

# Big data for internet applications

# Introduction to Big Data

Based on “Big Data: Hype or Hallelujah?” by Elena Baralis  
[http://dbdmg.polito.it/wordpress/wp-content/uploads/2010/12/BigData\\_2015\\_2x.pdf](http://dbdmg.polito.it/wordpress/wp-content/uploads/2010/12/BigData_2015_2x.pdf)

# Big data



# Google Flu trends




- February 2010
  - Google detected flu outbreak two weeks ahead of CDC data (Centers for Disease Control and Prevention – U.S.A)
  - Based on the analysis of Google search queries




# Data on the Internet...


## ■ Internet live stats


- <http://www.internetlivestats.com/>


  
**4,159,950,855**  
Internet Users in the world


  
**1,756,781,513**  
Total number of Websites


  
**185,422,700,748**  
Emails sent [today](#)


  
**4,678,531,893**  
Google searches [today](#)


  
**4,436,695**  
Blog posts written [today](#)


  
**539,254,323**  
Tweets sent [today](#)


  
**4,991,515,761**  
Videos viewed [today](#)  
on YouTube


  
**57,798,739**  
Photos uploaded [today](#)  
on Instagram


  
**96,105,235**  
Tumblr posts [today](#)


  
**2,436,676,131**  
Facebook active users

  
**682,245,690**  
Google+ active users

  
**344,445,798**  
Twitter active users

  
**278,737,344**  
Pinterest active users

  
**235,338,560**  
Skype calls [today](#)

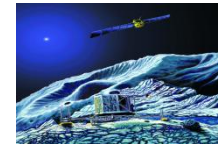
  
**87,939**  
Websites hacked [today](#)

# Who generates big data?

- User Generated Content (Web & Mobile)
  - E.g., Facebook, Instagram, Yelp, TripAdvisor, Twitter, YouTube

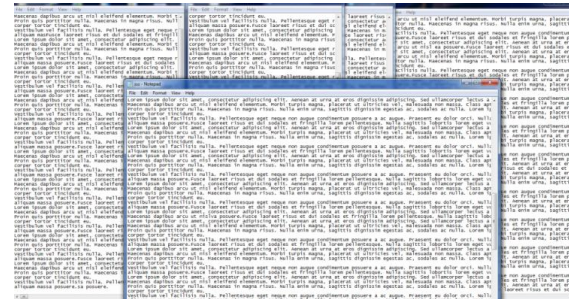


- Health and scientific computing

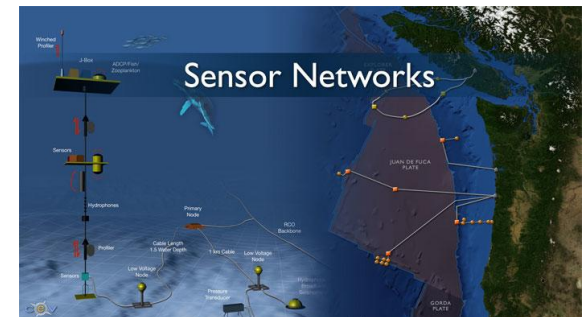
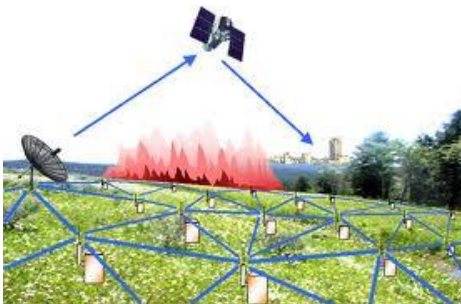


