

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

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



Google Flu trends

google.org Flu Trends

[Google.org home](#)

Flu Trends

Select country/region

[Home](#)

[How does this work?](#)

[FAQ](#)

Flu activity

Intense
High
Moderate
Low
Minimal

Explore flu trends around the world

We've found that certain search terms are good indicators of flu activity. Google Flu Trends uses aggregated Google search data to estimate flu activity. [Learn more](#)



- February 2010
 - Google detected flu outbreak two weeks ahead of CDC data (Centers for Disease Control and Prevention)


Nowcasting




Data on the Internet...


- Internet live stats


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



4,485,508,861
Internet Users in the world



1,752,142,970
Total number of Websites



193,688,718,339
Emails sent *today*



5,222,289,027
Google searches *today*



4,990,992
Blog posts written *today*



572,159,945
Tweets sent *today*



5,348,093,035
Videos viewed *today*
on YouTube



62,832,046
Photos uploaded *today*
on Instagram



107,175,151
Tumblr posts *today*



2,435,900,914
Facebook active users



795,537,418
Google+ active users



357,398,865
Twitter active users



278,573,312
Pinterest active users


287,236,096
Skype calls *today*


110,341
Websites hacked *today*


5,664,059,486 GB
Internet traffic *today*


3,065,544 MWh
Electricity used *today*
for the Internet

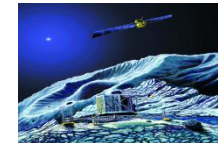

2,507,959 tons
CO₂ emissions *today*

Who generates big data?

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

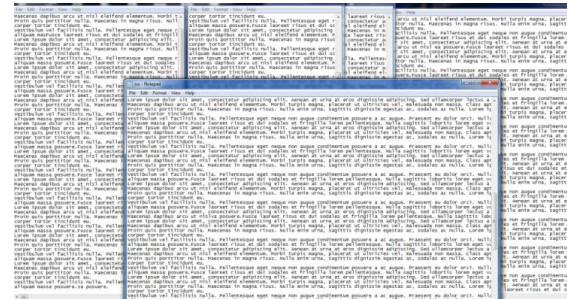


- Health and scientific computing

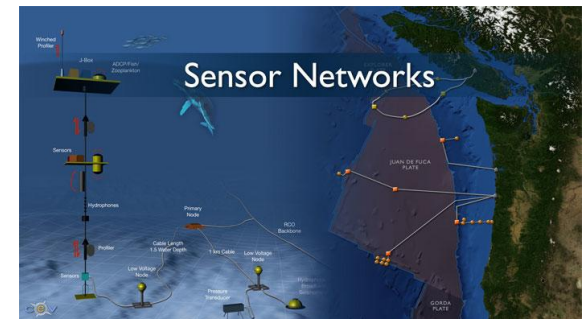
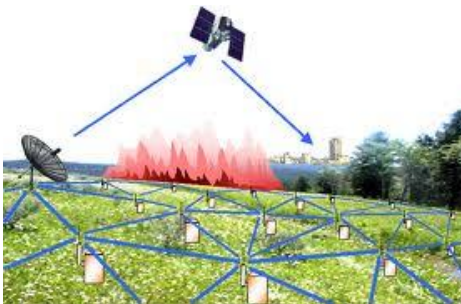


Who generates big data?

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

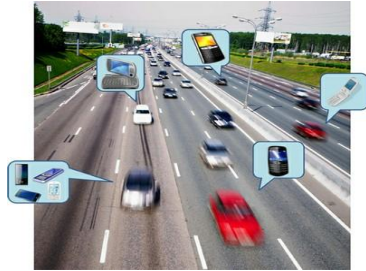


- Internet Of Things (IoT)
 - Sensor networks, RFIDs, smart meters



An example of Big data at work

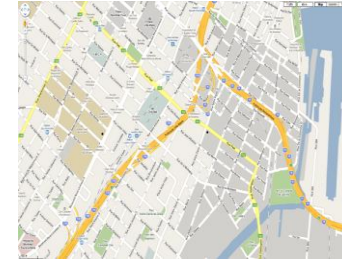
Crowdsourcing



Sensing



Map data



Computing




Real time traffic info

Travel time forecast/nowcast

The Vs of big data

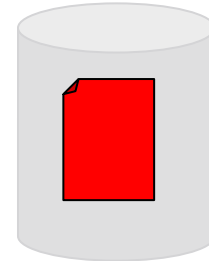
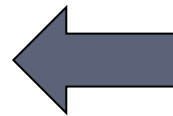
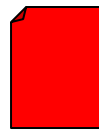
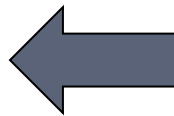
- The 3Vs of big data
 - **V**olume: scale of data
 - **V**ariety: different forms of data
 - **V**elocity: analysis of streaming data
- ... but also
 - **V**eracity: uncertainty of data
 - **V**alue: exploit information provided by data

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

