

# BM 과정 AI Architecture & Business Model

2019

에스튜비즈  
대표컨설턴트  
이 이 백



## 목 차

- I. Big Data Architecture와 AI 시스템의 Position
- II. AI 구축 방법론
- III. Industry 4.0
- IV. AI 시스템 구축을 위한 기본 지식과 기술
- V. Deep Learning Revisit
- VI. Business Model과 AI Topic 발굴

## I

## Big Data Architecture와 AI 시스템의 Position

1. Big Data Concept Making
2. Big Data Architecture
3. Machine Learning & Deep Learning
4. Big Data 시스템에서의 AI의 위치와 역할



## 1. Big Data 기반 Analytics Concept Making

- 분석 패러다임의 변화

### 1) Analysis to Analytics

		직관/경험 중심	Analytics 중심
문제 해결 범위	특정 상황의 단위별 과제 해결	복합 과제 해결을 통해 일반화	
Input	작업자 선정한 소수의 Parameter	모든 Parameter와 파생 변수	
Output	Rule(If-Then)	<ul style="list-style-type: none"><li>▪ Label = Features(X1, X2, .....Xn)</li><li>▪ Rule(If-Then)</li></ul>	
지식 유형	유효 기간이 있는 암묵지	연속된 자기 학습으로 개선하는 형식지	
개선의 양상	연속적 점진적 개선	불연속 급진적 개선	

# I. Big Data Architecture와 AI 시스템의 Position

## 1. Big Data 기반 Analytics Concept Making

- 분석 패러다임의 변화

### 1) Analysis to Analytics



## 1. Big Data 기반 Analytics Concept Making

- 분석 패러다임의 변화

### 2) Technology Shift : Data Warehouse to Data Lake

	전통적인 Data Warehouse	Data Lake
Technology	▪ DBMS	▪ Big Data Ecosystem
Data Type	▪ Structured Data ▪ Semi-Structured의 제한적 수용	▪ Structured ▪ Semi-Structured ▪ Unstructured / Free Format
Data Capability	▪ 상업용 고가의 Software & Repository ▪ 확장 제약(inflexibility)	▪ Open Source Software & 저비용 Repository ▪ Scale-Out
Data Access	▪ DBMS 접근은 속련자 Only ▪ 일반 사용자는 제공하는 Report 내에서 한정된 정보 접근 ▪ 분석가는 개인PC는 Mart를 통해 사전에 정의된 한정된 데이터 접근	▪ 비숙련자의 정보 접근 확대 수용 ▪ 일반 사용자는 다양한 시각화 도구로 확대된 정보 접근 ▪ 분석가를 위한 개인화 대용량 공간 제공

## 1. Big Data 기반 Analytics Concept Making

- 분석 패러다임의 변화

### 2) Technology Shift : Data Warehouse to Data Lake

#### Unstructured data

The university has 5600 students.  
John's ID is number 1, he is 18 years old and already holds a B.Sc. degree.  
David's ID is number 2, he is 31 years old and holds a Ph.D. degree. Robert's ID is number 3, he is 51 years old and also holds the same degree as David, a Ph.D. degree.

#### Semi-structured data

```
<University>
  <Student ID="1">
    <Name>John</Name>
    <Age>18</Age>
    <Degree>B.Sc.</Degree>
  </Student>
  <Student ID="2">
    <Name>David</Name>
    <Age>31</Age>
    <Degree>Ph.D. </Degree>
  </Student>
  ...
</University>
```

