



WORKSHOP ON LARGE-SCALE TOMOGRAPHY

BIG DATA: SIZE DOES MATTER!

LAJOS RODEK
BIG DATA ARCHITECT, EPAM SYSTEMS, SZEGED
LAJOS_RODEK@EPAM.COM

JANUARY 26, 2016

DISCLAIMER



NO SCIENCE TODAY!

AGENDA

- 1 Introduction to Big Data
- 2 Big Data in practice
- 3 Technologies & tools
- 4 Conclusions





INTRODUCTION TO BIG DATA

DEFINITION OF BIG DATA



*“... a new generation of technologies and architectures designed to extract **value** economically from very large **volumes** of a wide **variety** of data by enabling **high-velocity** capture, discovery, and/or analysis.” (IDC, 2012)*

*“... **high-volume**, **-velocity** and **-variety** information assets that demand cost-effective, innovative forms of information processing for enhanced **insight** and **decision** making.” (Gartner, 2013)*

THE 3 V'S

Volume

Scale of data

Large & expanding

Many data sources

Velocity

Rate of data arrival

Rate of processing: offline
(batch) vs low-latency vs real-
time (stream)

Rate of changes

Variety

Structured vs unstructured vs
semi-structured data

Text vs binary data

“Dark data”

(Doug Laney, META Group / Gartner, 2001)

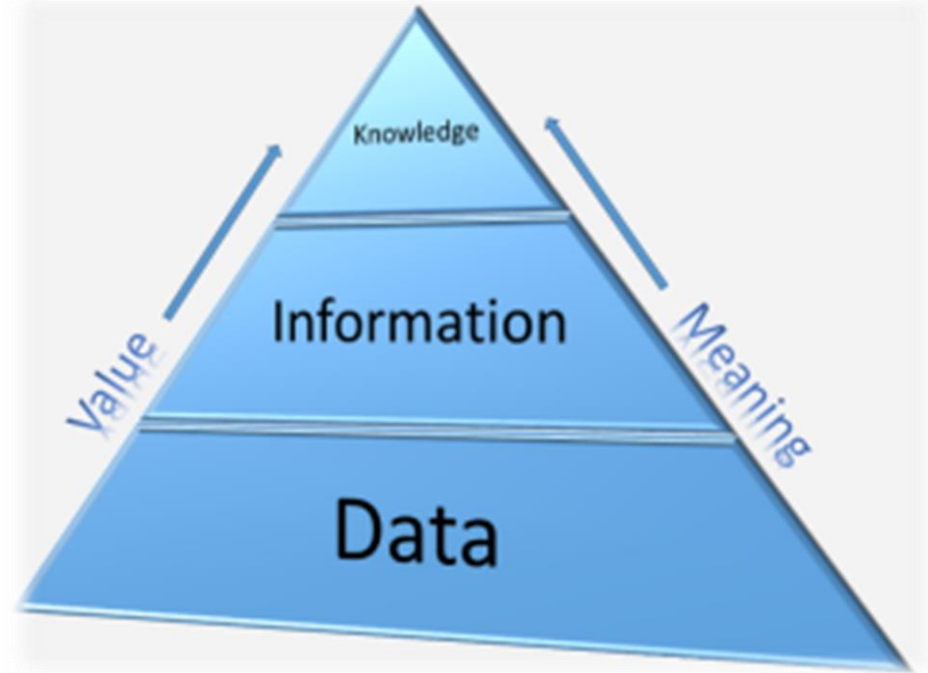
ONE MORE IMPORTANT V

Value

Relevance

Outcome

Actions



USE CASES

TOP BIG DATA USE CASES

Customer Financial Marketing Retail Security Pharma

Customer Analytics
48%

45%
Experience Analytics

Threat Analysis
30%

Risk Analysis
37%

28%
Regulatory Compliance Analysis

Campaign Optimization
26%

23%
Location-based Targeting

Fraud Analysis
22%

16%
Brand Sentiment Analysis

Product Placement Optimization
16%

9%
Other

Drug Discovery
1%



Big Data use cases across all industries

Financial Services

- Fraud detection
- Risk management
- 360° View of the Customer



Utilities

- Weather impact analysis on power generation
- Transmission monitoring
- Smart grid management

Transportation

- Weather and traffic impact on logistics and fuel consumption



IT

- Transition log analysis for multiple transactional systems
- Cybersecurity

Health & Life Sciences

- Epidemic early warning system
- ICU monitoring
- Remote healthcare monitoring



Retail

- 360° View of the Customer
- Click-stream analysis
- Real-time promotions

Telecommunications

- CDR processing
- Churn prediction
- Geomapping / marketing
- Network monitoring



Law Enforcement

- Real-time multimodal surveillance
- Situational awareness
- Cyber security detection



© 2012 IBM Corporation

A wide-angle photograph of a mountain range. The foreground shows a steep, rocky slope covered in green grass and small shrubs. In the middle ground, there are more mountain ridges and valleys, some with patches of snow or light-colored rock. The background features a series of jagged, snow-capped peaks that fade into a hazy, blue-tinted distance under a sky with soft, white clouds.

