# CS5344
# Big Data Analytics Technology

# Class Information

- **Lecturer: Lee Mong Li**

  - Email: leeml@comp.nus.edu.sg

- **Tutors:**

  - **Gao Qiao** (email: gaoqiao@comp.nus.edu.sg)

  - **Suman Bhoi** (email: e0267909@u.nus.edu)

- **Lectures on Tuesday 1830 – 2030 hrs**

- **Course website @ IVLE**

- **Reference text**

  - Mining of Massive Datasets by J. Leskovec, A. Rajaraman and J.D. Ullman (available online: http:///www.mmds.org)

# Course Focus

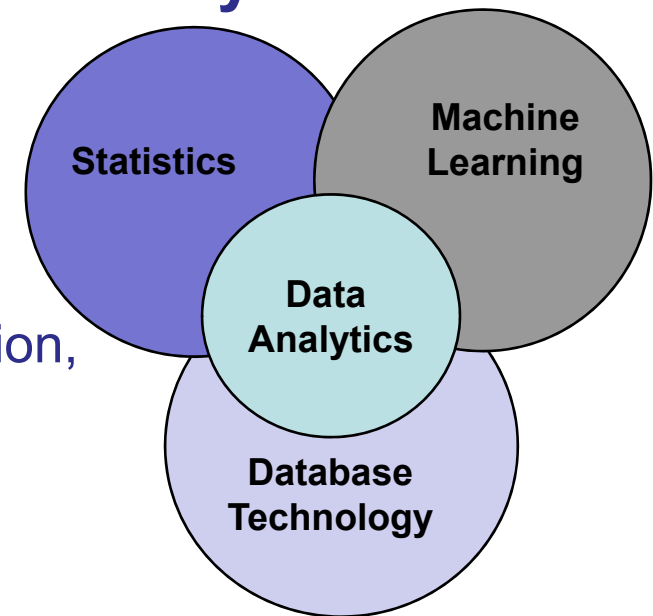- **Handle data that cannot fit in main memory**
  - Scalability of algorithms
  - Computing architecture
- **Real world problems**
  - Market basket analysis, Market segmentation, Recommender systems, Spam detection
- **Tools and Techniques**
  - MapReduce/Hadoop, Spark
    - create parallel algorithms to operate on large amount of data
  - Google's PageRank

# Assessment

- **100% CA**

- **Lab Assignments (30%)**

- **Written Assessments (30%)**

- **Team-based Project (40%)**

*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.*
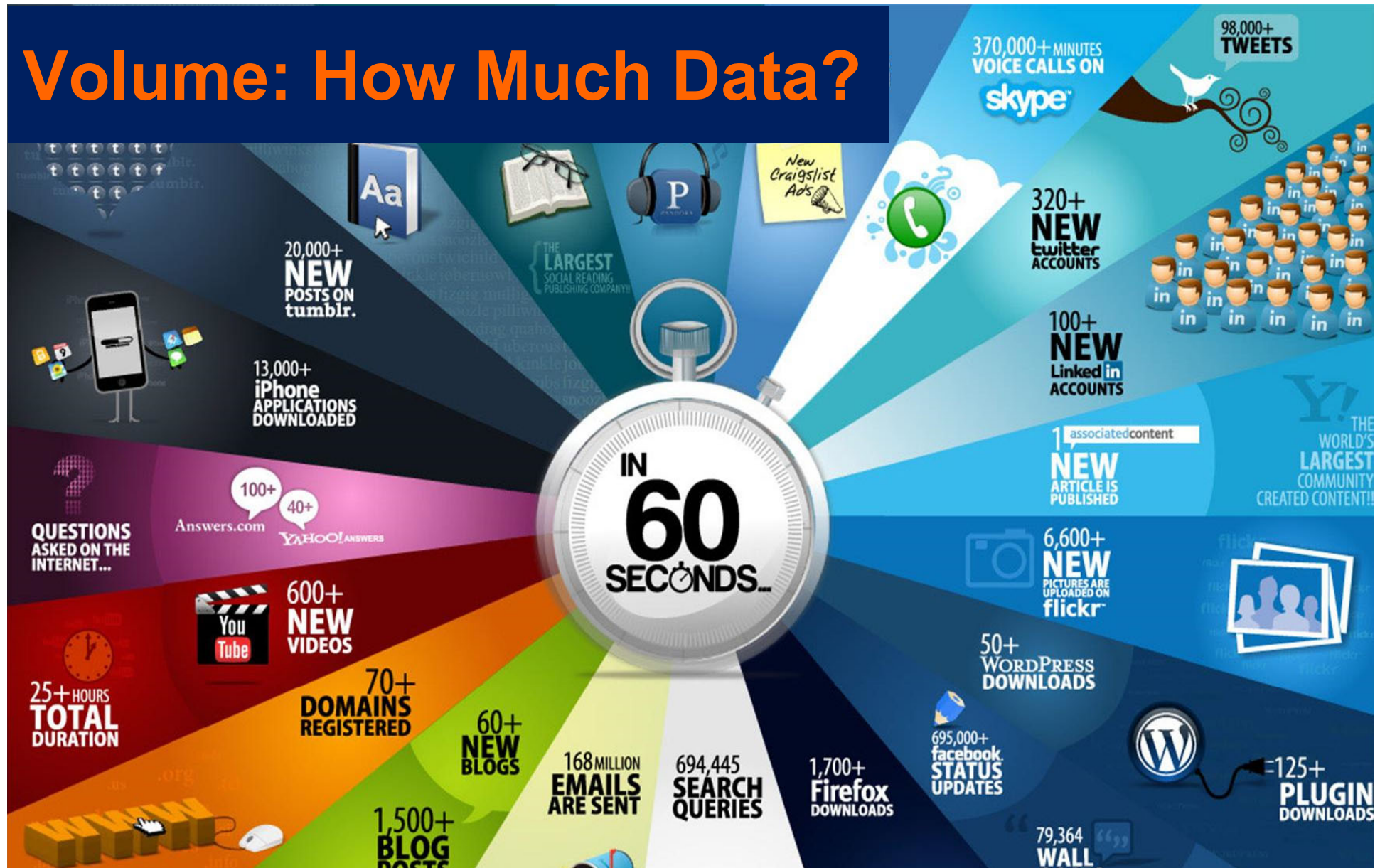
# What is Big Data?

