# CS5344:Big Data Analytics
# Lesson 1: Introduction

https://canvas.nus.edu.sg/courses/38824

**Anthony Tung**
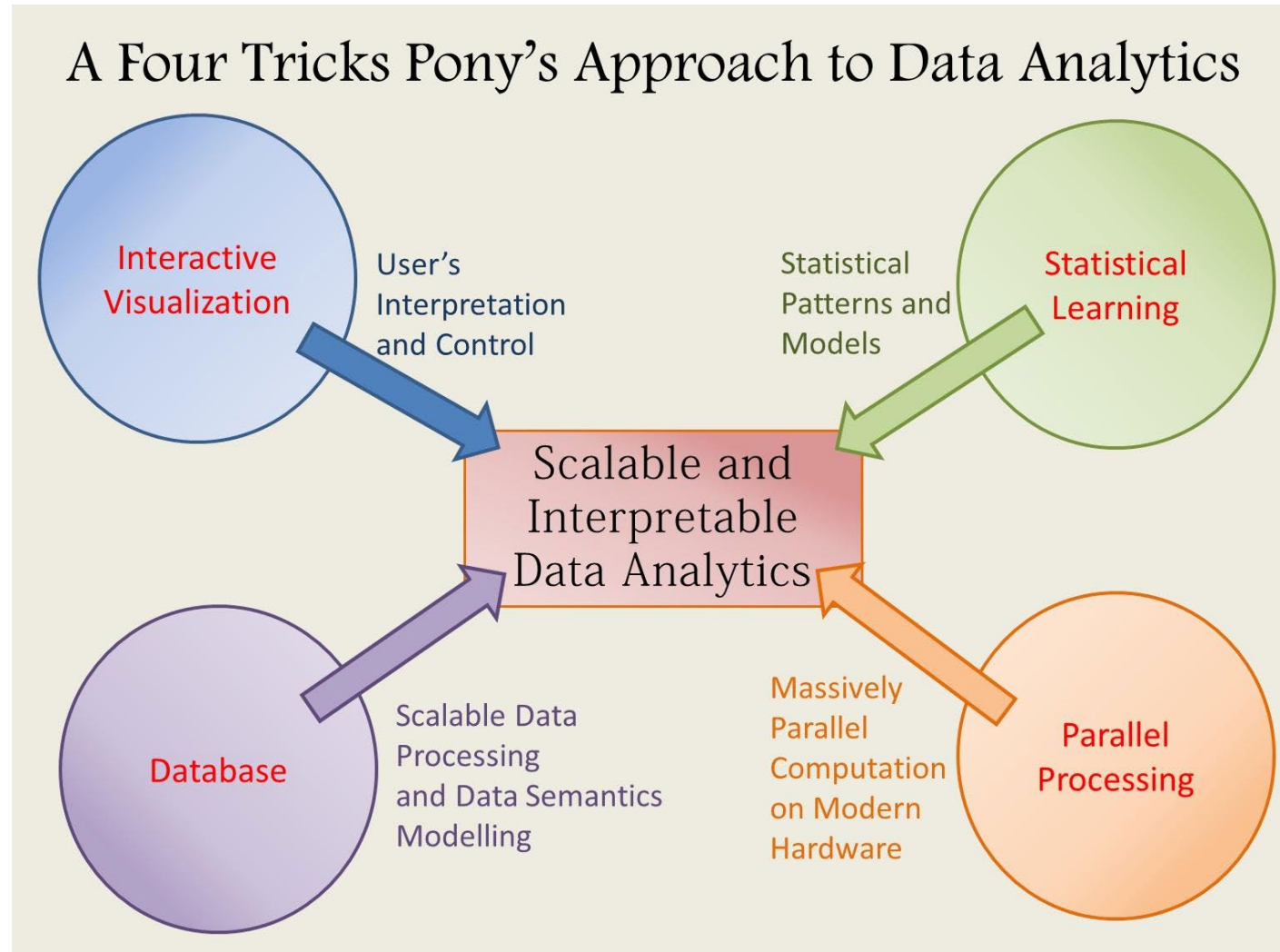**Department of Computer Science**
**anthony@comp.nus.edu.sg**

# Right Infringements on NUS Course Materials

**All course participants (including permitted guest students) who have access to the course materials on LumiNUS or any approved platforms by NUS for delivery of NUS modules are not allowed to re-distribute the contents in any forms to third parties without the explicit consent from the module instructors or authorized NUS officials**
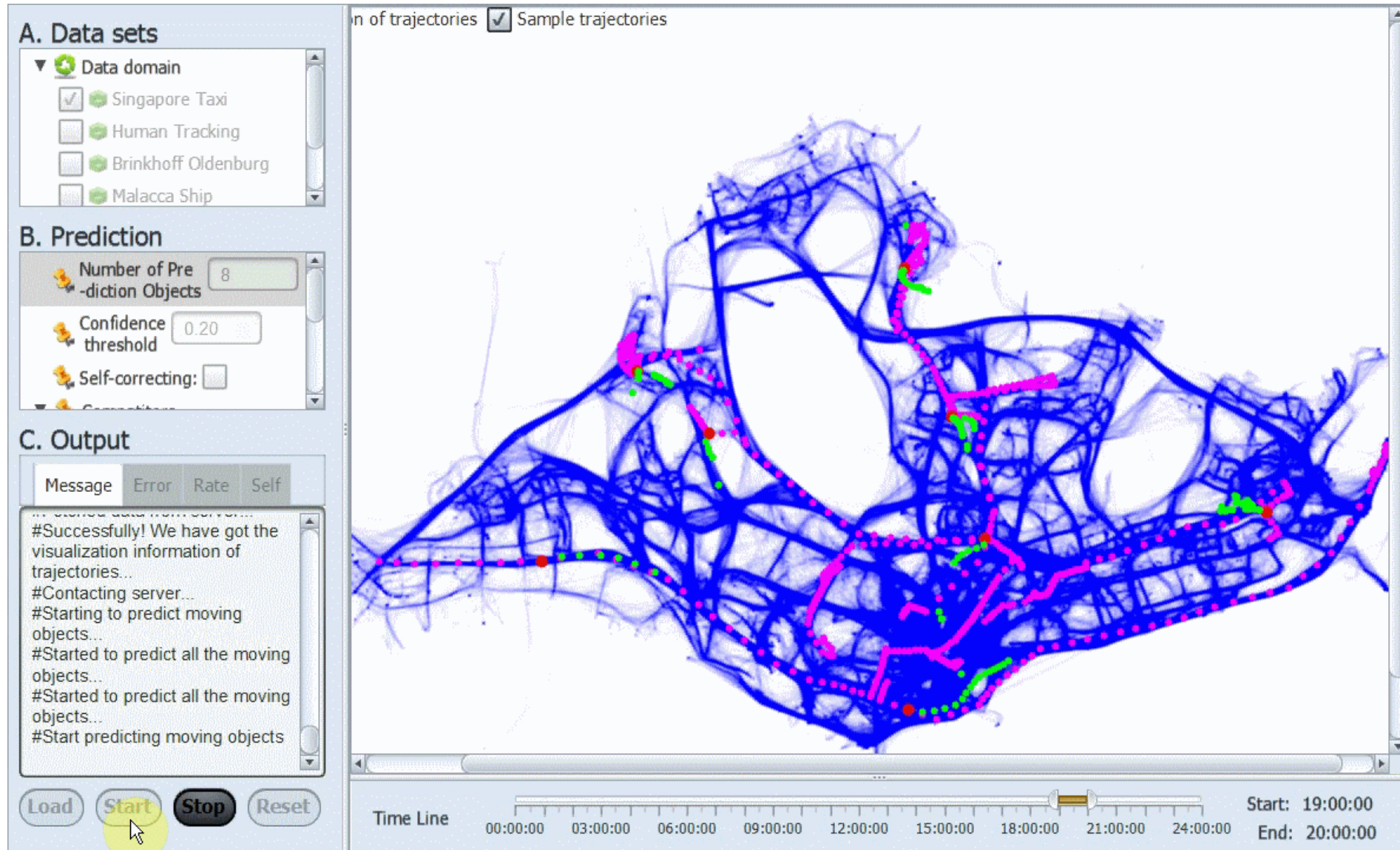
# Class Information

- **Lecturer: Anthony Tung**
    - Email: anthony@comp.nus.edu.sg
- **Tutors:**
    - Wu Shengqiong swu@u.nus.edu
    - Xu Danni dannixu@u.nus.edu
- **Lectures on Monday 1830 – 2030 Physical F2F only**
- **Office hours/Project Consultation: Monday 2035 - 2200hrs**
- **Course website: https://canvas.nus.edu.sg/courses/61715**
- **Reference text(Do NOT need to buy)**
    - Mining of Massive Datasets by J. Leskovec, A. Rajaraman and J.D. Ullman (available online: http:///www.mmds.org)
    - Introduction to Data Mining (Second Edition) by Anuj Karpatne , Michael Steinbach, Pang-Ning Tan, Vipin Kumar

