# Hadoop and Spark

History, Concepts and offerings

14 setp 2017

Click to edit Master subtitle sty

**Christophe Menichetti**
Lead Architect (HPDA/AI) – Cognitive Systems
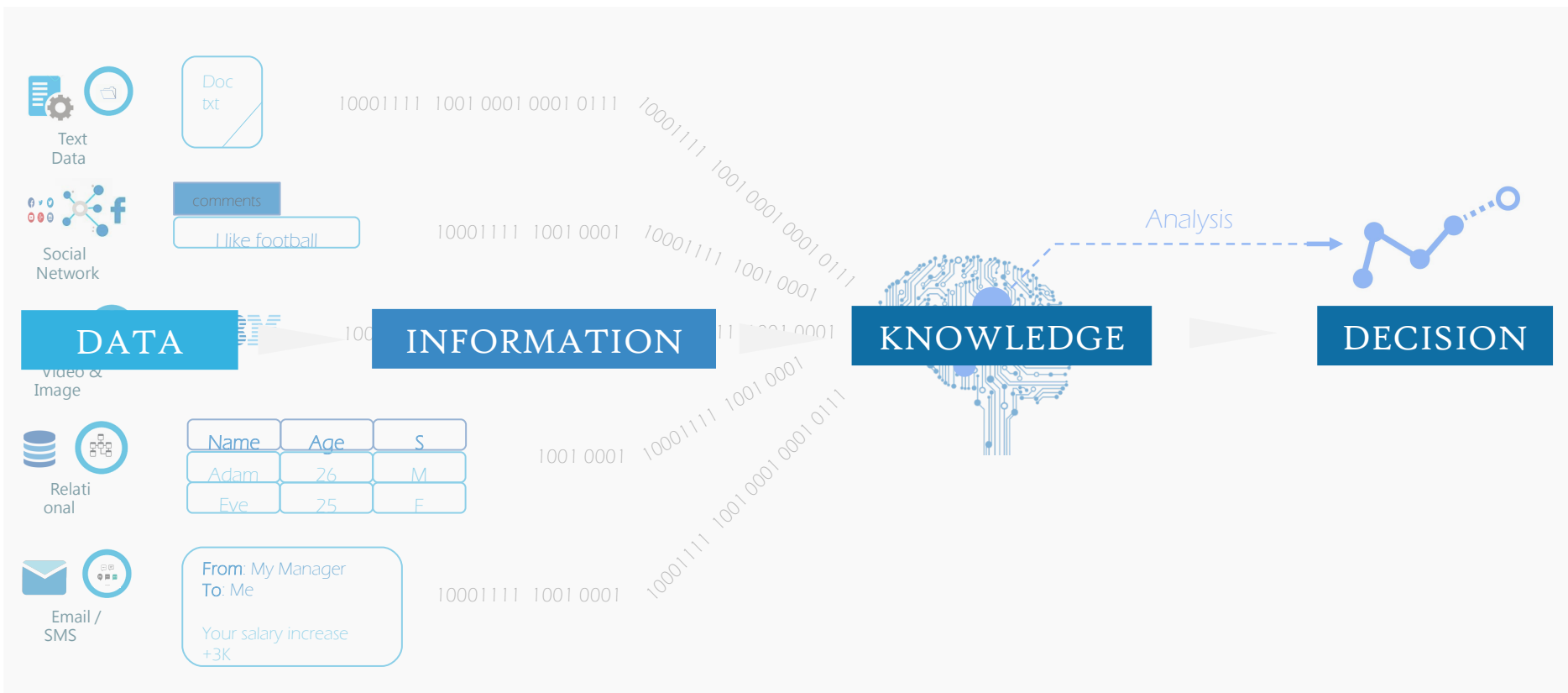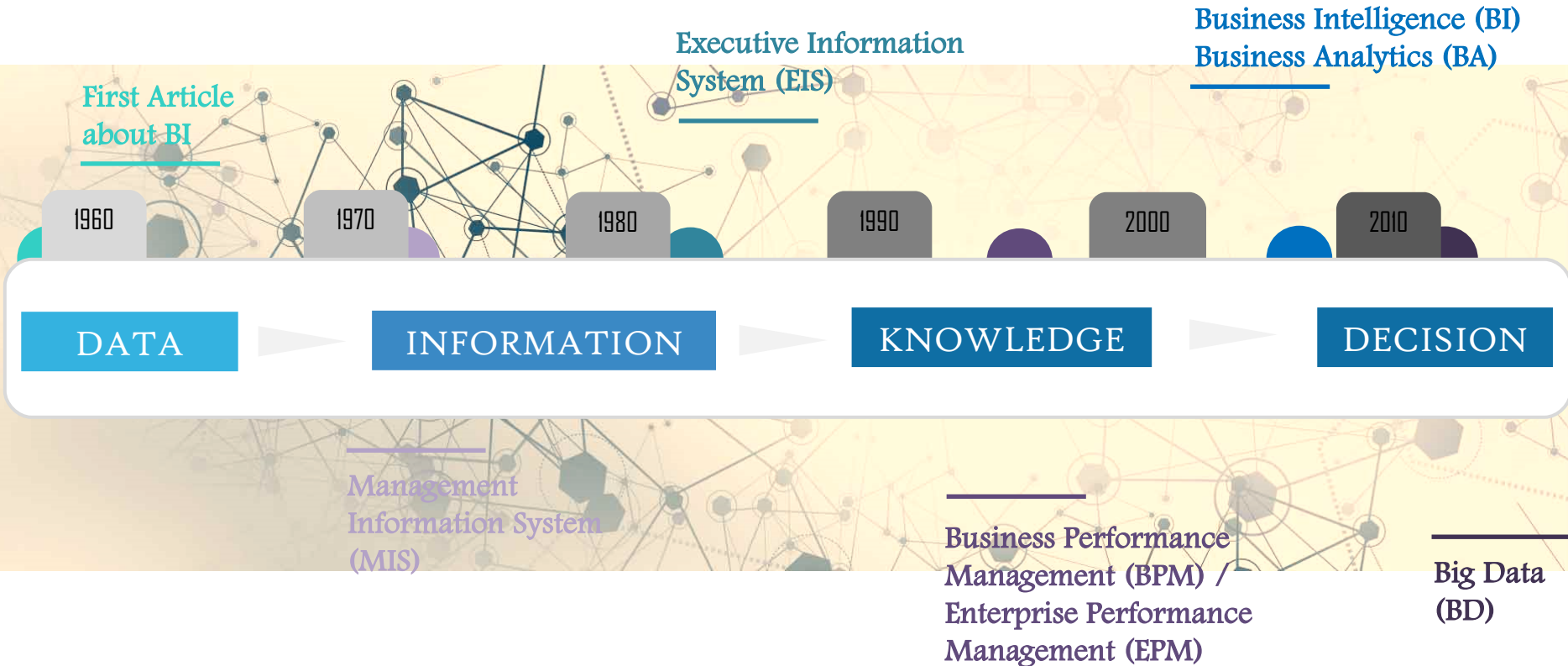
IBM Montpellier Client Center (IBMCCMPL)