- **Gartner's Definition**

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

- **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.**
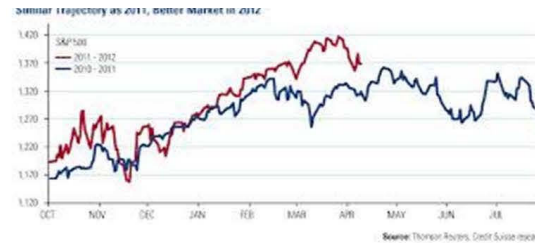
# Volume: How Much Data?



In 60 Seconds...

- 98,000+ TWEETS
- 370,000+ MINUTES VOICE CALLS ON skype
- 320+ NEW twitter ACCOUNTS
- 100+ NEW LinkedIn ACCOUNTS
- 1 NEW associatedcontent ARTICLE IS PUBLISHED
- 6,600+ NEW PICTURES ARE UPLOADED ON flickr
- 50+ WORDPRESS DOWNLOADS
- 695,000+ facebook STATUS UPDATES
- =125+ PLUGIN DOWNLOADS
- 79,364 WALL
- 1,700+ Firefox DOWNLOADS
- 694,445 SEARCH QUERIES
- 168 MILLION EMAILS ARE SENT
- 1,500+ BLOG POSTS
- 60+ NEW BLOGS
- 70+ DOMAINS REGISTERED
- 25+ HOURS TOTAL DURATION
- 600+ NEW VIDEOS
- 100+ 40+ QUESTIONS ASKED ON THE INTERNET... Answers.com YAHOO! ANSWERS
- 13,000+ iPhone APPLICATIONS DOWNLOADED
- 20,000+ NEW POSTS ON tumblr.
- THE LARGEST SOCIAL READING PUBLISHING COMPANY!!
- New Craigslist Ads
- THE WORLD'S LARGEST COMMUNITY CREATED CONTENT!!

- 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.

# 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

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.

# Velocity: At What Speed?

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

# Fourth V - Veracity

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

- **Bias, inconsistencies**

- **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?

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

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

- **Can create Value**

  - 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)

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

$600 to buy a disk drive that can store all of the world's music

# Why Big Data?

*Data contains Value and Knowledge*

# Big Data Analytics

- **From raw data to actionable information**

- **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?*

**Data Mining ≈ Big Data ≈
Predictive Analytics ≈ Data Science**

# Data Analytics/ Data Mining

- **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 Mining 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: **Recommender systems**

# Big 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*

# Data Integration and Cleaning

## Garbage in ➡ Garbage out

- **Quality of results relates directly to quality of data**

- **50% to 70% of analytics process effort is spent on data integration and cleaning**

- **Problems: duplicate records, entity resolution, conflict resolution, missing values, outliers, etc**

# Data from Different Sources



- **Different name representations**

# Data from Different Sources



- **Erroneous attribute values**

# Data from Different Sources



**No opening hours**

**Incomplete information**

# Data from Different Sources



- **Ambiguous references**

# Application of Big Data Analytics

# Acquiring Better Customers



Source: https://www.youtube.com/watch?v=BfoJgoItd4M

# Improving Customer Experience



Source: https://www.youtube.com/watch?v=BfoJgoItd4M

# Summary

- Lots of Buzz

- With good reason
  - Great potential
  - Many challenges

# Lab 1 (5%)

- **Get started with Spark**

- **Compile and execute a simple Spark program**

  - WordCount

- **Write your own Spark program**

  - Count the number of words that begin with each letter

- **Due: Tuesday, 29 January**

- **Submit to IVLE Lab 1 Folder by 6 PM**

# Apache Spark

- **Big Data is diverse and messy.**

- **Typical pipeline**
  - MapReduce for data loading and batch processing
  - Exploratory SQL-like queries
  - Iterative machine learning

- **Specialized engines create complexity and inefficiency**
  - Users must stitch together disparate systems

- **Spark is a unified engine for distributed data processing**
  - Programming model similar to MapReduce
  - Data-sharing abstraction called **Resilient Distributed Databases (RDDs)** to capture range of processing workloads that previously need separate engines (SQL, machine learning, graph processing)