# Who generates big data?

- Log files
  - Web server log files, machine system log files



- Internet Of Things (IoT)
  - Sensor networks, RFID, smart meters





# An example of Big data at work

## ■ Crowdsourcing



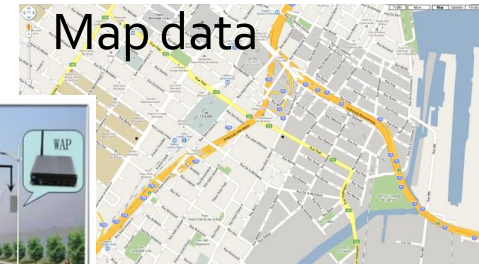
Computing



Real time traffic info



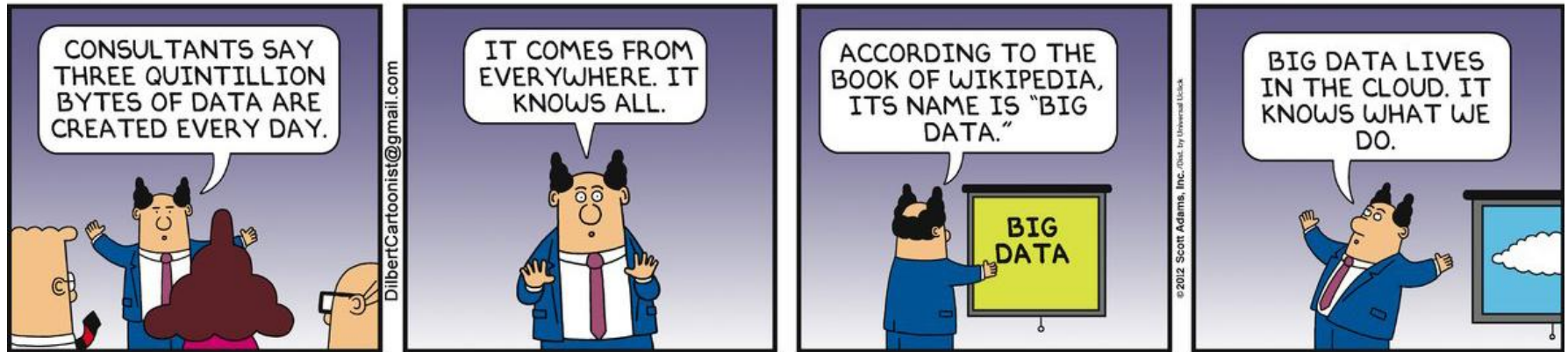
Sensing



Map data

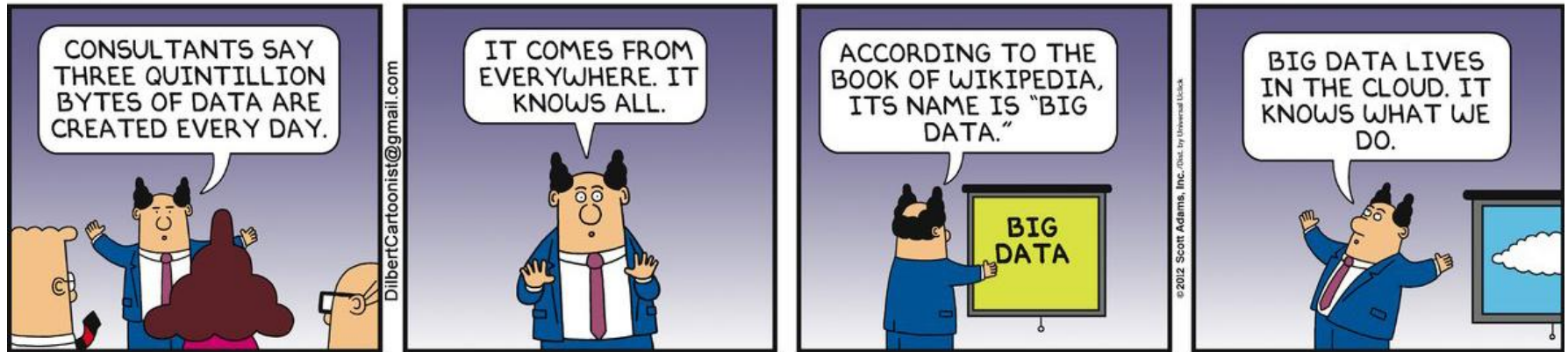


# What is big data?



- Many different definitions
  - "Data whose scale, diversity and complexity require new architectures, techniques, algorithms and analytics to manage it and extract value and hidden knowledge from it"

