

# **Big Data**

### overview, issues, challenges and opportunities

C. Onime (onime@ictp.it)





# Outline

- Interactive session
  - Introduction to Big-Data
  - Issues/challenges
  - Taxonomy classifications
- Conclusion

- Opportunities and future



### Pre-exercise

- Before providing a formal definition, let's try answer the questions:
  - What exactly is Big-Data?
  - Can you identify it?





# Definition(s)

- The term Big-Data by definition is used for data that is "massive" in one of the following areas:
  - Volume: quantity
  - Velocity: generated at high speed
  - Variety: wide spread from diverse sources and types.
  - Variability: constantly changing meaning
  - Veracity: making data accurate (removing bad data)
  - Visualization: presenting and conveying meaning
  - Value: applying findings and taking action





# **Big-Data examples**

- Astronomical Image data from a telescope exceeds 1TB/day
- Environmetal monitoring
- Government: Census, National Health Records/Systems, etc.
- Industry: Amazon, Google, Ebay...





### World wide storage







### Another forecast



- 0.076 ZB = 76 EB
- 76 EB = 76M PB
- Current estimate is that 82% of global IP traffic will be video by 2020



### Preamble

- So what is driving Big Data?
  - Mainly industry related paradigms & applications
    - Data mining, Business Intelligence, Knowledge Management and now Big Data Management





# Data Mining

 A process of analyzing data from different perspectives and summarizing it into useful information, [...] which allows users to analyze data from many different dimensions or angles, categorize it, and summarize the relationships identified.





# **Business Intelligence**

 A process of finding, gathering, aggregating and analyzing information for decisionmaking. It makes use of a set of technologies that allow the acquisition and analysis of data to improve company decision making and work flows.





# Knowledge Management

- A business process that formalizes the management and use of an enterprise's intellectual assets." KM promotes a collaborative and integrative approach to the creation, capture, organization, access and use of information assets, including the tacit, un-captured knowledge of people.
- A systematic process of finding, selecting, organizing, distilling and presenting information in a way that improves an employee's comprehension in a specific area of interest which supports an organization to gain insight and understanding from its own experience.





IAE.

### Big Data Management







# Other drivers

- Scientific Research
  - High Performance Computing (LHC, SKA, Genomics)
- Improvements in hardware technology
  - Heading towards Nano-circuits, clocking resolutions, etc
- Improvements in computing platforms
  - Networks: always connected devices, capacity; Clouds: anytime, anywhere on-demand metered access to resources
- Every user is a now a provider/consumer
  - Social networking





# Issues and challenges

- Perspectives backgrounds, use cases
- Taxonomies, ontologies, schemas, workflow
- Bits raw data formats and storage methods
- Cycles algorithms and analysis
- Infrastructure (screws) to support Big Data

- From presentation by Michael Cooper & Peter Mell of NIST



IAE.



#### Clement Onime - onime@ictp.it

CI



### Six dimensional Taxonomy



IAEA





# Data Mapping examples







### Compute infrastructure



IAEA





CI

### Overview of Hadoop MapReduce



IAEA





### Hadoop 2.0 Ecosystem











IAEA









- Tuple
  - Key-value pairs
- Streams
  - Sequence of tuples pairs
- Spout
  - Source of streams
- Bolt
  - Processing element
  - (filers, join, transform, e.t.c)









### Storm topology



- Graph of Computation
  - Network of spouts and bolt
  - Parallel & cyclic execution
- Groupings
  - Shuffle, all, Global, fields
- Example:
  - Twitter analytics: spout, bolts: parse, count, ranks, report





### Storage infrastructure



IAEA



![](_page_24_Picture_0.jpeg)

# Infrastructure

![](_page_24_Picture_2.jpeg)

IAEA

### mapping

![](_page_24_Figure_4.jpeg)

BATCH

#### NEAR-REAL-TIME

#### **REAL-TIME**

![](_page_25_Picture_0.jpeg)

# Storage complexity/size

![](_page_25_Picture_3.jpeg)

IAEA

![](_page_25_Figure_4.jpeg)

![](_page_26_Picture_0.jpeg)

### Analytics

![](_page_26_Picture_3.jpeg)

![](_page_26_Picture_4.jpeg)

![](_page_26_Figure_5.jpeg)

![](_page_27_Picture_0.jpeg)

Statistics	Machine learning
Model	Network, Graphs
Data point	Examples/instances
Response	Label
Parameters	Weights
Covariate	Feature
Fitting/Estimation	Learning
Test set performance	Generalization
Regression/Classification	Supervised Learning
Density estimation, Clustering	Unsupervised Learning

![](_page_28_Picture_0.jpeg)

### Visualisation

![](_page_28_Picture_3.jpeg)

IAEA

![](_page_28_Picture_4.jpeg)

![](_page_29_Picture_0.jpeg)

# Mixed Reality Environments

![](_page_29_Picture_2.jpeg)