#### Structured data

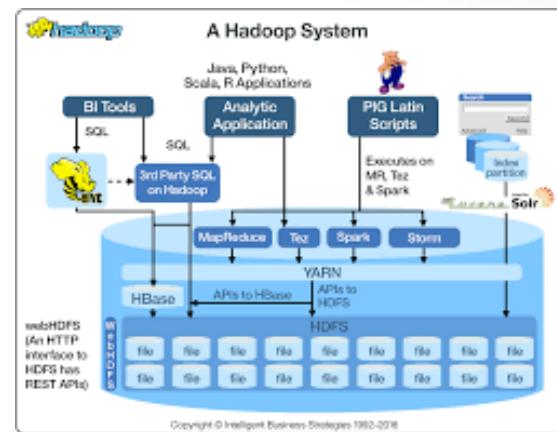
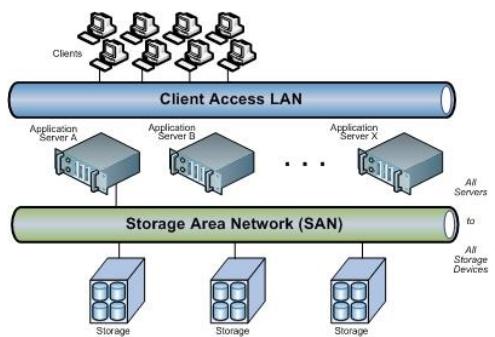
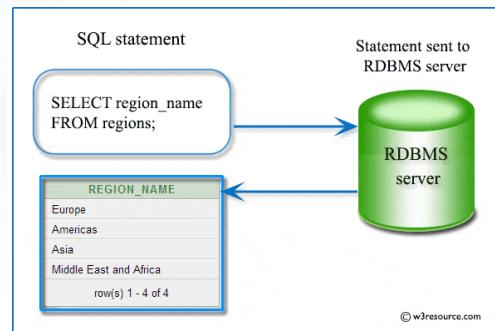
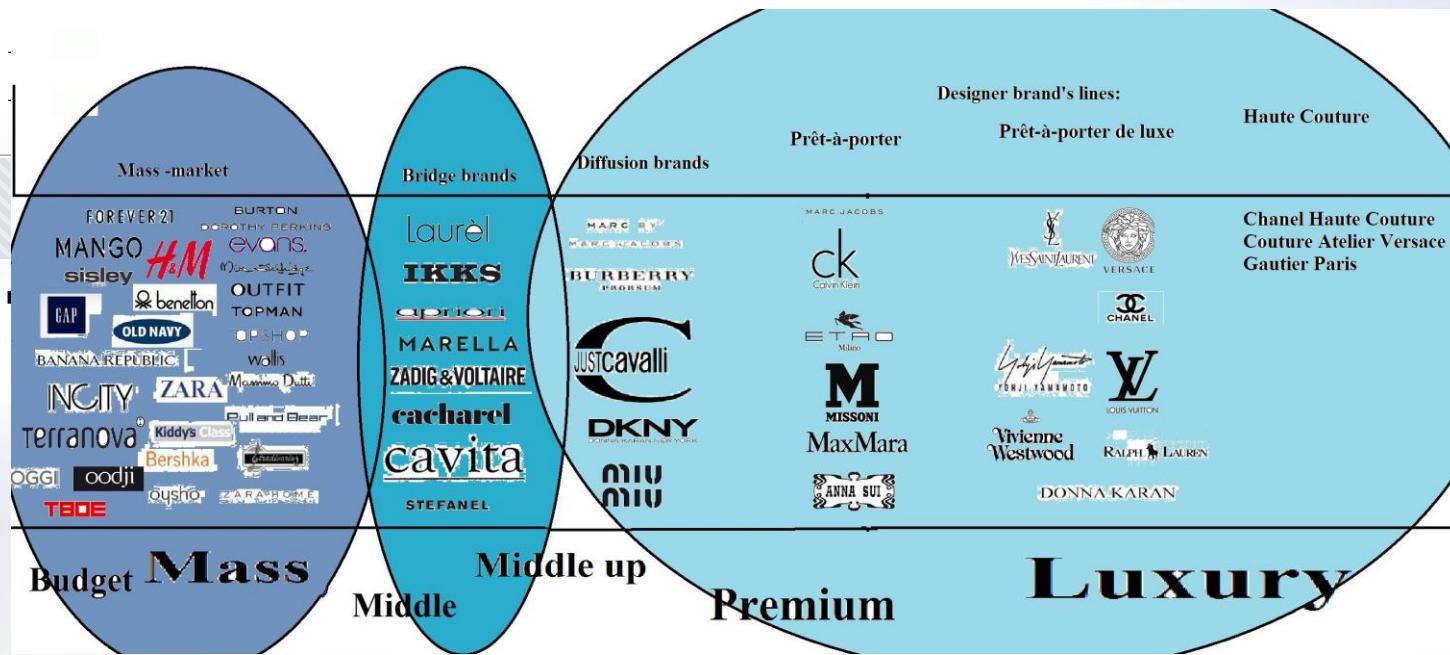
ID	Name	Age	Degree
1	John	18	B.Sc.
2	David	31	Ph.D.
3	Robert	51	Ph.D.
4	Rick	26	M.Sc.
5	Michael	19	B.Sc.

## 1. Big Data 기반 Analytics Concept Making

- 분석 패러다임의 변화

### 3) Small Data to Big Data

	Small Data	Big Data
Capability	개인 PC, Workstation, Single Server	<b>Big Data Infrastructure</b>
Tools	Excel, R, Python	Python, Scala, R, Scala, Java
Data	<ul style="list-style-type: none"><li>▪ Software dependent File<ul style="list-style-type: none"><li>• excel</li></ul></li><li>▪ Flat File<ul style="list-style-type: none"><li>• txt, csv</li></ul></li><li>▪ DBMS<ul style="list-style-type: none"><li>• oracle, mysql, etc</li><li>• sql, stored procedure, api for dbma</li></ul></li></ul>	<ul style="list-style-type: none"><li>▪ Software Independent</li><li>▪ Any File<ul style="list-style-type: none"><li>• txt, csv</li><li>• json, parquet</li><li>• hdfs on Hadoop</li></ul></li><li>▪ Any DBMS<ul style="list-style-type: none"><li>• sql, orc</li><li>• hive, impala</li><li>• real time</li></ul></li><li>▪ NoSQL<ul style="list-style-type: none"><li>• hbase, mongodb, kasandra</li></ul></li></ul>