# What is big data?



- Many different definitions
  - “Data whose **scale**, **diversity** and **complexity** require new architectures, techniques, algorithms and analytics to manage it and extract value and hidden knowledge from it”

# What is big data?



- Many different definitions
  - "Data whose scale, diversity and complexity require new **architectures**, **techniques**, **algorithms** and **analytics** to manage it and extract value and hidden knowledge from it"

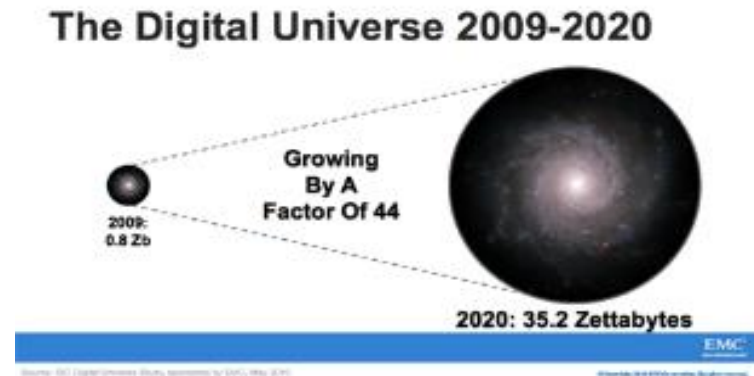
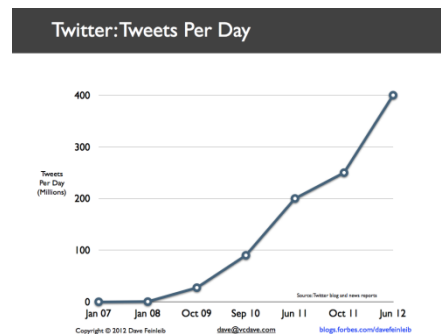
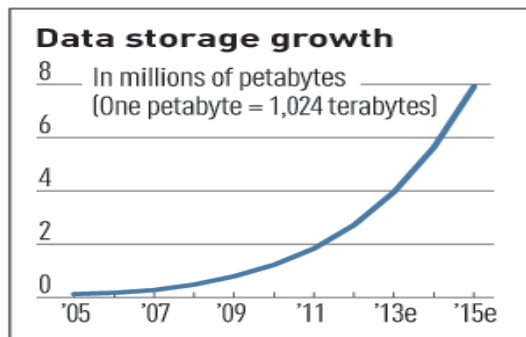
# The Vs of big data

- The 3Vs of big data
  - **V**olume: scale of data
  - **V**ariety: different forms of data
  - **V**elocity: analysis of streaming data
- ... but also
  - **V**eracity: uncertainty of data
  - **V**alue: exploit information provided by data