# About Myself



A Four Tricks Pony's Approach to Data Analytics

# What is Big Data?

- **Gartner's Definition**

*"Big data" is high-volume, -velocity and -variety information assets that demand cost-effective, innovative forms of information processing for enhanced insight and decision making.*

- **Information assets characterized by 3Vs**

  - **High-volume (Terabytes → Zettabytes)**

  - **High-velocity (Batch → Streaming data)**

  - **High-variety (Structured → Semistructured & unstructured)**

> **Data becomes BIG when the volume, velocity or variety EXCEEDS the abilities of our IT systems to ingest, store, analyze and process it to derive actionable intelligence in a TIMELY manner.**

# Volume: How Much Data?



- Amount of data we create every day, every minute
- 90% of the data in the world today has been created in one year alone
- Data comes from everywhere e.g. sensors gather climate data, posts to social media, digital pictures and videos, purchase transaction records, cell phone GPS signals etc.
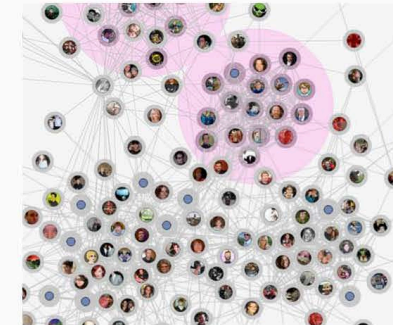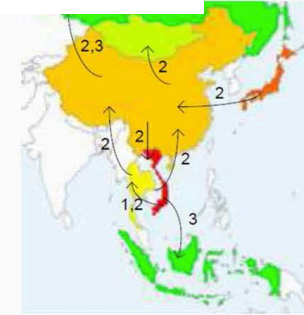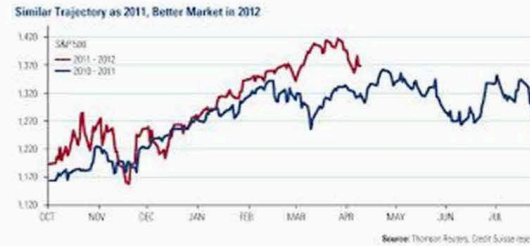
# Volume: How Much Data?

- **Facebook**
  - >250 billion photos (>600 petabyte)
  - 6 billion messages per day (5-10 terabyte)
  - >1500 million users (2 trillion connections?)
- **Sloan Digital Sky Survey**
  - 35% of the sky mapped
  - >1 billion objects classified
  - 100 terabyte of data available

# Velocity: At What Speed?



- Pace at which data flows in from sources
- Bursts of activities
- Real-time analysis
  - Late decisions → Missed opportunities

# Variety: What Kind of Data?

- Relational databases
- Transactional databases
- XML databases
- Spatial databases
- Temporal databases
- Text databases and multimedia databases
- Graph databases

Relationships between people

23

**Do not fit into a data warehouse, into neat tables of columns and rows.**
**Better place in Hadoop Distributed File System (HDFS) or in non-relational NoSQL databases.**

# Fourth V - Veracity

- **How accurate or trustworthy is the data?**

- **Bias, inconsistencies, half truth**

- **Reliability of data source**



DATA QUALITY    ACCURACY          INTEGRITY    VALIDITY

---

**Dennis** @brit_newsman · 7 Mar 2014
BREAKING: Malaysian flight **MH370** aircraft found at **Nanning**, China. Emergency landing. Waiting comfirmation from airline

**Zaim Aizzat** @zaimaizzat · 7 Mar 2014
RT @saupee Aircraft found at **Nanning**, China. Emergency landing. Waiting comfirmation frm MAS. #prayforMH370 #**MH370**

**Nota Kembara** @NotaKembara · 7 Mar 2014
Alhamdulillah. **MH370** Aircraft Emergency landing at **Nanning**, China.

▬▬▬ · now
MAS CEO confirms SAR ops and says airline is working to verify speculation that MH370 may have landed in Nanning.