## 1. Big Data 기반 Analytics Concept Making

- 분석 패러다임의 변화

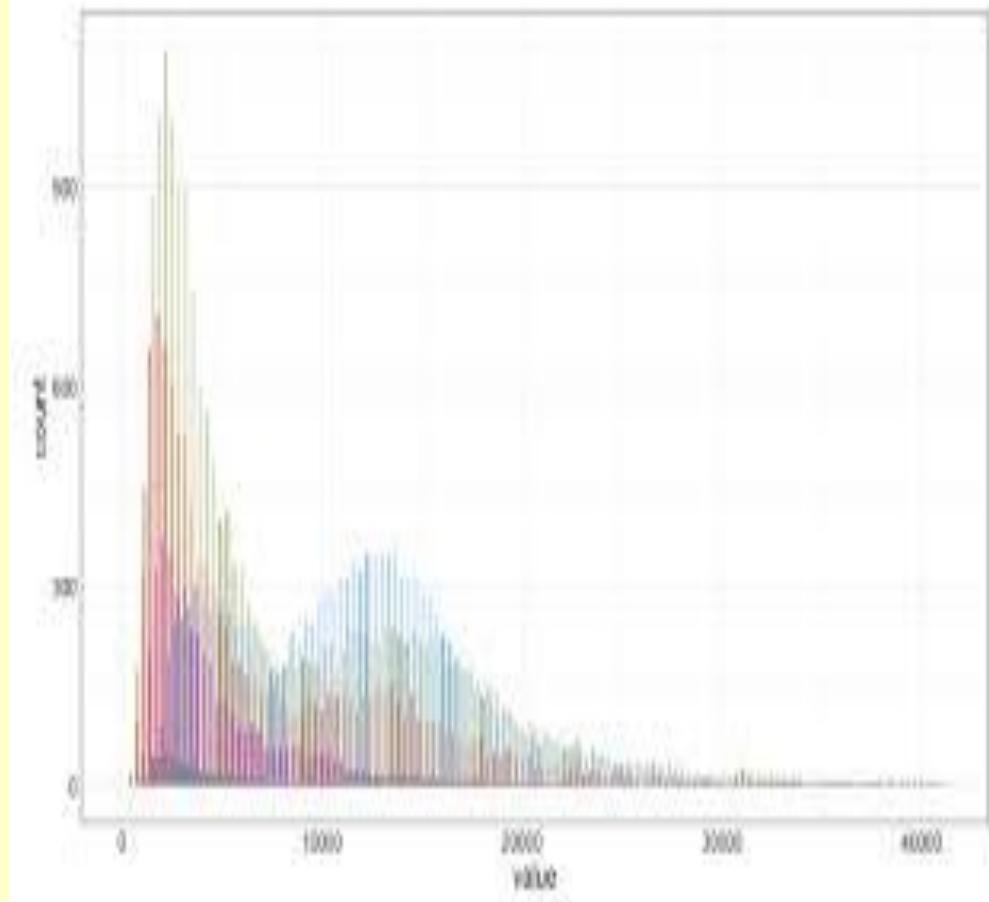
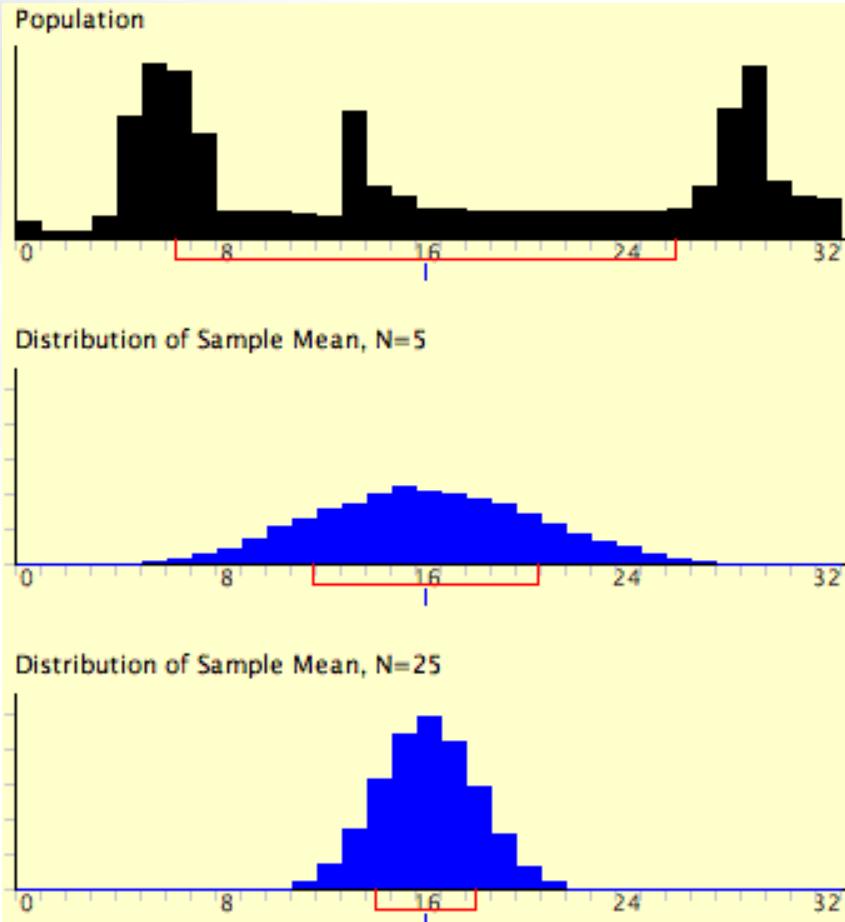
### 4) Sampling Methodology to Big Data Methodology