# What is Analytics ?

**DATA**

Text Data

Social Network

Video & Image

Relational

Email / SMS

Doc txt

comments

I like football

| Name | Age | S |
|------|-----|---|
| Adam | 26 | M |
| Eve | 25 | F |

From: My Manager
To: Me

Your salary increase +3K

**INFORMATION**

10001111  1001 0001 0001 0111

10001111  1001 0001

1001 0001

10001111  1001 0001

**KNOWLEDGE**

Analysis

**DECISION**

First Article about BI

Executive Information System (EIS)

Business Intelligence (BI)
Business Analytics (BA)

1960  1970  1980  1990  2000  2010

| DATA | INFORMATION | KNOWLEDGE | DECISION |

Management Information System (MIS)

Business Performance Management (BPM) / Enterprise Performance Management (EPM)

Big Data (BD)

Analytics Appli & Tools

Data Ware house

ETL

| Advanced & Predictive Analysis | Corporate Performance Management | Planning, Budgeting & Forecasting |

By Product Line

By Region

By Customer

Predictive, Data Mining

Query & Reporting

Dash boards, Scorecards, Visualization

**?**

**Data change**

# Volume

# Variety

# Velocity

**Analog data**

**Digital data**

100%
90%
80%
70%
60%
50%
40%
30%
20%
10%
0%

99%

1986    2016

Transaction and Application data

Social media,

20% Structured data

80% Unstructured data

Machine data

Enterprise content

6 600+ pictures uploaded

80 000+ posts

In 60 sec

98 000+ tweets

168 Millions mails sent

# Value

# Veracity

CPU speed x10

RAM speed x9

Network speed x100

Disk speed x1.2

Management

How to improve performance over traditional BI solutions with disk bottleneck?

Query & Reporting

Dash boards, Scorecards, Visualization

By Product Line

By Region

By Customer

Reduce disk use

Parallelise disk use

In Memory Architecture

Hadoop/MPP Architecture

How to fix scalability and cost issues of traditional BI solutions?

NoSQL

Hadoop is an open-source software framework for storing data and running applications on clusters of commodity hardware. It provides massive storage for any kind of data, enormous processing power and the ability to handle virtually limitless concurrent tasks or jobs

Doug Cutting is one of the creators. His son's toy was a yellow elephant, becoming the icon of Hadoop

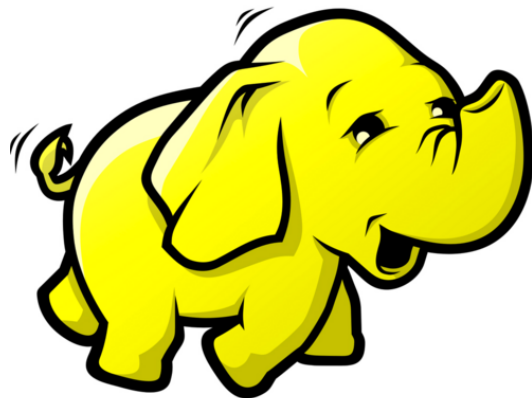The current version of Hadoop is 2.8 (Hadoop v3 is coming)

Today, Hadoop's framework and ecosystem of technologies are managed and maintained by the non-profit Apache Software Foundation (ASF), a global community of software developers and contributors

**Fault tolerance.** Multiple copies of all data are stored automatically. If one node goes down, jobs are automatically redirected to another one up

**Flexibility.** Unlike traditional relational databases, you don't have to pre-process data before storing it. You can store as much data as you want and decide how to use it later.