# The Vs of big data

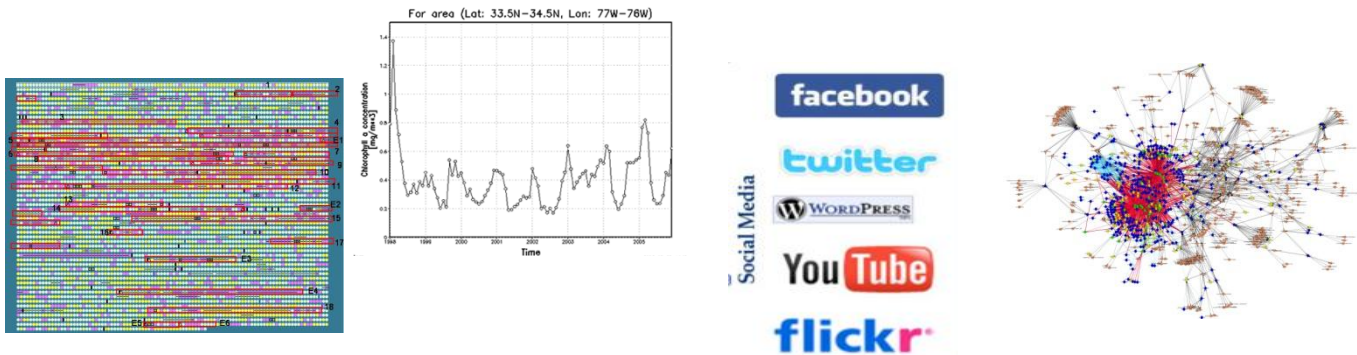
## ■ Volume

- Data volume increases exponentially over time
- 44x increase from 2009 to 2020
  - Digital data 35 ZB in 2020



# The Vs of big data

- **V**ariety
  - Various formats, types and structures
    - Numerical data, image data, audio, video, text, time series

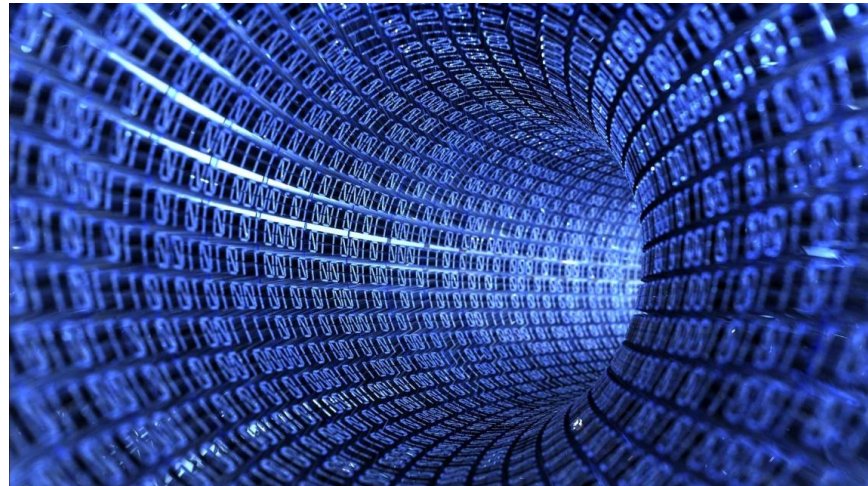


- A single application may generate many different formats
  - Heterogeneous data
  - Complex data integration problem



# The Vs of big data

- **V**elocity
  - Fast data generation rate
    - Streaming data
  - Very fast data processing to ensure timeliness