	Sampling Methodology	Big Data Methodology
일반화 접근	<ul style="list-style-type: none"><li>▪ Parametric Approach</li><li>▪ Small Sample의 반복 적용</li><li>▪ Normal Distribution</li><li>▪ Linearity(선형성) 중족</li></ul>	<ul style="list-style-type: none"><li>▪ Non-Parametric의 수용</li><li>▪ Big Data의 반복 적용</li><li>▪ Combination between Linearity and Non-Linearity</li></ul>
Main Stream	<ul style="list-style-type: none"><li>▪ Statistics</li><li>▪ Simple Gaussian<ul style="list-style-type: none"><li>• Central Tendency</li><li>• Linear Transformation</li><li>• Correlation</li><li>• Regression(Logistic 포함)</li><li>• PCA</li></ul></li><li>▪ Data Mining<ul style="list-style-type: none"><li>• Knowledge Data Discovery</li><li>• Decision Tree, Clustering</li></ul></li></ul>	<ul style="list-style-type: none"><li>▪ Machine Learning &amp; Deep Learning</li><li>▪ Multiple Gaussian Fitting<ul style="list-style-type: none"><li>• (looping) Linear Transformation + Non-Linear Transformation</li></ul></li><li>▪ Perceptron<ul style="list-style-type: none"><li>• Label &amp; Features</li><li>• Multi-Layers</li><li>• Activations</li></ul></li></ul>

## 1. Big Data 기반 Analytics Concept Making

- 분석 패러다임의 변화

### 4) Sampling Methodology to Big Data Methodology



## 1. Big Data 기반 Analytics Concept Making

- 분석 패러다임의 변화

### 4) Sampling Methodology to Big Data Methodology

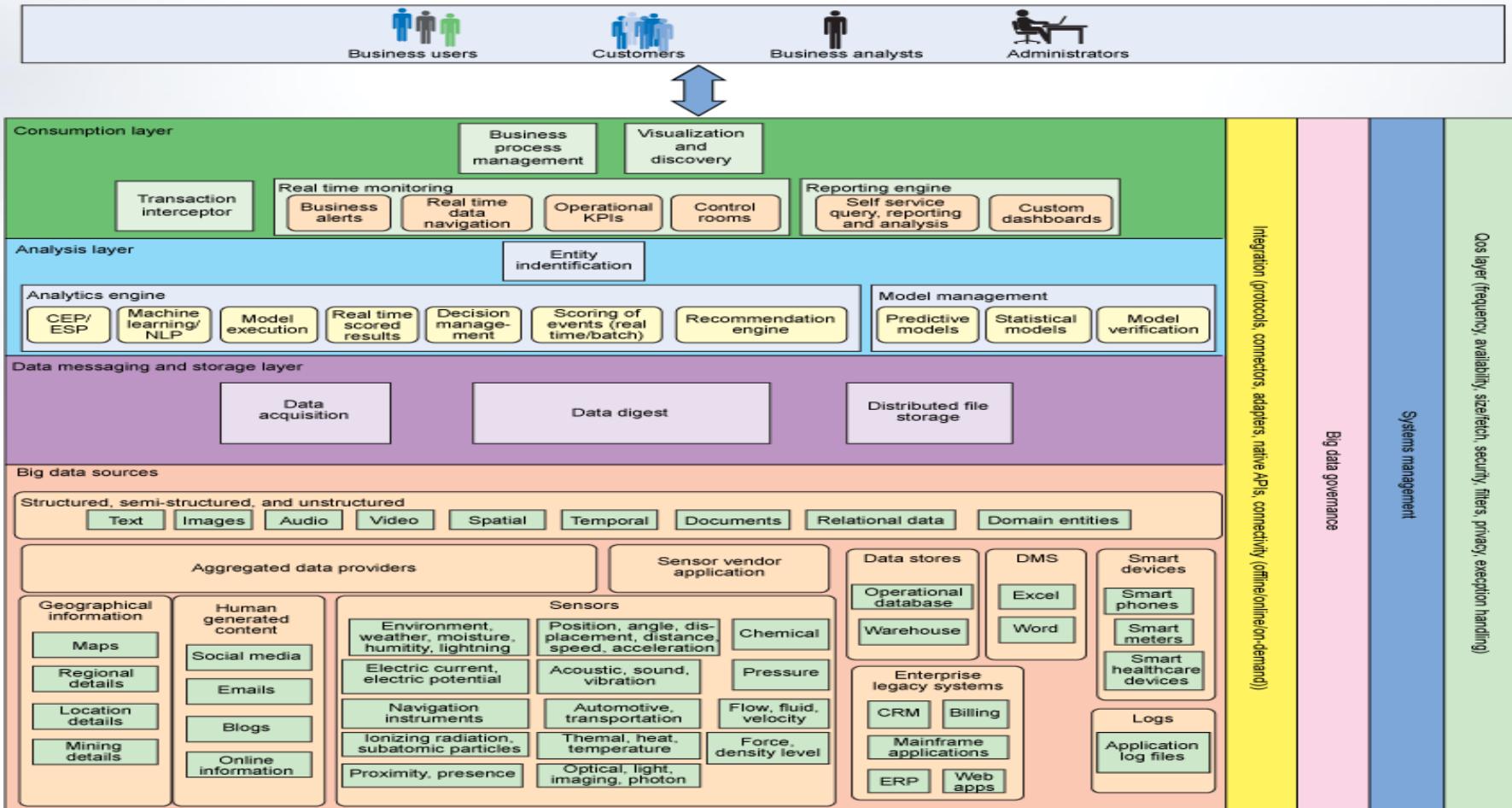


# I. Big Data Architecture 와 AI 시스템의 Position

## 1. Big Data 기반 Analytics Concept Making

### ▪ Logical Layers of Big Data Solution

- Big Data Sources, Data Management and Store Layer, Analysis Layers, Consumption Layer
- [Point] 아래의 표를 이용하여 가치 창출을 요하는 Big Data 기반 Business Model의 Layer 정의