**Affordable.** The open-source framework is free and uses commodity hardware to store large quantities of data.
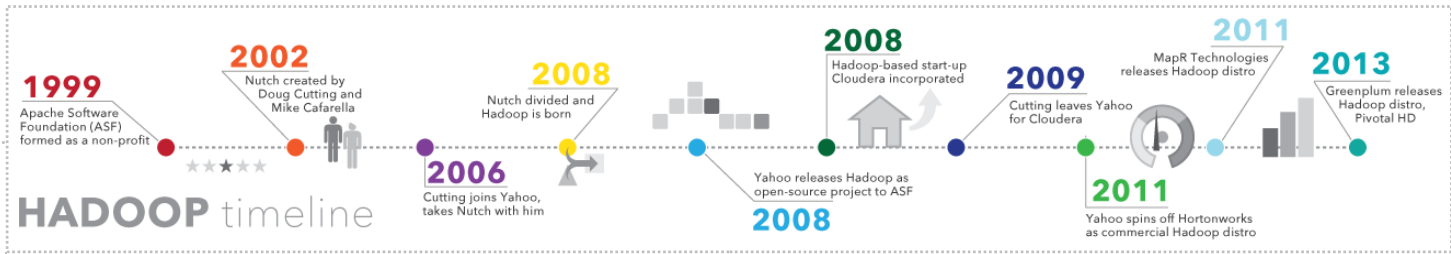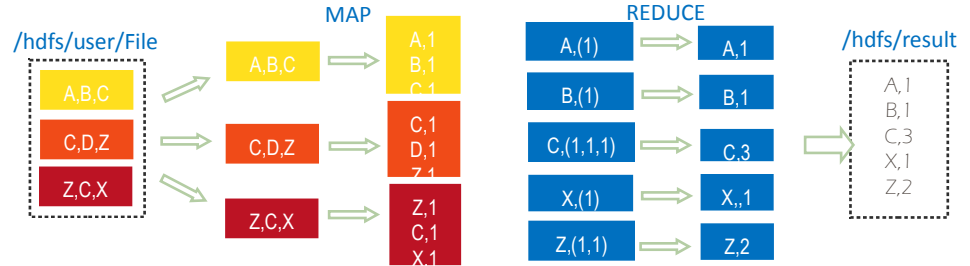
**Scalability.** You can easily grow your system to handle more data simply by adding nodes. Little administration is required.

## Map Reduce

From Yahoo Nutch web crawler

MapReduce is the JAVA framework for writing applications that process large amounts of structured and unstructured data stored in the HDFS cluster. We call that "function to data" instead of data to function
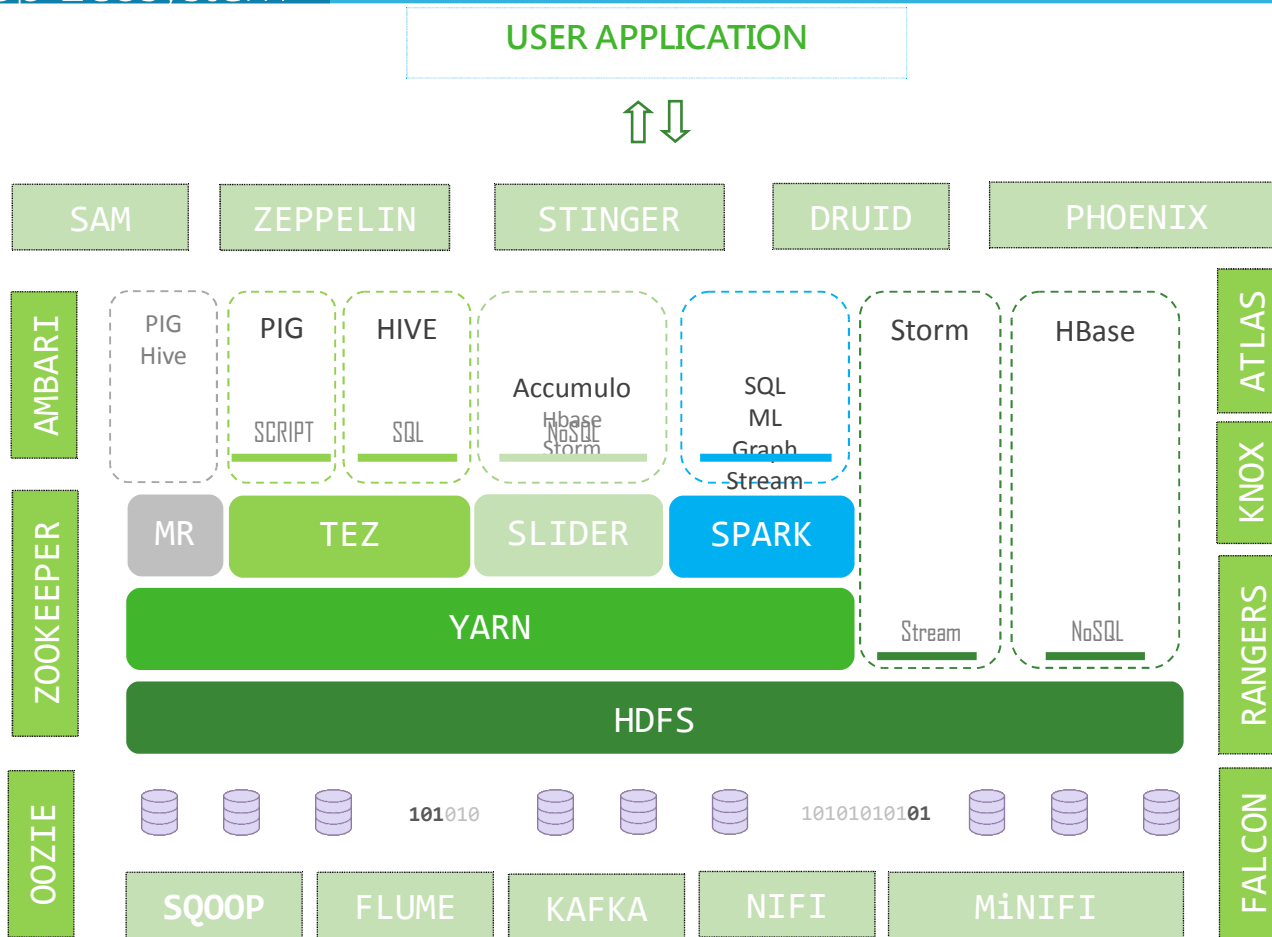
/hdfs/user/File