Mystery of **MH370**

# Why Big Data?

**$5 million vs $500**
Price of fastest supercomputer in 1975 and iPhone with comparable performance

- **Can collect cheaply, due to automation**

- **Can store cheaply, due to falling media prices**

- **Can create Value**  $600 to buy a disk drive that can store all of the world's music

  - Turn 12 terabytes of tweets created each day into improved product sentiment analysis

  - Convert 350 billion meter readings to better predict power consumption

  - Find communication patterns of successful projects in emails

  - Analyze elevator logs to predict vacated real estate

  - Scrutinize 5 million trade events created each day to identify potential fraud (time-sensitive, sometimes 2 minutes is too late)

  - Monitor 100's of live video feeds from surveillance cameras to target points of interest (new insights when you link and analyse different data types together)

# Why Big Data?

*Data contains Value and Knowledge*

# Big Data Analytics

- **From raw data to actionable information**

  - Complex process of examining large and varied datasets to uncover information (hidden patterns, unknown correlations, market trends, customer preferences) that can help organizations make informed business decisions

- **Data needs to be**
  - **Stored**
  - **Managed**
  - **and ANALYZED**

  *Discover - Do we really know what we have?*

  *Explore - How do different data relate to each other?*

  *Iterative - What are the actual relationships?*

# Big Data Analytics

**≈ Data Mining ≈ Data Science**

- **Discover patterns and models that are**
  - **Valid:** hold on new data with some certainty
  - **Useful:** should be possible to act on the item
  - **Unexpected:** non-obvious to the system
  - **Understandable:** humans should be able to interpret the pattern
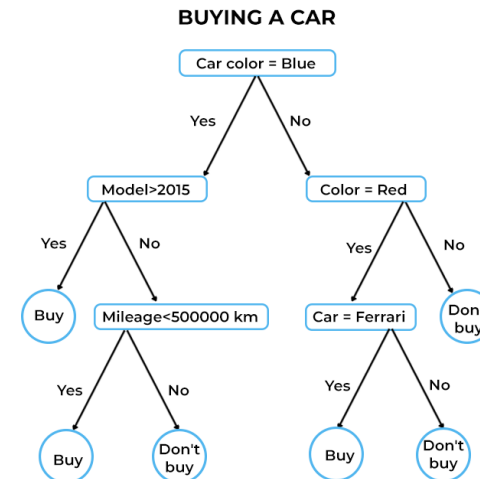
# Data Analytics Tasks

**Descriptive methods**

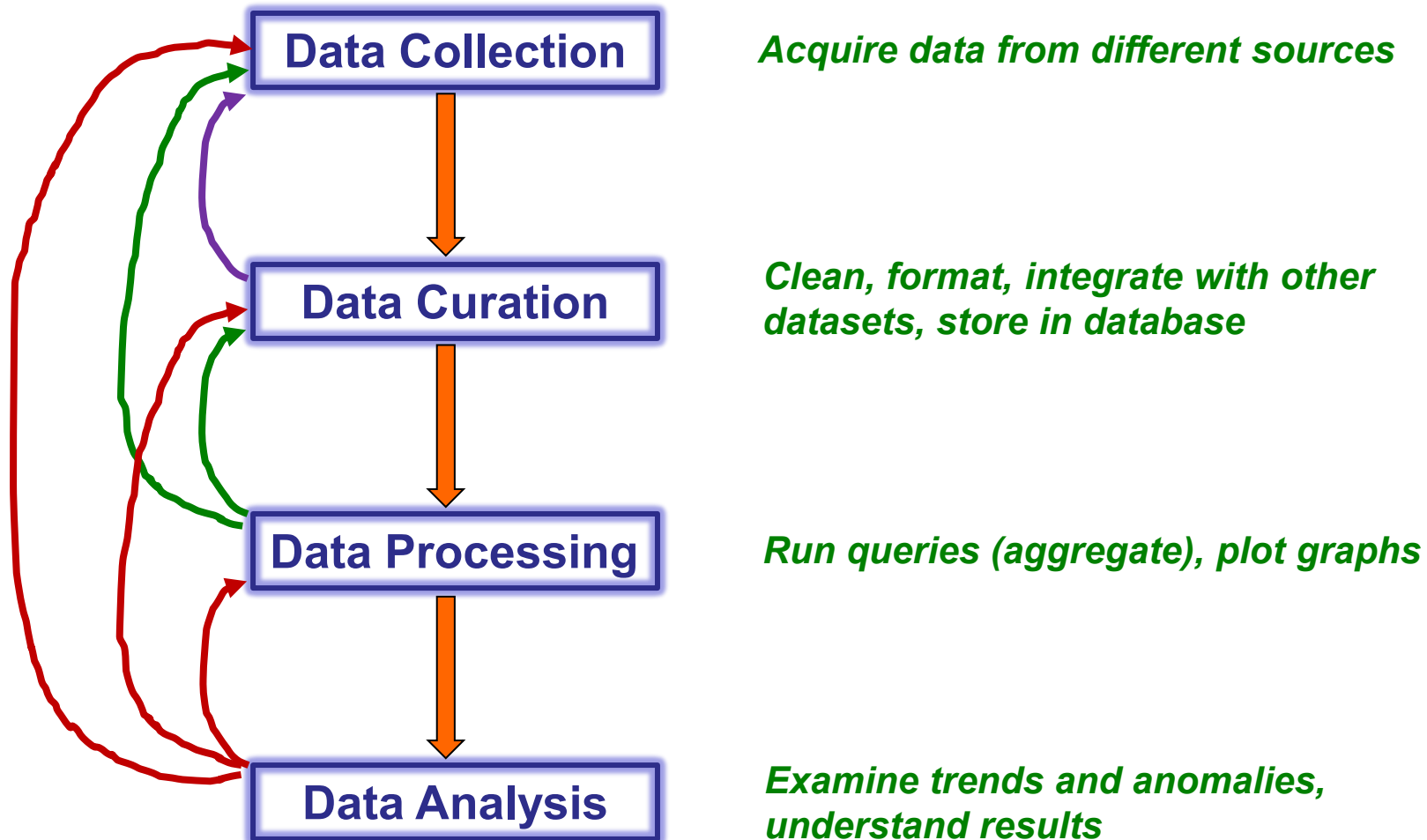- Find human-interpretable patterns that describe the data

