



Big Data: Hype or Hallelujah?



Elena Baralis
Politecnico di Torino



Big data hype?



2

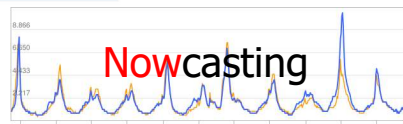


Google Flu trends



- February 2010
 - detected flu outbreak two weeks ahead of CDC data

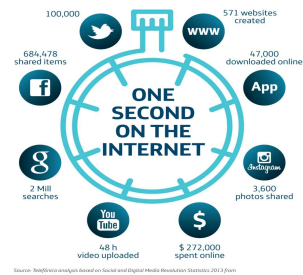
Discontinued!



3



On the Internet...



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



4



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"



5



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"*



6

What is big data?



CONSULTANTS SAY THREE QUINTILLION BYTES OF DATA ARE CREATED EVERY DAY.

IT COMES FROM EVERYWHERE. IT KNOWS ALL.

ACCORDING TO THE BOOK OF WIKIPEDIA, ITS NAME 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"

DBG

7

What is big data?



CONSULTANTS SAY THREE QUINTILLION BYTES OF DATA ARE CREATED EVERY DAY.

IT COMES FROM EVERYWHERE. IT KNOWS ALL.

ACCORDING TO THE BOOK OF WIKIPEDIA, ITS NAME 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"

DBG

8

The Vs of big data

- The 3Vs of big data
 - Volume: scale of data
 - Variety: different forms of data
 - Velocity: analysis of streaming data
- ... but also
 - Veracity: uncertainty of data
 - Value: exploit information provided by data

DBG

9

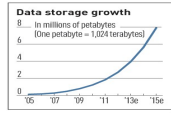
The Vs of big data

terabytes petabytes exabytes zettabytes

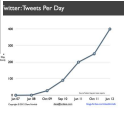
↑
the amount of data stored by the average company today

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

The Digital Universe 2009-2020



Data storage growth
8... in millions of petabytes (One petabyte = 1,024 terabytes)



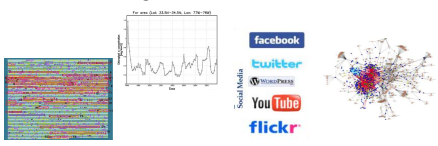
2009: 1.4 ZB
2020: 35.2 Zettabytes
Growing By A Factor Of 44

DBG

10

The Vs of big data

- Variety
 - Various formats, types and structures
 - Numerical data, image data, audio, video, text, time series



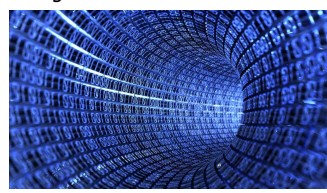
- A single application may generate many different formats

DBG

11

The Vs of big data

- Velocity
 - Fast data generation rate
 - Streaming data
 - Very fast data processing to ensure timeliness




DBG

12

The Vs of big data

- **Veracity**
 - Data quality

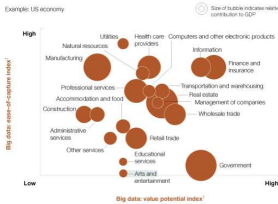


Reliability
Format Sufficiency Flexibility Timeliness Accuracy Currency Consistency Precision Relevance Completeness Precision Relevance Consistency Precision Relevance

DBG 13

The Vs of big data


- **Value**
 - Translate data into business advantage



DBG 14

Who generates big data?

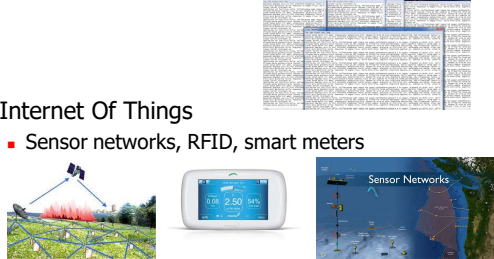
- **User Generated Content (Web & Mobile)**
 - E.g., Facebook, Instagram, Yelp, TripAdvisor, Twitter, YouTube
- **Health and scientific computing**



DBG 15

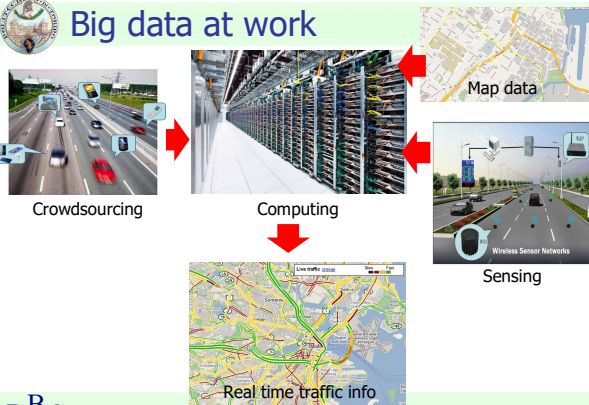
Who generates big data?