| A,B,C |
| C,D,Z |
| Z,C,X |

MAP

| A,B,C |
| C,D,Z |
| Z,C,X |

| A,1 B,1 C,1 |
| C,1 D,1 Z,1 |
| Z,1 C,1 X,1 |

REDUCE

| A,(1) | A,1 |
| B,(1) | B,1 |
| C,(1,1,1) | C,3 |
| X,(1) | X,,1 |
| Z,(1,1) | Z,2 |

/hdfs/result

A,1
B,1
C,3
X,1
Z,2

**HADOOP timeline**

**1999**
Apache Software Foundation (ASF) formed as a non-profit

**2002**
Nutch created by Doug Cutting and Mike Cafarella

**2006**
Cutting joins Yahoo, takes Nutch with him

**2008**
Nutch divided and Hadoop is born

**2008**
Yahoo releases Hadoop as open-source project to ASF

**2008**
Hadoop-based start-up Cloudera incorporated

**2009**
Cutting leaves Yahoo for Cloudera

**2011**
Yahoo spins off Hortonworks as commercial Hadoop distro

**2011**
MapR Technologies releases Hadoop distro

**2013**
Greenplum releases Hadoop distro, Pivotal HD

## HDFS

From Google File System project

Hadoop Distributed File System (HDFS) is a JAVA-based distributed file system that provides scalable, reliable (3 copies), high-throughput access to application data stored across servers cluster

/hdfs/user/File
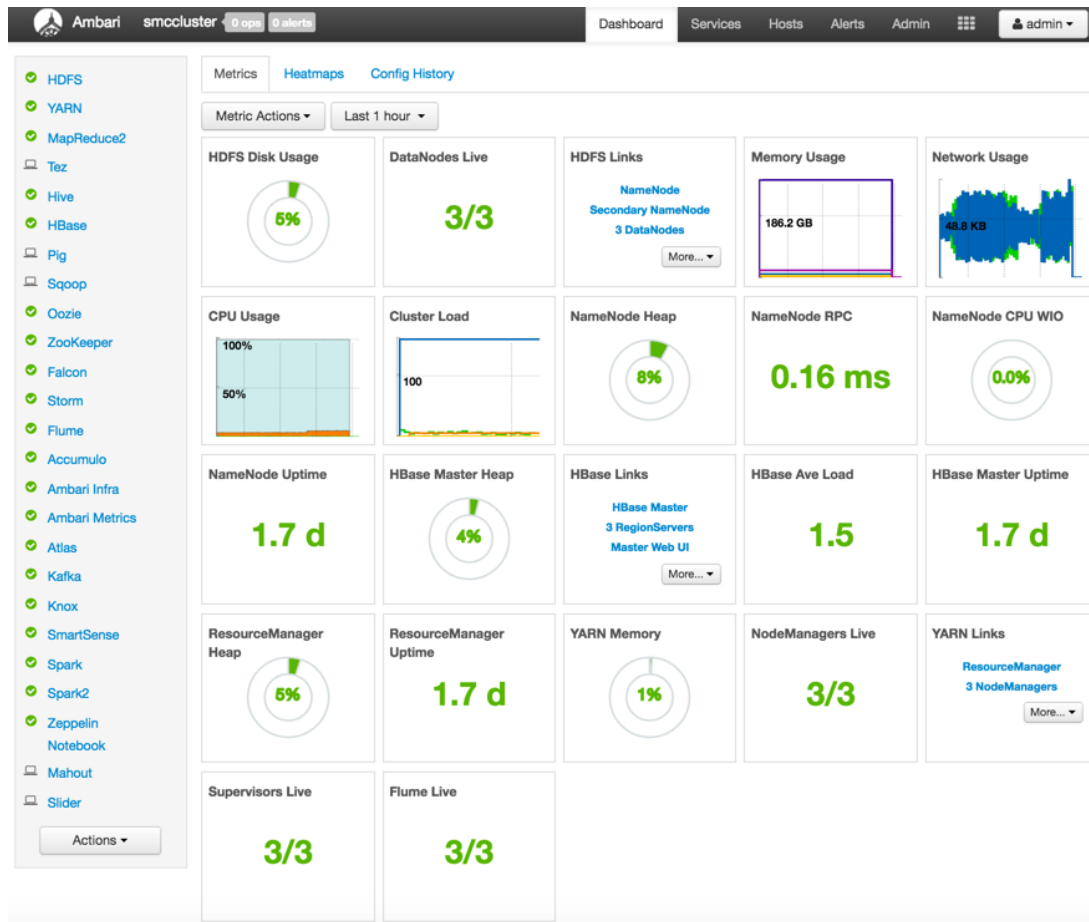
| line1 |
| line2 |
| line3 |