- Example: **Clustering**

**Predictive methods**

- Use some variables to predict unknown or future values of other variables
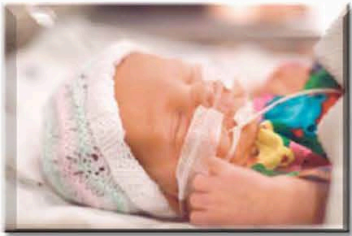
- Example: **Classification**



**BUYING A CAR**

# Data Analytics Pipeline

**Data Collection** — *Acquire data from different sources*

**Data Curation** — *Clean, format, integrate with other datasets, store in database*

**Data Processing** — *Run queries (aggregate), plot graphs*

**Data Analysis** — *Examine trends and anomalies, understand results*

# Big Data Applications

# Big Data Applications(I): Logistic



What does end-to-end logistics planning look like?

**Supplier → Manufacture →Distributor→Customer**

**Transport capacity**

**The delivery time is affected by traffic and weather conditions.**

**Storage capacity and price**

**Accuracy of the plan**

# Big Data Applications(II): Transportation

**Early Warning of Human Crowds Based on Query Data from Baidu Map: Analysis Based on Shanghai Stampede**

Jingbo Zhou, Hongbin Pei and Haishan Wu[*]

Baidu Research – Big Data Lab, Beijing, China

[Media Report: MIT Technology Review, Wall Street Journal, South China Morning Post ]



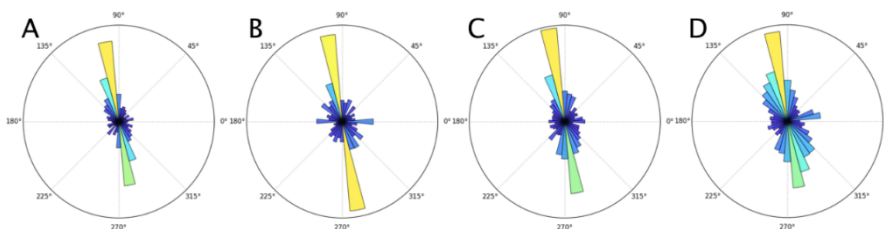Figure 2: Human population density between 23:00-24:00 on Dec. 31th 2014.