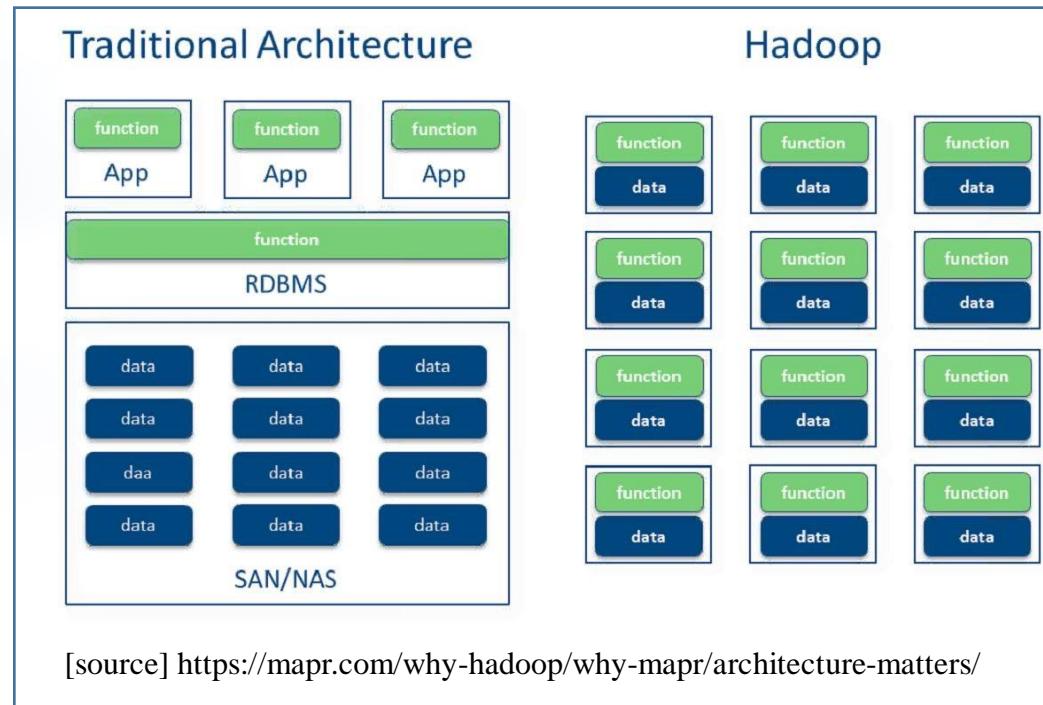


## 2. Big Data Architecture

- Hadoop System의 특징

- Data Node의 **유연한 확장** – Tera, Peta
- Map Reduce : Key-Value Pair

- Scale-Up
- Row/Column
- SQL
  - Conditional Record Selection
- OLTP
- Structured Data