11

USER APPLICATION

SAM · ZEPPELIN · STINGER · DRUID · PHOENIX

AMBARI · ATLAS · KNOX · RANGERS · FALCON · ZOOKEEPER · OOZIE

PIG Hive

PIG · SCRIPT

HIVE · SQL

Accumulo · Hbase · NoSQL · Storm

SQL ML Graph Stream

Storm · Stream

HBase · NoSQL

MR · TEZ · SLIDER · SPARK

YARN

HDFS

101010 · 10101010101

SQOOP · FLUME · KAFKA · NIFI · MiNIFI

# Ambari

- **Dashboard**

  display KPI

- **Health Checks**

  display current
  service levels

- **Alerts**

  collect and send
  service metrics

PIG

SCRIPT

HIVE

SQL

HBase

NoSQL

Interactive

processing data with continual
exchange of information
between cluster and user

jobs always running unless
manually stopped, ingesting a
continuous data stream into
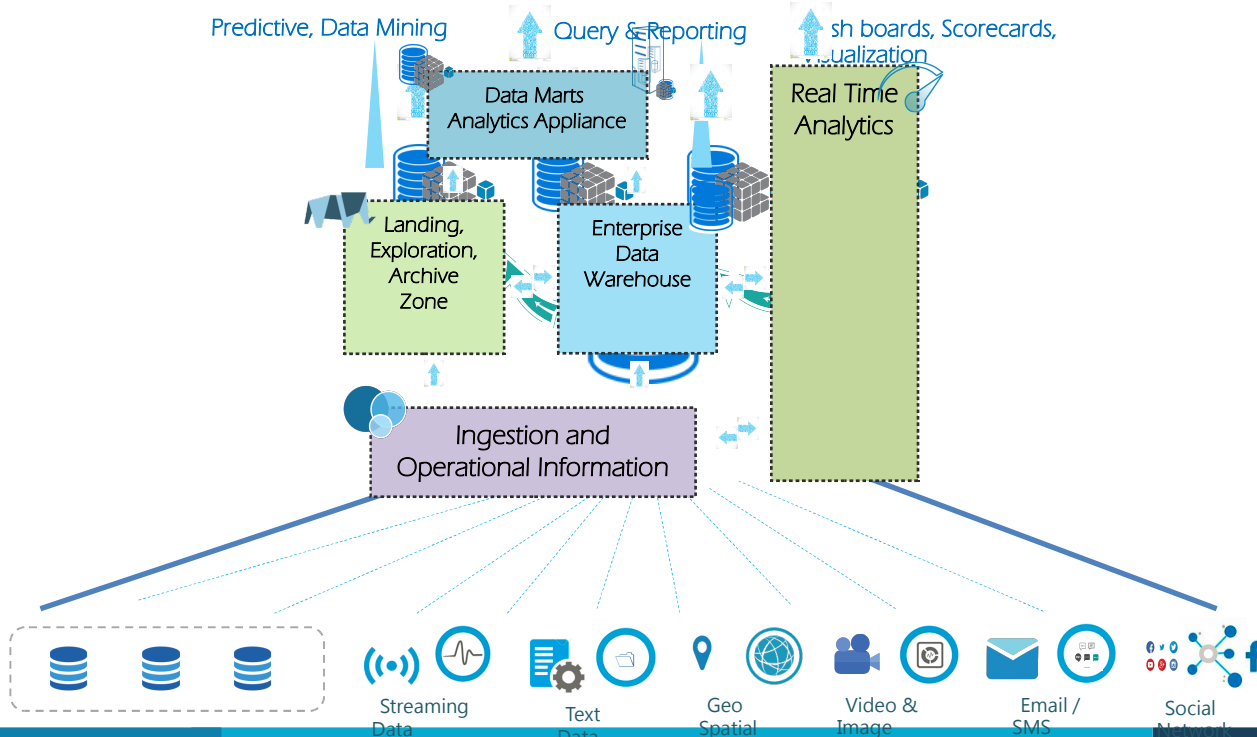mem, process it, and then
output it to storage.

Batch jobs are run
sporadically or periodically
from seconds to multiple
hours

Real Time

BATCH

Storm

Stream

SPARK

Multi

KAFKA

Stream

Map reduce

COMPUTE

# Big Data – high volume, high velocity, high variety - creates opportunities to extend Analytics for higher value

Zetta

10 Exa

100 PB

10 PB

1 PB

100 TB

**Untructured Data**
**Social Content Data**
**External Data**

?

**Structured Data**
**Operational Data**
**Internal Data**

*Multi-channel customer sentiment and experience a analysis*

*Detect life-threatening conditions at hospitals in time to intervene*

*Predict weather patterns to plan optimal wind turbine usage, and optimize capital expenditure on asset placement*

*Make risk decisions based on real-time transactional data*

*Identify criminals and threats from disparate video, audio, and data feeds*

# The Client Journey

BUSINESS VALUE

HIGHE

LOW

Commodity server myths

Most of them are HERE

New Business Models

**OPERATIONS**
- Low cost storage

**DATA WAREHOUSE**
- Data lake
- Data offload
- ETL offload
- queryable archive

**LoB / Analytics**
- 360 view customer
- Data Exploration
- Security
- Operations Analysis
- Social Media Analytics

**INNOVATION**
- Machine Learning
- Cognitive
- predictive analytics
- Recommandations
- Search Optimisation
- ......
- All data driven innovation for business