![](_page_29_Figure_3.jpeg)

$$E_{MR} = \int (R+V) \quad \text{where} \quad E_{MR} = \begin{cases} E_R, & \text{if } V = 0\\ E_{AR}, & \text{if } R > V\\ E_{AV}, & \text{if } R < V\\ E_{VR}, & \text{if } R = 0 \end{cases}$$

![](_page_30_Picture_0.jpeg)

International Centre for Theoretical Physics

### VR and AR

![](_page_30_Picture_3.jpeg)

#### Virtual Reality (VR) CAVE

- Computer generated virtual environment
- Creates a completely virtual environment that is without real objects
- Portable
  - Headsets, wearable devices
  - Custom and typically not cost effective

#### Augmented Reality(AR)

- Real-time integration of computer generated information into a 3D world.
- Blends into real world and supports real objects
- Mobile
  - Commodity devices: smartphones and tablets
  - Cost effective

![](_page_31_Picture_0.jpeg)

### Some Examples

![](_page_31_Picture_3.jpeg)

![](_page_31_Picture_4.jpeg)

#### **VR Environments**

![](_page_31_Picture_6.jpeg)

#### **AR Environments**

![](_page_31_Picture_8.jpeg)

![](_page_31_Picture_9.jpeg)

![](_page_31_Picture_10.jpeg)

![](_page_32_Picture_0.jpeg)

### **AR** Cubicle

![](_page_32_Picture_2.jpeg)

![](_page_32_Picture_3.jpeg)

![](_page_32_Picture_4.jpeg)

![](_page_32_Picture_5.jpeg)

![](_page_32_Picture_6.jpeg)

![](_page_32_Picture_7.jpeg)

180° horizontal by 3 markers on walls and 90° vertical by marker on floor

![](_page_33_Picture_0.jpeg)

# Security and privacy

![](_page_33_Picture_3.jpeg)

![](_page_33_Picture_4.jpeg)

![](_page_33_Figure_5.jpeg)

![](_page_34_Picture_0.jpeg)

# Public Key Cryptography

![](_page_34_Picture_2.jpeg)

- Asymmetric cryptography
  - A pair of keys: one public and the other private
  - Useful for authentication and encryption
  - Depends mainly on the impracticability of computing the equivalent private key from its public component.
  - Public key may be freely exchanged without secure channels such as public key servers, etc..
  - Computationally intensive mathematical algorithms

![](_page_35_Picture_0.jpeg)

International Centre for Theoretical Physics

### Digital Certificates

![](_page_35_Picture_3.jpeg)

- Similar to travel passport
  - Provides forgery resistant identifying information
    - Name of holder
    - Serial number
    - Expiration date
    - Copy of holder's public key (used for encryption)
    - Digital signature of issuing authority (CA)

![](_page_36_Picture_0.jpeg)

### SSL Transport

![](_page_36_Picture_3.jpeg)

![](_page_36_Picture_4.jpeg)

![](_page_36_Figure_5.jpeg)

![](_page_37_Picture_0.jpeg)

### Data colouring

![](_page_37_Picture_2.jpeg)

![](_page_37_Figure_3.jpeg)

![](_page_38_Picture_0.jpeg)

![](_page_38_Picture_1.jpeg)

### Conclusion

- The potentiality of Big Data is now all around us in our everyday lives. Every device will be connected and constantly generating data.
- Good mapping of big-data is fundamental to understanding/selecting infrastructure (compute & storage), analytics, visualization and protection (security and privacy).
- New frontiers such as "Data Science" is bringing many of the ideas/techniques from Big-Data Analytics to almost any field or discipline.

![](_page_39_Picture_0.jpeg)

![](_page_39_Picture_1.jpeg)

# 2017 opportunities @ ICTP

- Workshop on Open Source Solutions for the Internet of Things – June 28 – July 7<sup>th</sup>, 2017
- The CODATA RDA Advanced School of Research Data Science for Extreme sources of Data, Bioinformatics and IoT/Big-Data Analytics – July 3rd -28<sup>th</sup>, 2017
- Two other CODATA/RDA schools on Data Science
  South Africa and Brazil , maybe a HPC school in Mexico
- Masters degree in HPC, Trieste, Italy
- Graduate studies @ East African Institute for Fundamental Research (EAIFR), Kigali, Rwanda

![](_page_40_Picture_0.jpeg)

![](_page_40_Picture_1.jpeg)

### References

- Michael Cooper & Peter Mell, "Tackling Big Data", NIST Information Technology Laboratory, 2010
- Big Data Working Group, "Big Data Taxonomy", Cloud Security Alliance, 2014
- M. Bornschlegl et al, "IVIS4BigData: A Reference Model for Advanced Visual Interfaces Supporting Big Data Analysis in Virtual Research Environments", 2016
- S. Rajendran, Apache Storm: A scalable distributed & fault tolerant real time computation system, 2015

![](_page_41_Picture_0.jpeg)

![](_page_41_Picture_1.jpeg)

![](_page_41_Picture_2.jpeg)

## That's all folks!!

questions

Clement Onime - onime@ictp.it