# The Vs of big data

- **V**eracity
  - Data quality

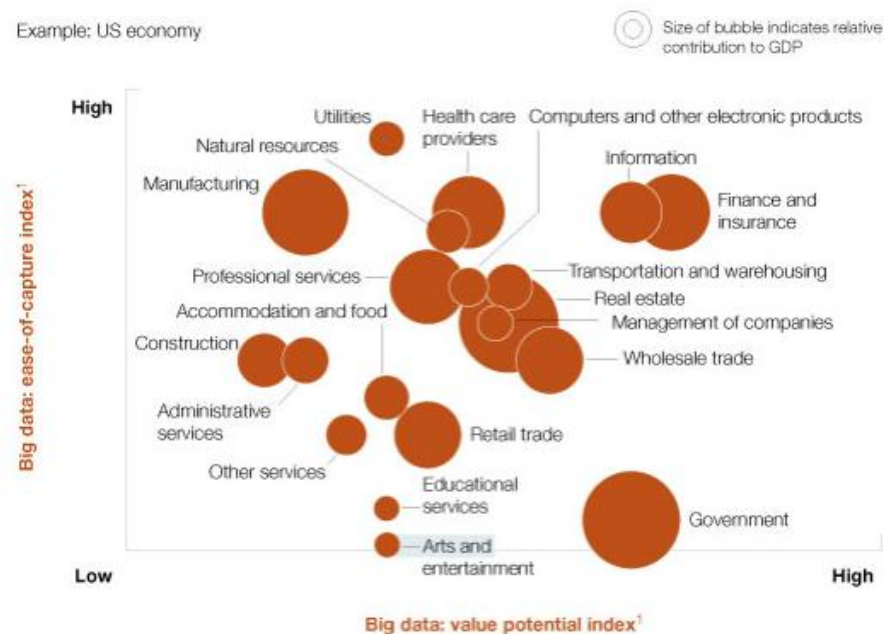


**Reliability**  
**Accuracy**  
**Timeliness**  
**Completeness**  
**Consistency**  
**Relevance**

Format  
Sufficiency  
Flexibility  
Conciseness  
Currency  
Comparability  
Scope  
Level-of-detail  
Precision  
Efficiency  
Quantitativeness  
Interpretability  
Understandability  
Usefulness  
Usableness  
Clarity  
Content  
Importance  
Informativeness  
Freedom from bias

# The Vs of big data

- Value
  - Translate data into business advantage



<sup>1</sup> For detailed explication of metrics, see appendix in McKinsey Global Institute full report *Big data: The next frontier for innovation, competition, and productivity*, available free of charge online at [mckinsey.com/mgi](http://mckinsey.com/mgi).

Source: US Bureau of Labor Statistics; McKinsey Global Institute analysis

# Big data value chain



## ■ Generation

- Passive recording
  - Typically structured data
  - Bank trading transactions, shopping records, government sector archives
- Active generation
  - Semistructured or unstructured data
  - User-generated content, e.g., social networks
- Automatic production
  - Location-aware, context-dependent, highly mobile data
  - Sensor-based Internet-enabled devices

# Big data value chain



- Acquisition
  - Collection
    - Pull-based, e.g., web crawler
    - Push-based, e.g., video surveillance, click stream
  - Transmission
    - Transfer to data center over high capacity links
  - Preprocessing
    - Integration, cleaning, redundancy elimination

# Big data value chain



## ■ Storage

- Storage infrastructure
  - Storage technology, e.g., HDD, SSD
  - Networking architecture, e.g., DAS, NAS, SAN
- Data management
  - File systems (HDFS), key-value stores (Memcached), column-oriented databases (Cassandra), document databases (MongoDB)
- Programming models
  - Map reduce, stream processing, graph processing




# Big data value chain



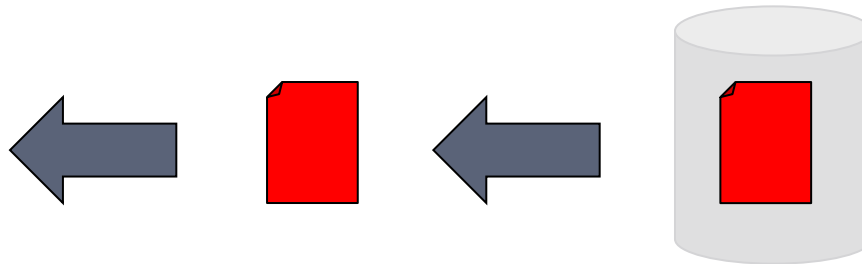
- Analysis
  - Objectives
    - Descriptive analytics, predictive analytics, prescriptive analytics
  - Methods
    - Statistical analysis, data mining, text mining, network and graph data mining
    - Clustering, classification and regression, association analysis
  - Diverse domains call for customized techniques

# Big data challenges

- Technology and infrastructure
  - New architectures, programming paradigms and techniques are needed
- Data management and analysis
  - New emphasis on “data”
  -  Data science

# The bottleneck

- Processors process data
- Hard drives store data
- We need to transfer data from the disk to the processor



# The solution

- **Transfer the processing power to the data**
- Multiple distributed disks
  - Each one holding a portion of a large dataset
- Process in parallel different file portions from different disks