LOW                    HIGHE

BIG DATA MATURITY

1010101011111  10101010101001010101010101  IIIIIIIIIIII10101010101110

IIIIIIIII10101010101110

1010

1010

## Asterisk Call Detail Records

### List

Show 25 entries

| Source | Destination | Disposition | Answered | End | Duration |
|---|---|---|---|---|---|
| 646514562 | 100 | ← | 2011-06-24 08:13:58 | 2011-06-24 08:15:50 | 2min, 5s |
| 75382242 | 1358256400 | | 1970-01-01 01:00:00 | 2011-06-24 08:13:58 | 13s |
| 945590461 | 455450500 | ← | 2011-06-24 07:26:09 | 2011-06-24 07:26:46 | 1min, 37s |
| 1837728189 | 90723991 | ← | 06-24 07:14:18 | 2011-06-24 07:16:07 | 2min, 49s |
| 300 | 1769654229 | ← | 06-23 21:45:54 | 2011-06-23 21:48:25 | 3min, 47s |
| 1943684150 | 864446667 | | 01-01 01:00:00 | 2011-06-23 21:05:50 | 17s |
| 1380717900 | 100 | | 01-01 01:00:00 | 2011-06-23 21:05:50 | 17s |

| Phone Son | Phone Mum | | Monday 1st | 13:45 end | 3 seconds |
| | | | | 13:45 end | 2 seconds |
| | | | | 13:46 end | 3 seconds |

Call Details Records (CDR)

0101000111010101

Our sincere apologizes for the inconvenience, we offer you **30% discount** for the next month

Mum's birthday !

Call Mum

Son

Hi Mum, Happy ... **1**

Hi Mum, I wish ... **2**

Mum I said... **3**

Son ? It's you ?

What ? ? ?

I don't hear you

| CarID | Date | Speed | Location | Temp |
|-------|------|-------|----------|------|
| 1 | 10oct16 | 30 | 2 "34'15º23 | |
| 3 | 10oct16 | 34 | 2 "38'15º23.5 | |
| 4 | ............ | | | |

AUTO

Increase revenue

TEMPERATURE

GPS

PRESSURE

SPEED

Weather Forecast company

$$$$$$$$$$$$$$

Car rent company

$$$$$$$$$$$$$$

Gas Station company

$$$$$$$$$$$$$$

Traditional — Structured & Repeatable Analysis

Big Data — Iterative & Exploratory Analysis

ANALYTICS APPROACH

BIG DATA APPROACH

Data | Repository | Analysis | Insight

Small amount of carefully considered information

Video & Image
Streaming Data
Social Network
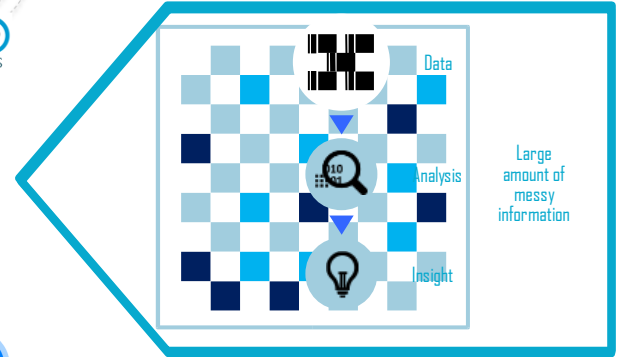Email / SMS