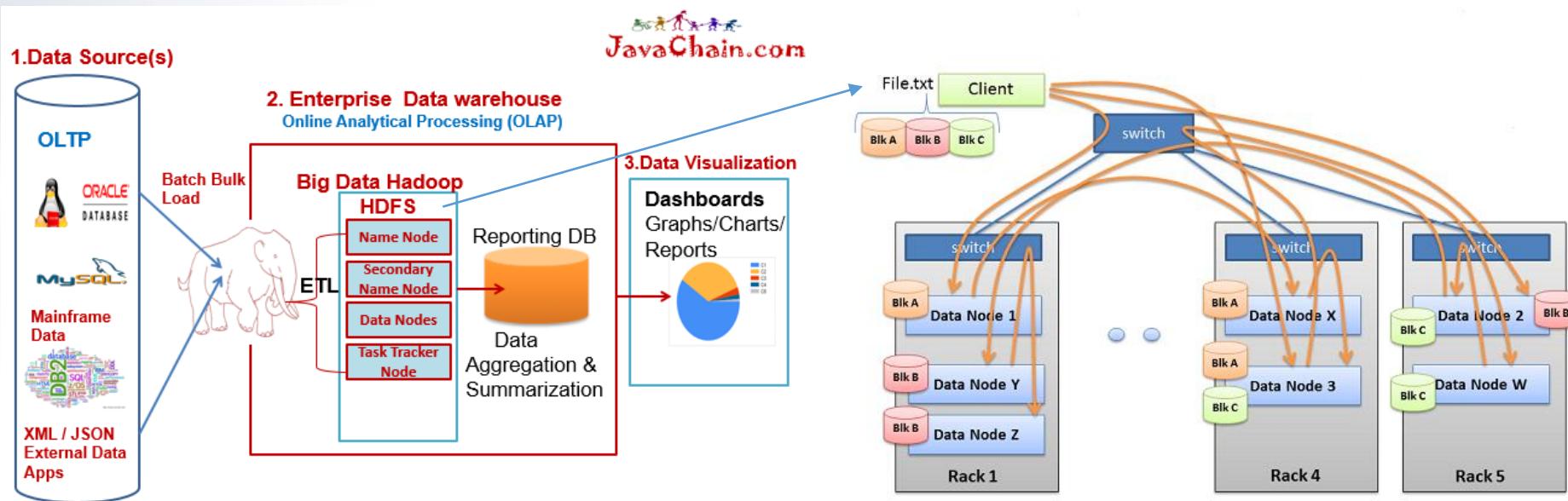


- Scale-Out
  - Tera → Peta
- Key-Value Pair
- Map Reduce
  - Mapper, Reducer, Driver
- Batch
- Any Types
  - Structured
  - Semi-Structured
  - Unstructured

- 분석하고자 하는 대상 데이터의 집합소 → 전수 데이터 기반 분석 접근
- CSV, Text, Image, Voice, Video 등 모든 유형의 데이터 접근
- ☞ Real World Data의 제한없는 확장으로 실질적 분석 작업 인프라 실현

## 2. Big Data Architecture

- HDFS : Hadoop Distributed File System
  - CSV, TXT 외의 다양한 Flat File Format(XML, JSON, Parquet 등)의 Hadoop에 저장
  - 아무리 큰 데이터라도 단위블록으로 분할하여 Data Node에 분산 저장
  - 데이터의 완전성 보증 : 하나의 단위 블럭을 3개의 Data Node에 저장



[source] <http://www.javachain.com/big-data-hadoop-in-data-warehouse/>

[source] <https://hadoopabcd.files.wordpress.com/2015/03/multi-block-replication-pipeline1.png/>

AI 관점

- Analytic Input Data의 누락없음의 보증
- CSV 중심으로 Input Data를 구성하더라도 XML, JSON, Parquet Format의 추가 가능
  - XML, JSON, Parquet 등의 데이터 접근이 가능한 Library 설치 및 Script 작성

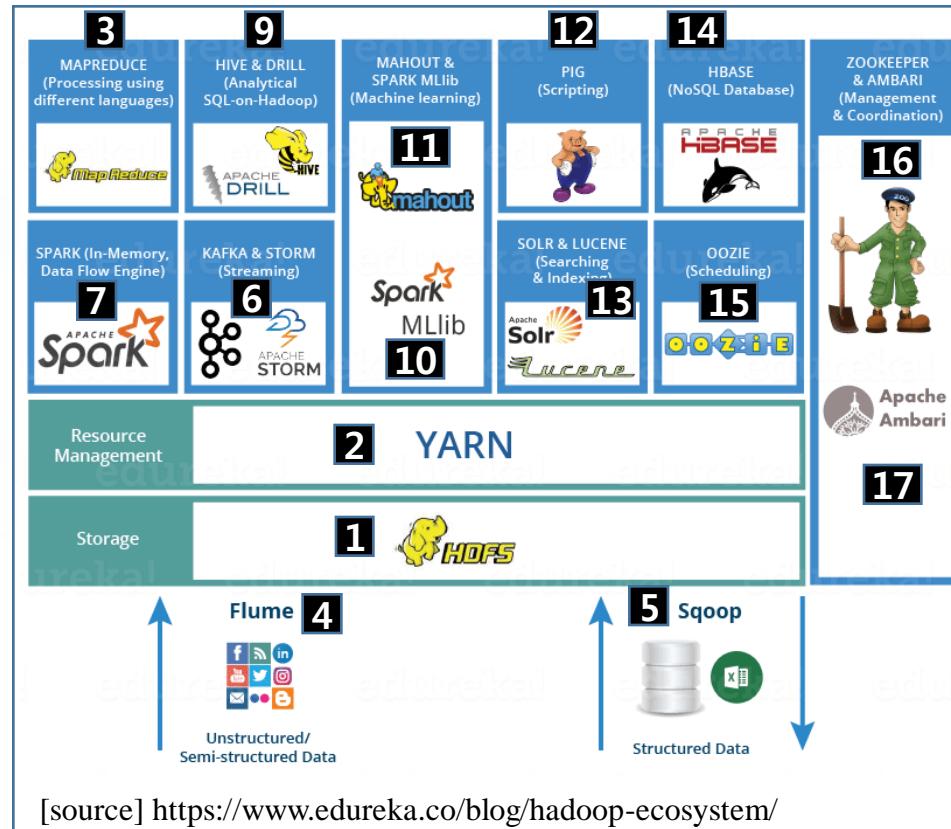
## 2. Big Data Architecture

11

- Hadoop Ecosystem

- Open Source로 제공하는 Enterprise Data Components 집합 상태계

- 2 → 1 Hadoop 자원 관리
- 1 → 3 Hadoop 파일 시스템에서의 데이터 프로세싱 방식
- 4 → 1 외부 데이터(SNS 포함)의 HDFS로 Import
- 5 ← 1 DBMS/CSV to HDFS
- 1 → 9 SQL on Hadoop
- 5 → 9 Any DBMS table to Hive Table
- 15 → 1 Hadoop Scheduling
  - Workflow 제공
- 14 HDFS 기반 NoSQL database
- 6 실시간 스트리밍 연산 처리
- 12 HDFS processing 전용 Scripter
- 13 검색 엔진
- 16 Hadoop Storage 확장, 분산 처리 관리 코디네이터
- 17 Hadoop 모니터링 Browser