Figure 5: Human flow direction distribution in Chenyi Square (the specific disaster area of 2014 Shanghai Stampede) from 22:00 to 24:00 in: A – a common weekend (Aug. 23th 2014); B – the eve of the Mid-Autumn Festival (Sept. 7th 2014); C – the China's National Day (Oct. 1st 2014) and D – New Year's Eve of 2014

## Integration of transportation data

Multiple sources: car, taxi, bus, pedestrian, sensor

Multiple organizations: telecom corporation, taxi company, bus company, government

Data sharing and integrating

## Transportation planning

Construction of new roads

Location of transport junction

Answer "what-if" questions
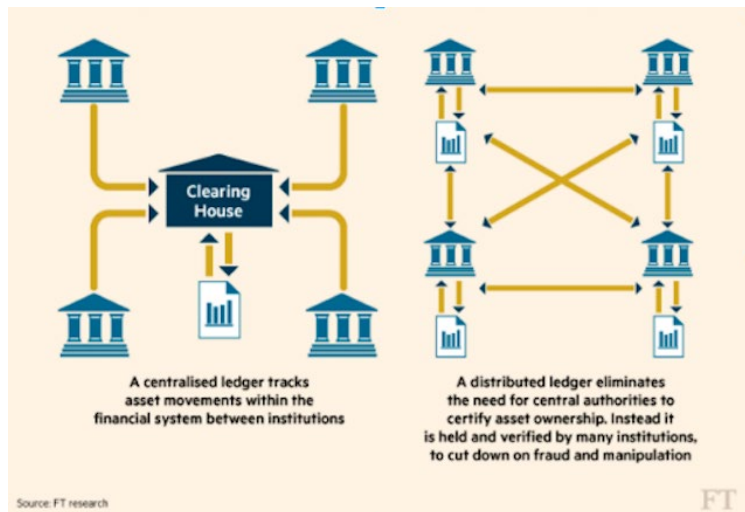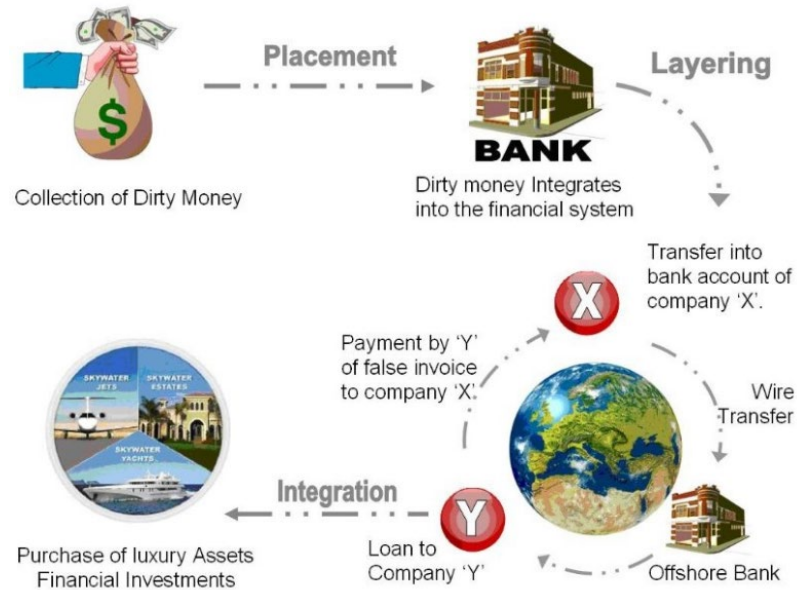
## Transportation management

Prevent traffic jam

Optimize traffic lights

Direct human crowds

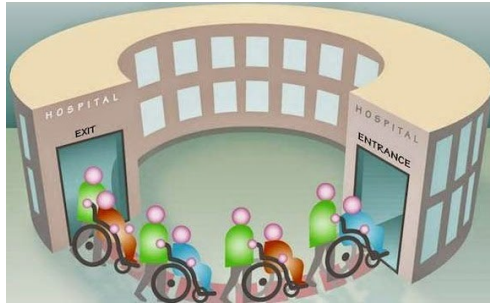Identify bottlenecks
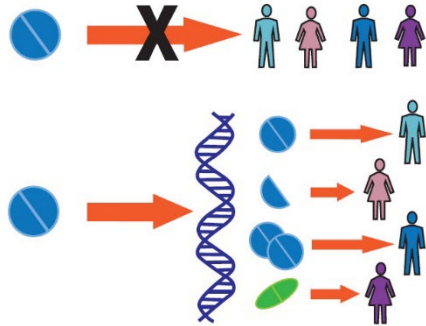
# Big Data Applications(III): Finance