Data
Analysis
Insight

Large amount of messy information

Analyze data *after* it's been processed and landed in a warehouse or mart

Carefully cleanse information *before* any analysis

Analyze data *in motion* as it's generated, in real-time

Analyze information as is, cleanse as needed

Centralized Architecture

SCHEMA Architecture

Distributed Architecture

SCHEMA Architecture

**VARIETY**

VOLUME

UNSTRUCTU RED

## Distributed No SQL

Landing/Exploration zone

Schema less (no SQL), scalable, cost competitive

_____

Cloudera MapR Cloudant
Hortonworks Cassandra Redis
BigInsights MongoDB

## Streaming

Real Time zone

Micro seconds latency, designed for millions of events per sec

_____

EsperTech IBM Streams
Oracle Complex Event Processing

Big Data

SPARK

STRUCTUR ED

## Standard SQL

Data Warehouse zone

Traditional, mature, consistent and structured technology

_____

IBM DB2 Oracle
MySQL SQL Server

## In Memory

Datamart / operation analytics
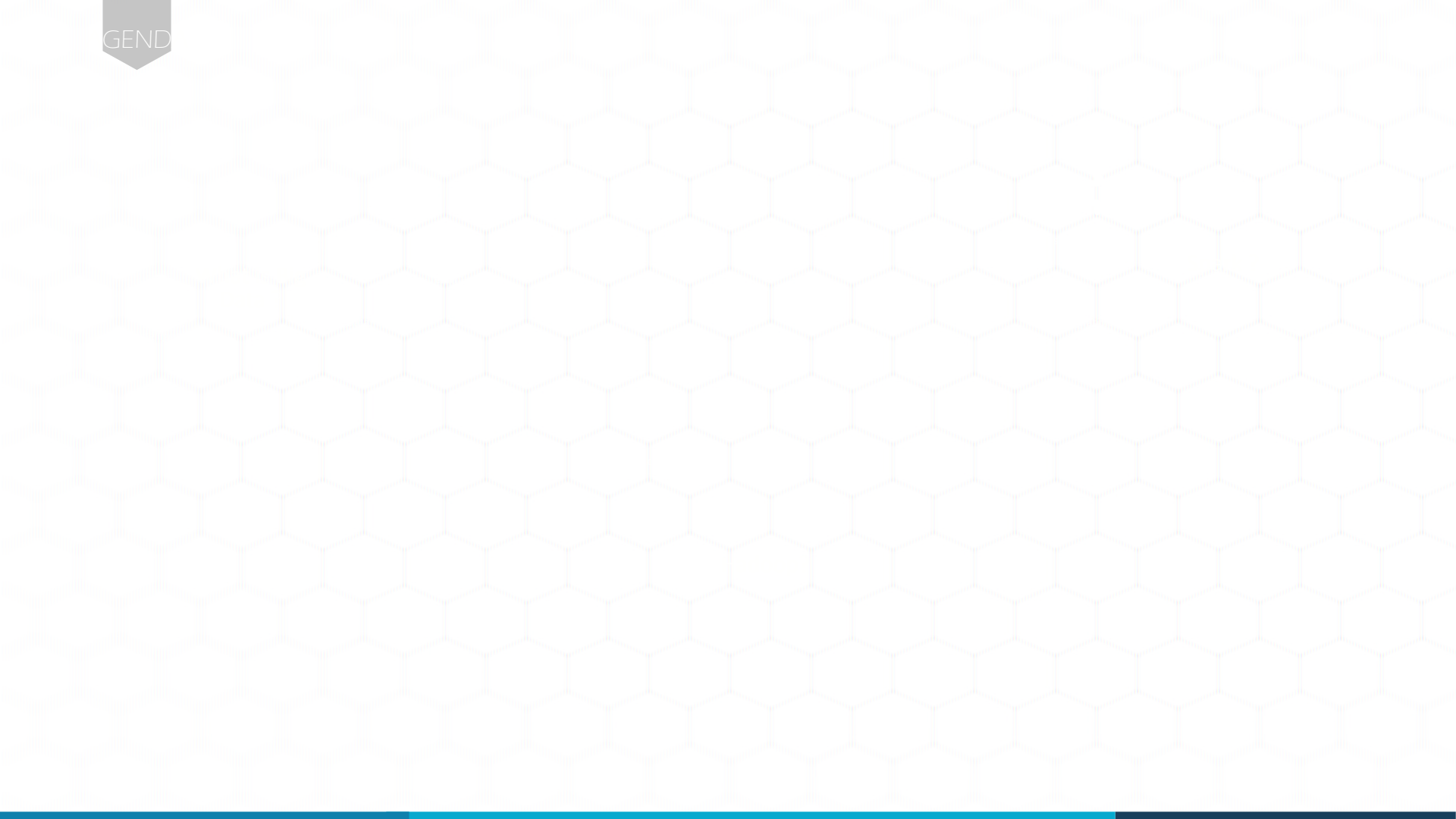
Extreme SQL performance, optimized data store

_____

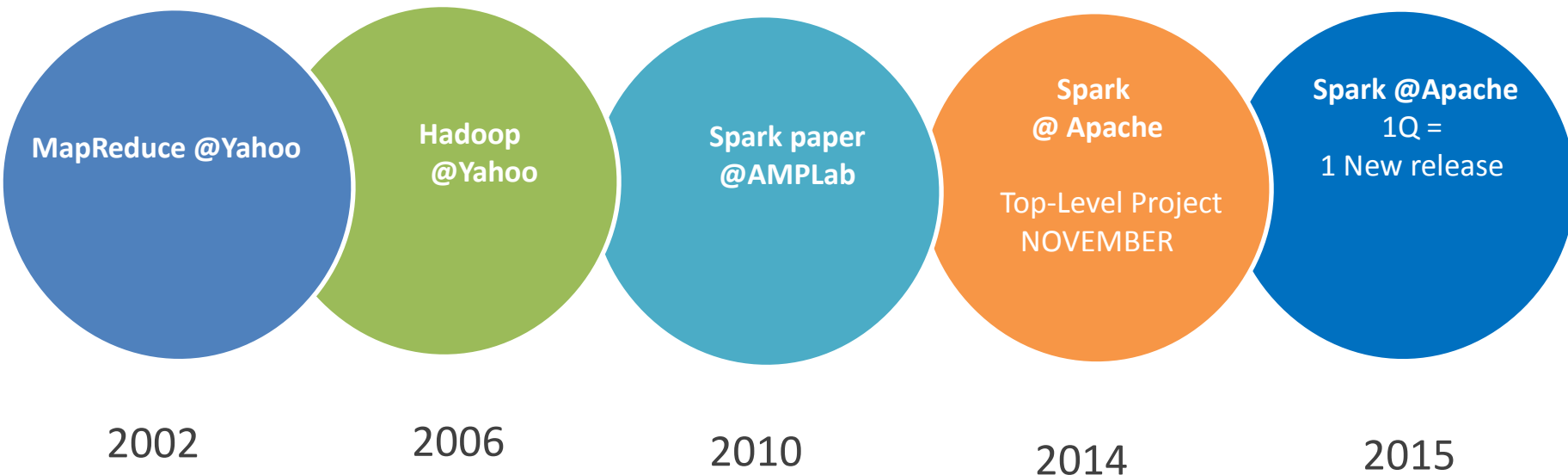SAP HANA, ORACLE Exalytics IBM DB2 memsql

