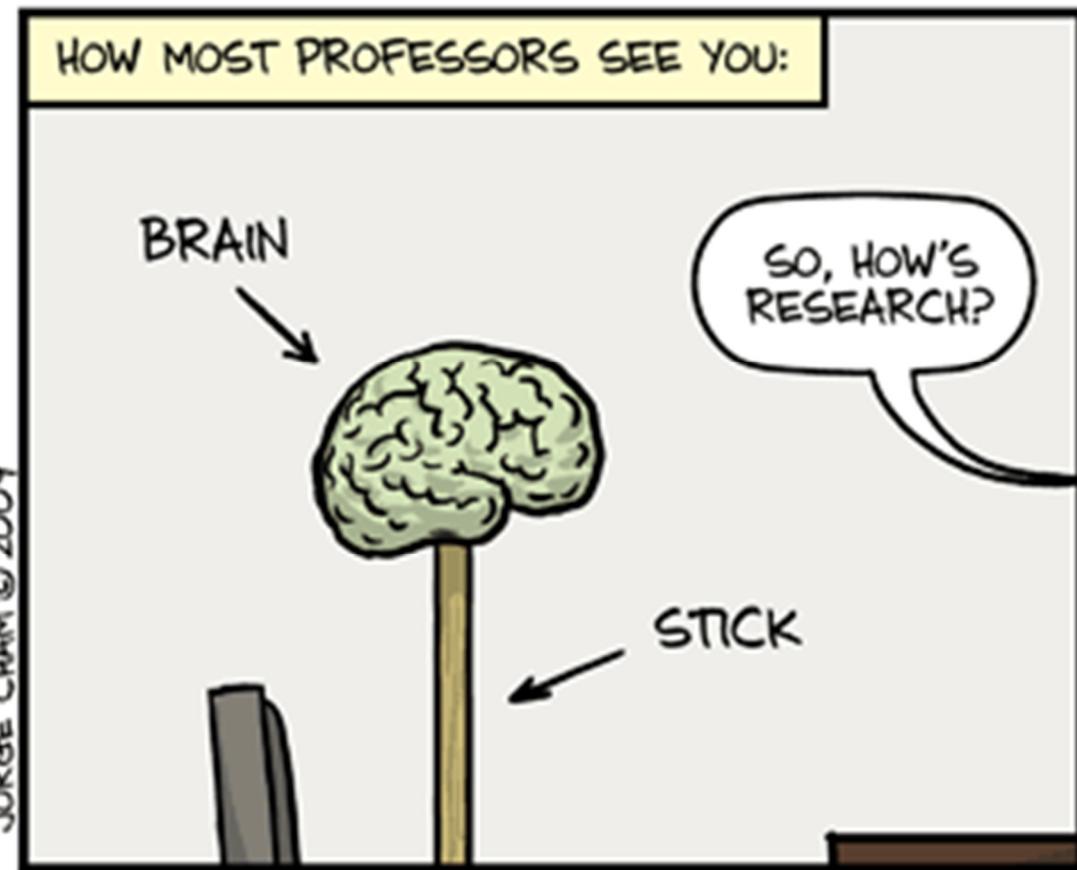
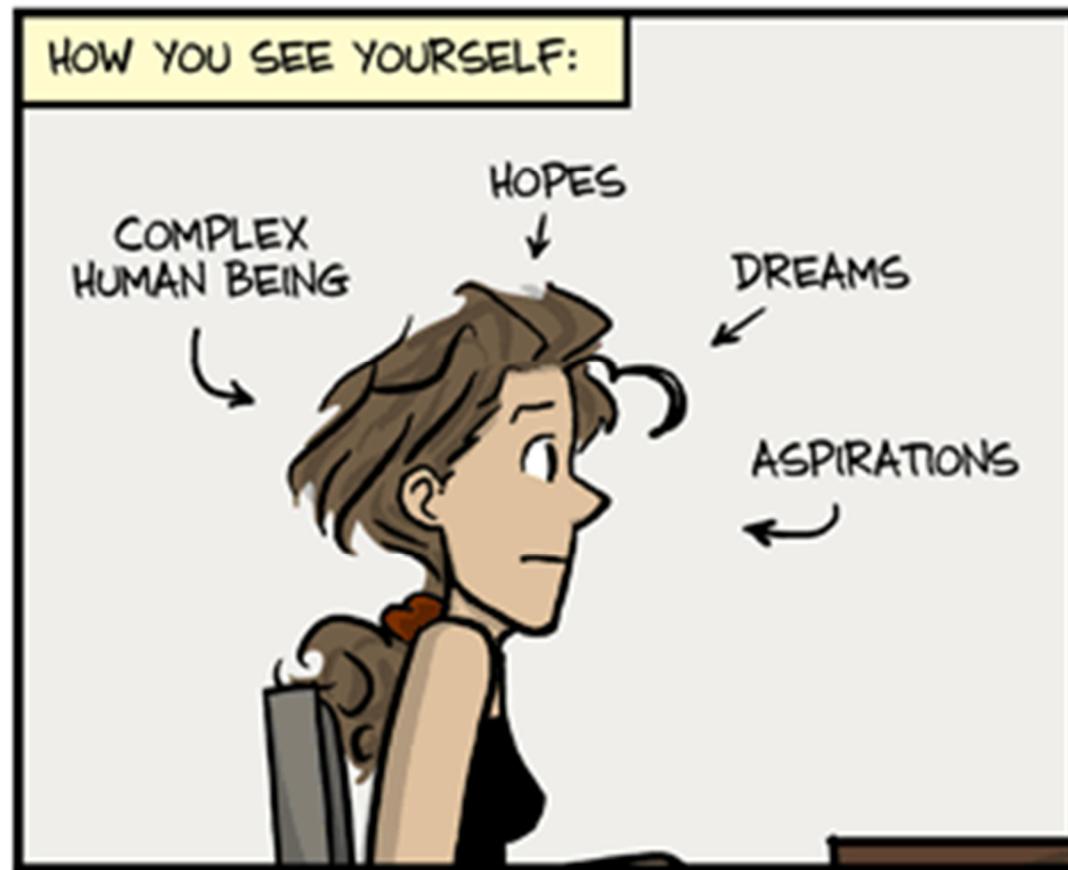


<https://www.youtube.com/watch?v=RUwS1uAdUcl>  
Hans Rosling: Debunking third-world myths with the best stats you've ever seen



<https://www.youtube.com/watch?v=fIK5-oAaeUs>  
Hans Rosling: Asia's rise - how and when

# I know you have other things to do, but...



# The first rule of CS 644 is...

- *Don't Plagiarize in CS 644*
- The second rule of CS 644 is...
  - *Don't Plagiarize in CS 644*
- There are plagiarism-detection systems
  - Work pretty well, seriously
  - Academic integrity tutorial: <https://goo.gl/ZX9oBe>
- The *penalty* would be: 0.0 in this course

# Reminder

- HW 1 is assigned, due by next lecture
- See Moodle for details

**Hope you will enjoy...**





*Thanks !* ☺

*Questions ?*