BIG DATA IN PRACTICE

TYPICAL TASKS

Distributed data storage

- Even geographically → Multiple data centers

Distributed data processing

- Collect
- Transform
- Query
- Analyze & understand

Distributed computing

PRINCIPLES 1.

Robustness & reliability on SW framework level

- Fault tolerance
- Redundant storage

“Keep everything”

- Including raw data

Linear (or better) scalability

- Horizontal (scale out) vs vertical (scale up)
- Scale down
- Dynamic / elastic / autoscaling

PRINCIPLES 2.

Efficiency

- High-throughput
- Low-latency

Data locality

- Execute computation where data are located → No unnecessary data transfers

Running on commodity HW

Dominated by open-source, community-driven SW (vs proprietary)

CHALLENGES 1.

Choosing the right tool

- Abundance of options 😊

Efficient data access

- Denormalization
- Graph schema
- Serialization

Testing

- Verification
- Debugging
- Performance measurement

CHALLENGES 2.

Enterprise integration

- Data hub / lake

Extremely large data size (exponential growth)

- Data federation / virtualization

Data governance

- Data sources, data integration / fusion, data catalogs, metadata management
- Data quality
- Security, privacy, legal compliance
- Retention policy

CHALLENGES 3.

High Availability (HA)

- No single point of failure (SPoF)
- Standby / fallback
- Replication / synchronization

Service Level Agreement (SLA)

- Availability
- Multi-tenancy
- Quotas
- Scheduler policy

CHALLENGES 4.

Administration / operation

- Installation, provisioning
- Monitoring
- Management
- Troubleshooting

Expenses

- Infrastructure
- Experienced workforce (e.g. Data Scientist, Data Engineer, Platform Engineer)
- Trainings, learning curve
- Commercial support / consultancy

An aerial photograph of a dramatic, winding mountain road. The road is paved and features multiple sharp turns, some with stone retaining walls. It snakes across steep, green hillsides. The sun is shining from the top center, creating a warm, golden glow and casting long shadows. A blue rectangular banner with white text is positioned horizontally across the middle of the image.