- Finance Prediction/Policy

  Cash flow

**Abnormity Detection**

  Fraud

  Money laundering

  Tax evasion

**Fintech (financial technology )**

  Blockchain

  P2P loan

# Big Data Applications(IV): Retail Analytics

# Big Data Applications(V): Medical







**Medicine control**

> Drug allergy

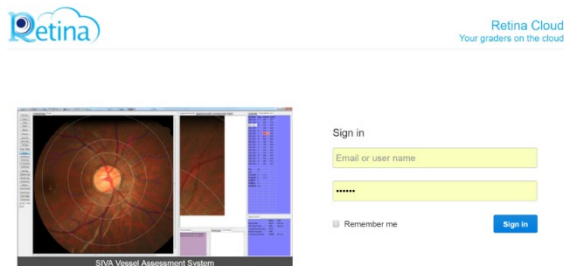> Drug and Poison Analysis

**Personal Medicine**

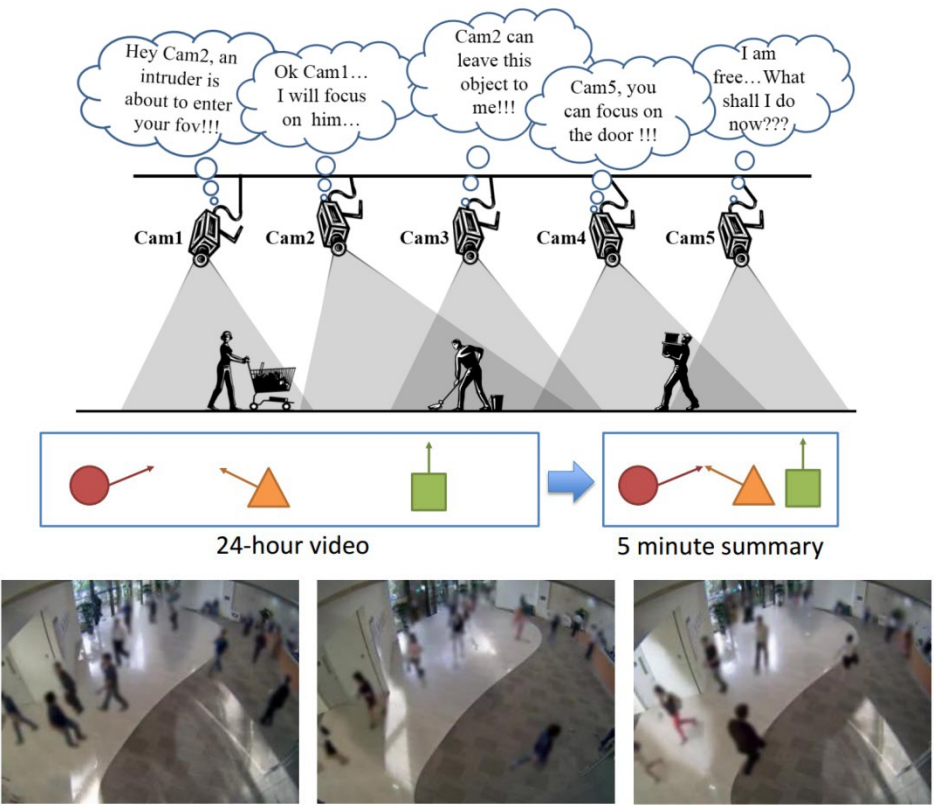**Hospital/Clinic management**

> Medical record

> Probability of re-hospitalization

**Doctor on the Cloud: Retinal-scan analysis**

> https://retinacloud.d1.comp.nus.edu.sg/users/sign_in

# Big Data Applications(VI): Security
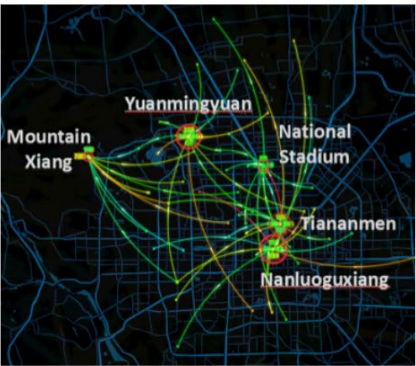


**Monitoring**
  CCTV
  IC cards
**Facial recognition to detect strangers**
**Exit passageway monitoring**
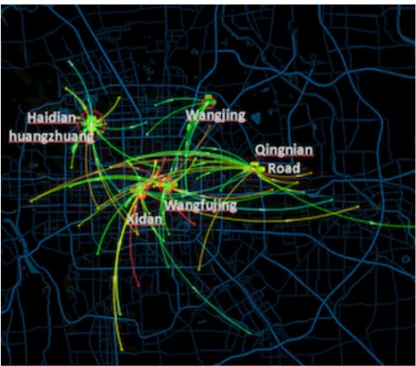**Crime analysis**