Analytics

VELOCITY

BATCH

REAL TIME

**MapReduce @Yahoo**

**Hadoop @Yahoo**

**Spark paper @AMPLab**

**Spark @ Apache**

Top-Level Project
NOVEMBER

**Spark @Apache**
**1Q =**
**1 New release**

2002

2006

2010

2014

2015

**Over 750 contributors from over 200 companies including IBM, Google, Amazon, SAP …**

Hadoop is an unlimited scale, extremely economic platform to "batch process" a wide range of data, especially unstructured data (80% of data). But Hadoop is

- **NOT Ease of development** (need of deep java expertise, few abstractions)
- **NOT Performant for interactive process** (slow disk write, suitable for bulk batch processing)
- **ONLY** suitable **for batch workloads**, Rigid processing model

## Performant



- **In-memory architecture** greatly reduces disk I/O, Anywhere from 20-100x faster for common tasks
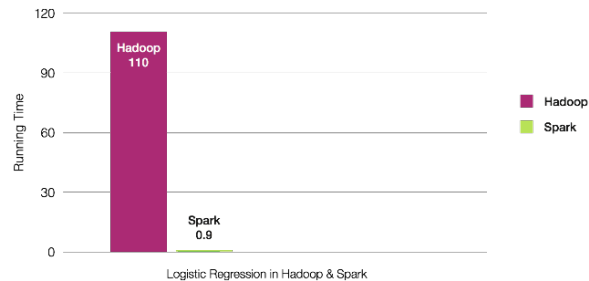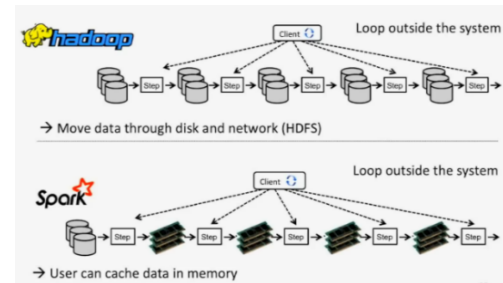
## Productive

- Concise and expressive syntax, **Single programming model** across a range of use cases
- Integrated with common programming languages – **Java, Python, Scala** / New tools continually reduce skill barrier for access (e.g. SQL for analysts)

## Protect existing investments

- Works well within **existing Hadoop ecosystem / Large and growing community**





Logistic Regression in Hadoop & Spark

# Apache Spark: a Analytics Framework

| | |
|---|---|
| Java / Python / Scala / R | Languages |

| | |
|---|---|
| **Spark SQL** Relational Operators | **Spark MLlib** Machine Learning | **Spark GraphX** Graph Processing | **Spark Streaming** Real-Time Streaming | Spark Libraries |

| | |
|---|---|
| Spark Core General Execution Engine | Spark Core |

| | |
|---|---|
| YARN | MESOS | Standalone | Cluster Manager |

| | |
|---|---|
| DB2 / HDFS / Cassandra / HBase / Oracle / JSON / Parquet / VSAM ... | Data Abstraction |

**Faster development**

**Easy of Use**

**Hight-level API**

**General purpose**

**Scheduling**

**Scalability**

**Fault tolerance**

## What Spark *isn't*

- **A data store** – Spark attaches to other data stores but does not provide its own

- ***Only* for Hadoop** – Spark can work with Hadoop (especially HDFS), but Spark is a separate, standalone system

- ***Only* for machine learning** – Spark includes machine learning and does it very well, but it can handle much broader tasks equally well

- **A replacement for Streams** – Spark Streaming is micro-batching, not true streaming, and cannot handle the real-time complex event processing that true streams do

# Common Spark use cases

**1** Interactive querying of very large data sets (e.g. BI)

**2** Running large data processing batch jobs (e.g. nightly ETL from production systems, primary Hadoop use case)

**3** Complex analytics and data mining across various types of data

**4** Building and deploying rich analytics models (e.g. risk metrics)

**5** Implementing near-realtime stream event processing (e.g. fraud / security detection)
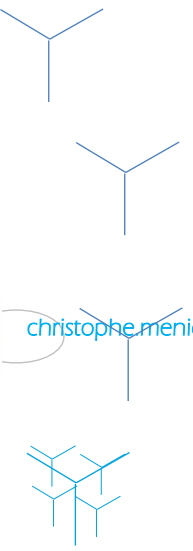
# Spark vs. Hadoop

## How is Spark SIMILAR to Hadoop?

- Similar divide-and-conquer architecture of breaking large jobs into smaller pieces

- General data processing platform suitable for batch analysis

- Can coexist within existing Hadoop environments and use Hadoop components such as HDFS

- Open source with extensive community support

## How is Spark DIFFERENT from Hadoop?

- In-memory architecture vs. file-based for Hadoop, generates up to 100x speed improvements

- Faster speed enables new use cases such as interactive or iterative analysis

- Simpler programming model, up to 5x less code

- Multiple programming languages supported, vs. only Java for Hadoop

- Single modular platform enables extension via libraries, not separate applications

- Specialized machine learning algorithms available

THANK YOU FOR YOUR ATTENTION

christophe.menichetti@fr.ibm.com

January@2017

More Details

LEARNING MORE

www.bigdatauniversity.com
www.ibm.com/bigdata

BRIEFING / CONFERENCES

Any Projects

WORKSHOPS

Big Data Solutions Sizing
Big Data Design Thinking

TEST / BENCHMARK

Related Subjects

COGNITIVE

HPC / GPU and FPGA acceleration