AI 관점

- Python + Tensorflow의 근시안에서 벗어나기
- 분산 시스템 프레임의 데이터 Flow의 기술적 이해는 Artificial Intelligence 전문가로 가는 길의 매우 충실한 축진제 역할

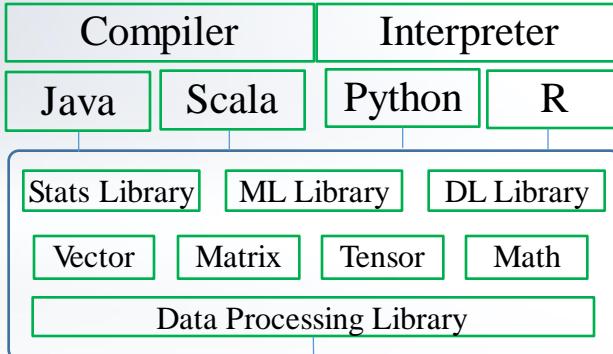
- Local file, hdfs, dbms 등 대부분의 데이터의 inmemory 기반 프로세싱(streaming 포함)
  - org.apache.spark.sql
- 7 → 1 hdfs data를 memory로 load 하여 101배 빠른 데이터 프로세싱
- 7 → 9 Hive on Spark
- 10 Spark Machine Learning Library
  - org.apache.spark.ml
  - org.apache.spark.mllib
- 11 Spark Machine Learning Library
- 12 Map Reduce 기반의 Machine Learning Package
  - org.apache.mahout

## 2. Big Data Architecture

- Hadoop Ecosystem

- Language를 이용한 Big Data 접근 → AI 구현의 첫단계

### 1 Analytics Language



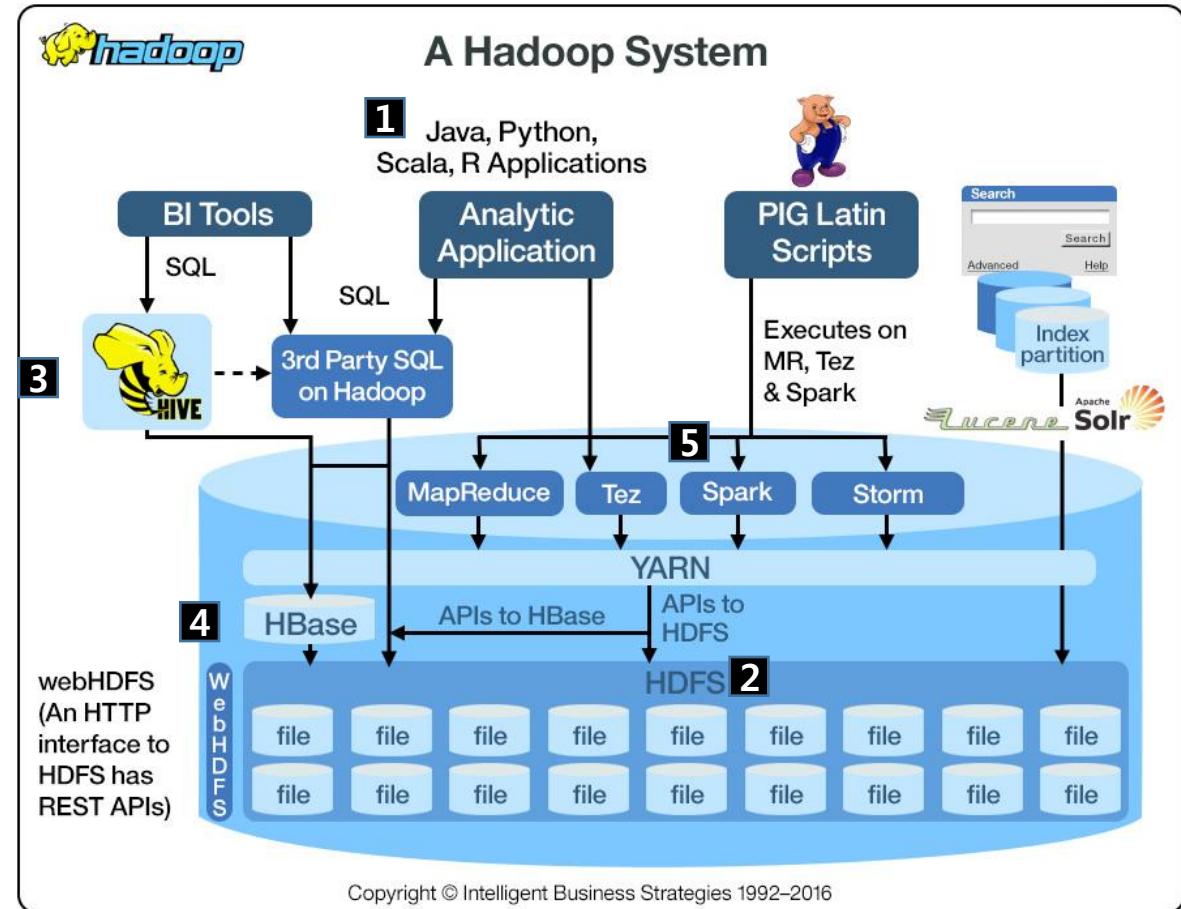
### 2 hdfs://<host:port>/<filepath>

- Flat File : csv, json, parquet, ...
- Image : jpg, png, gif, svg, ...
- Audio : mp3, wav, ogg, ...
- Video : mp4, avi, wmv, mpeg, ...