(a) all passengers      (b) visitors      (c) shoppers      (d) thieves
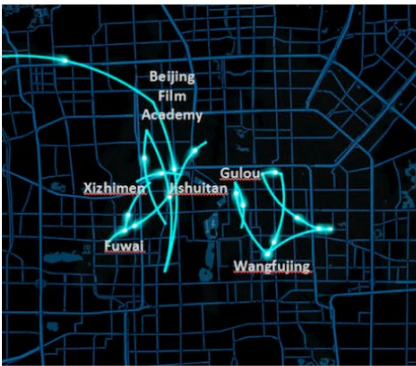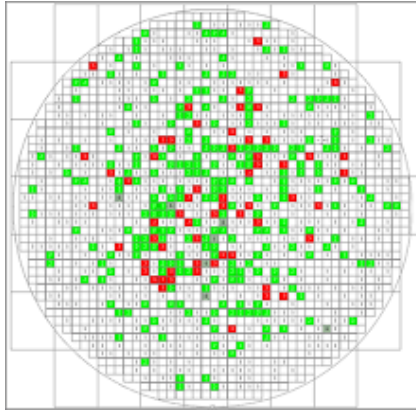
# Big Data Applications(VII): Manufacturing





**High returns products**

Wafer

Petroleum

**Manufacturing data**

Product imaging

Machine sensors

Machine repaired logs

**Usage/Applications**

Predictive Maintenance

Which machine affected the quality?

Which part of the machine needs to be repaired?

How to fully utilize the machines?

Product quality control

# Course Focus

- **Handle data that cannot fit in main memory**
  - Scalability of algorithms
  - Cluster computing architecture

- **Real world problems**
  - **Market basket analysis**
    - Finding **frequent itemsets**
  - **Customer segmentation**
    - **Clustering** large high dimensional data
  - **Recommender engines**
    - **Similarity Search**

# Course Focus

- **Tools and Techniques**

  - **Hadoop ecosystem**

    - Open source framework for **distributed processing** of large datasets

    - **Hadoop Distributed File System (HDFS)** for reliability and availability

    - **MapReduce, a Data-parallel programming model** to operate on large amounts of data

  - **Apache Spark**

    - Unified engine for distributed data processing

    - Fast in-memory processing and iterative processing with RDDs (Resilient Distributed Datasets)

  - **Search engine technology**

    - **Google's PageRank**, link-spam detection, hubs and authorities

# Assessment – 100% CA

- **Team-based Project (100%) --------- Max 2 members per team**

  - Project Proposal (20%)

    - *Proposal Presentation (10%)*

    - *Proposal Writeup(10%) ---------- 2 pages*

  - Project Updates (10%)

  - Final Project Presentation (20%)

  - Final Project Report(30%) ---------- 8 pages

  - Active Participation(20%)

*You are reminded **Plagiarism** is a very **SERIOUS** offence, and disciplinary action (including possibility of expulsion from the university) will be taken against any individual or team found plagiarizing.*

# Timetable(Approximate)

| Week No. | Date | Topic | Comments |
|---|---|---|---|
| 1 | 12th Aug | Introduction/Data | |
| 2 | 19th Aug | Hadoop /MapReduce | |
| 3 | 26th Aug | *Similarity Search | |
| 4 | 2nd Sep | *Frequent Items/Association Rules | Project Grouping Finalized 2$^{nd}$ Sep 23:59 |
| 5 | 9th Sep | *Clustering & Anomaly Detection(I) | |
| 6 | 16th Sep | *Clustering & Anomaly Detection(II) | |
| | BREAK | | |
| 7 | 30th Sep | Project Proposal(F2F Presentation) | |
| 8 | 7th Oct | *Classification/Regression I | |
| 9 | 14th Oct | *Classification/Regression II | |
| 10 | 21th Oct | Project Update (F2F Presentation) | |
| 11 | 28th Oct | *Graph Mining I | |
| 12 | 4th Nov | *Graph Mining II | |
| 13 | 11th Nov | Project Presentation (F2F Presentation) | Project Report, 17$^{th}$ Nov. 23:59 |

**\* Office Hour/Project Consultation Available from 20:35 to 22:00hrs**

# Desiderata of a Good Project

- **Innovation**
  - New Applications
  - New Algorithms
  - New Ways of looking at old problems
- **Complexity**
  - Application Complexity
  - Algorithm Complexity
  - Data Complexity
- **Technical Depth**
  - None trivial implementation
  - Thorough experiments and analysis