- **Log files**
 - Web server log files, machine syslog files
- **Internet Of Things**
 - Sensor networks, RFID, smart meters



DBG 16

Big data at work

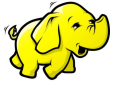


DBG 17

Big data challenges

- **Technology & infrastructure**
 - New architectures, programming paradigms and techniques are needed
 - *Transfer the processing power to the data*
 - Hadoop ecosystem
- **Data management & analysis**
 - New emphasis on "data"


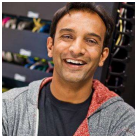
➔ **Data science**



DBG 18

Data science

"Extracting meaning from very large quantities of data"

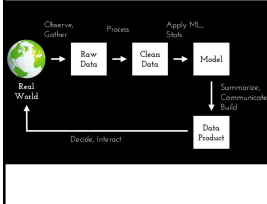

D.J. Patil coined the word *data scientist*

DBG

19

The data science process


AKA *KDD* process
Knowledge Discovery in Databases

DBG

20

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

DBG

21

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

DBG

22

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

DBG

23

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

DBG

24

Large scale data processing

- Traditional approach
 - Database and data warehousing systems
 - Well-defined structure
 - Small enough data
- Big data
 - Datasets not suitable for databases
 - E.g. google crawls
 - May need near real-time (streaming) analysis
 - Different from data warehousing
 - Different programming paradigm

DBG 25

Large scale data processing

- Traditional computation is *processor bound*
 - Small dataset
 - Complex processing
- How to increase performance?
 - New and faster processor
 - More RAM

DBG 26

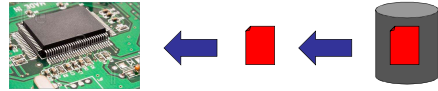
Large scale data processing

- Traditional data storage
 - On large SANs
 - Data transferred to processing nodes on demand at computing time
- Traditional distributed computing
 - Multiple machines, single job
 - Complex systems
 - E.g., MPI
 - Programmers need to manage data transfer synchronization, system failure, dependencies

DBG 27

The bottleneck

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



DBG 28

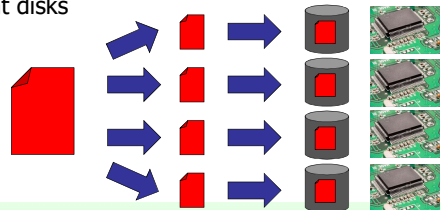
The bottleneck

- Hard drives evolution
 - Storage capacity increased fast in recent decades
 - E.g., from 1GB to 1TB
 - The transfer rate increased less
 - E.g., from 5MB/s to 100MB/s
- Transfer of disk content in memory
 - Few years ago: 3.33 min.
 - Now: 2.7 hours (if you have enough RAM)
- Problem: *data transfer from disk to processors*

DBG 29

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



DBG 30


Issues

- Need to manage
 - Hardware failures
 - Network data transfer
 - Data loss
 - Consistency of results
 - Joining data from different disks
 - Scalability
- Managed by Hadoop



DBG 31

Apache Hadoop



- Open source project by the Apache Foundation
- Based on 2 Google papers
 - Google File System (GFS), published in 2003
 - Map Reduce, published in 2004
- Reliable storage and processing system based on YARN (Yet Another Resource Negotiator)
 - Storage provided by HDFS
 - Different processing models
 - E.g., Map Reduce, Spark, Spark streaming, Hive, Giraph

DBG 32

Hadoop scalable approach

- Data distributed across nodes automatically
 - When loaded into the system
- Processing executed on local data
 - Whenever possible
- No need of data transfer to start the computation
- Data automatically replicated
 - For availability and reliability
- Developers only focus on the logic of the problem to solve

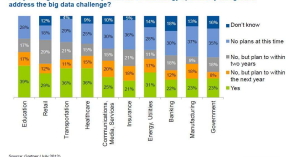
DBG 33

Conclusions

- Certainly not just hype

Big Data Investments by Industry

Has your organization already invested in technology specifically designed to address the big data challenges?



Source: Gartner (July 2012)

■ ... but not a panacea!

DBG 34