TECHNOLOGIES & TOOLS

BIG DATA OPEN-SOURCE LANDSCAPE

The Dataflog Open Source Landscape 2.0

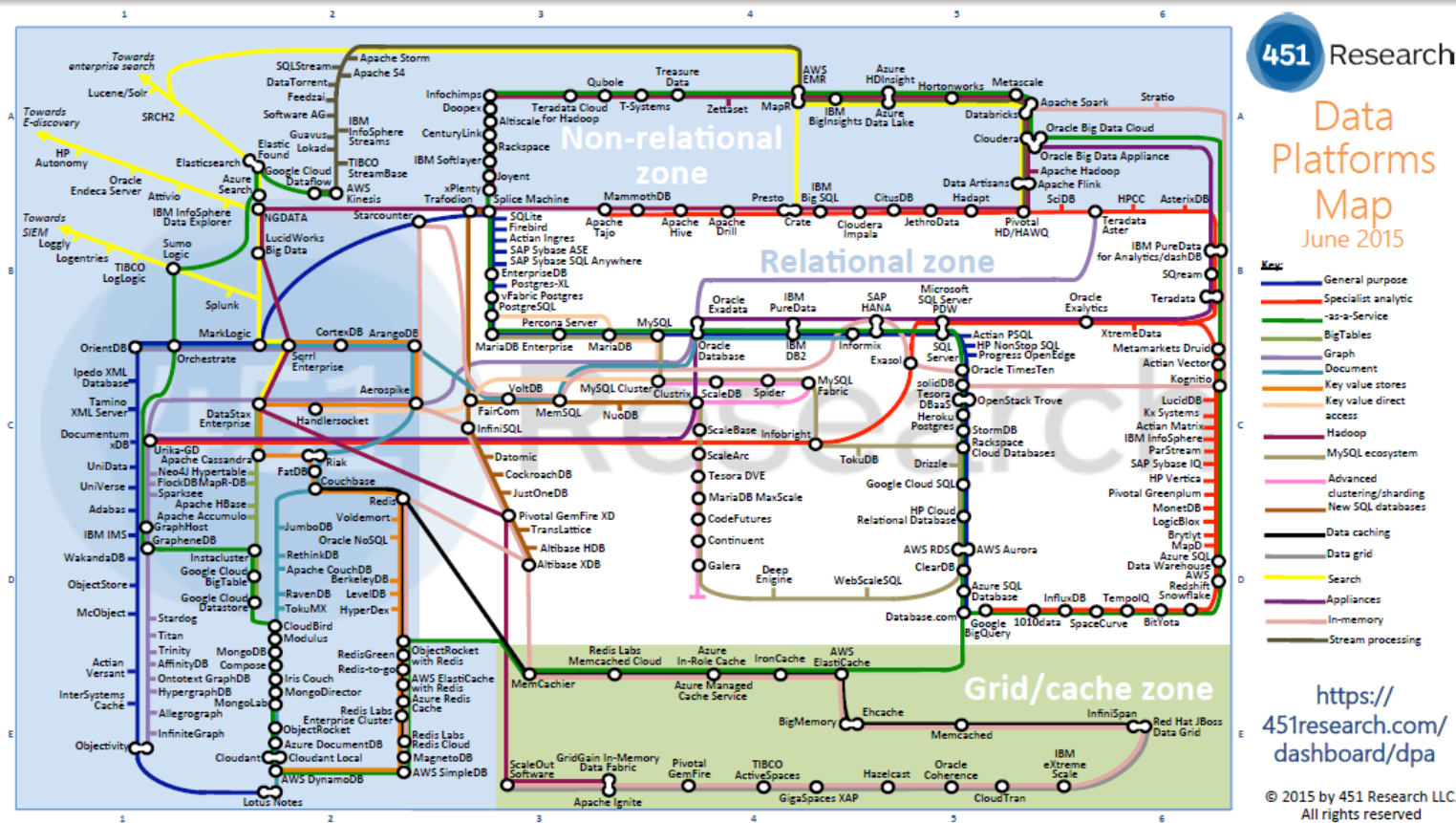


Created by: www.Dataflog.com

APACHE HADOOP AND ITS ECOSYSTEM



STORAGE: RDBMS, NEWSQL, NOSQL, GRID / CACHE

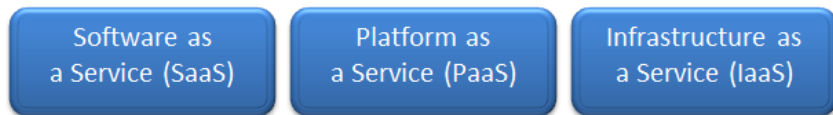


CLOUD

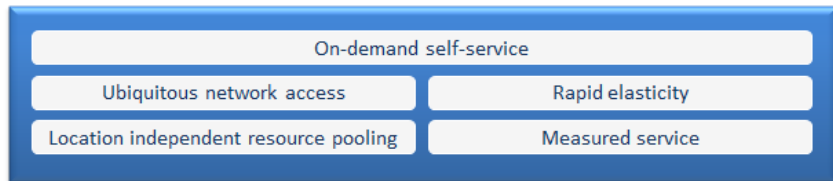
Deployment
Models



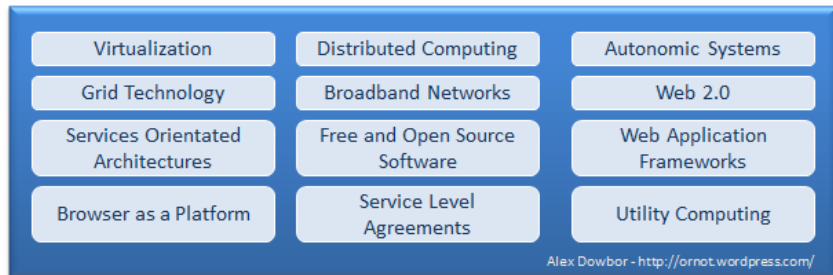
Delivery
Models



Essential
Characteristics



Foundational
Elements / Enablers



Based on the NIST Working Definition of Cloud Computing v14 and
<http://www.csrc.nist.gov/groups/SNS/cloud-computing/index.html>

Creative Commons Attribution-ShareAlike 3.0
Alexander Dowbor - <http://ornot.wordpress.com>



Amazon Elastic
MapReduce

APPLICATION DESIGN

Architecture

- Event-driven, reactive
- Lambda, Kappa
- Shared-nothing

Patterns

- MapReduce
- Actor model
- Data pipeline / flow

Algorithms

- Divide and conquer
- Concurrent / parallel

RELATED TOPICS: STORAGE

High-performance drives

- SSD
- RAID

Network storage

- SAN
- NAS

Network / distributed file systems

- NFS, Lustre, GlusterFS, GFS, HDFS, GPFS

“Fast data” (in-memory)

- Tachyon, GridGain / Apache Ignite file system

RELATED TOPICS: PROCESSING

High-performance networking

- InfiniBand, Fibre Channel, fiber-optics
- RDMA, zero-copy

Artificial intelligence

- Machine learning, NLP, data mining, dimension reduction

Analytics & statistics

- DWH, BI, data visualization

Data science

RELATED TOPICS: COMPUTING 1.

Parallel computing

- Multithreading, SMP, OpenMP
- GPGPU → OpenCL, CUDA
- SIMD, VLIW / MIMD, MPP, vector processors

Grid computing

- GigaSpaces XAP, GridGain / Apache Ignite, GemFire / Apache Geode, JPPF, HTCondor

HPC / supercomputers

- PVM, OpenMPI

RELATED TOPICS: COMPUTING 2.

Edge computing

- Sensor networks / IoT, P2P

“Fast data” (in-memory)

- Apache Spark, Apache Flink, SAP HANA



CONCLUSIONS

BIG DATA IS COMPLEX



POSSIBLE CONNECTIONS WITH TOMOGRAPHY

Storage

- Collect
- Query, retrieve
- Link with other data sources, associate metadata

Processing

- Transform, pre-process
- Analyze & understand
- Evaluate

Computing

- Reconstruct



THANK YOU!