3 import org.apache.hadoop.hive  
select \* from table\_a ...

4 import org.apache.hadoop.habse  
select \* from table\_a ...

5 import org.apache.spark  
select \* from table\_a ...



[source] <http://www.ibmbigdatahub.com/blog/what-hadoop/>

## 2. Big Data Architecture

12



### Scale @ Netflix

- 125M+ active members
- 190 countries
- 450B+ unique events/day
- 700+ Kafka topics



 **SPARK+AI**  
**SUMMIT 2018**  
ORGANIZED BY  databricks

 SPARK+AI  
SUMMIT 2018

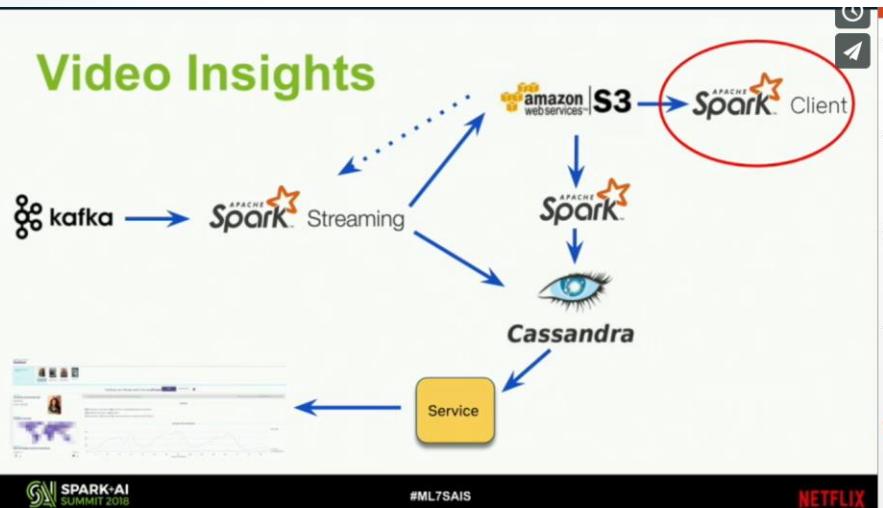
#ML7SAIS

NETFLIX

## 2. Big Data Architecture

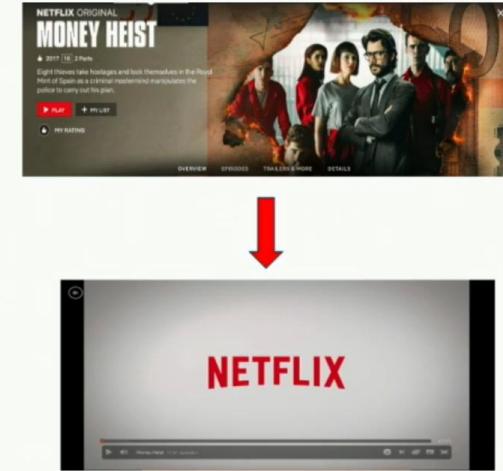
12

### Video Insights



### Required Data

- Impressions, Plays, etc.
- Attribution
- Explore/Exploit Metadata



### Billboard Recommendations

- Recommend a relevant title to each member
- Right time
- Respond quickly to member feedback



### Artwork Personalization

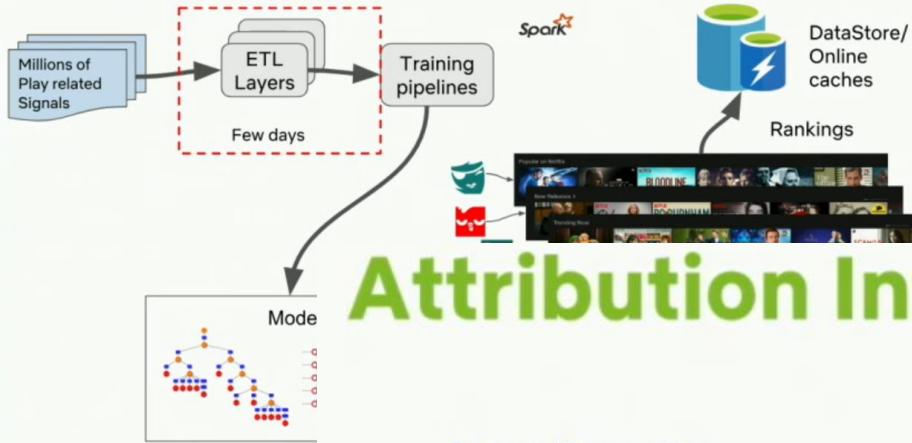
- Personalized Image
- Visual Evidence
- Quickly adapting - Title launches, member tastes
- Rapid learning - Cold start



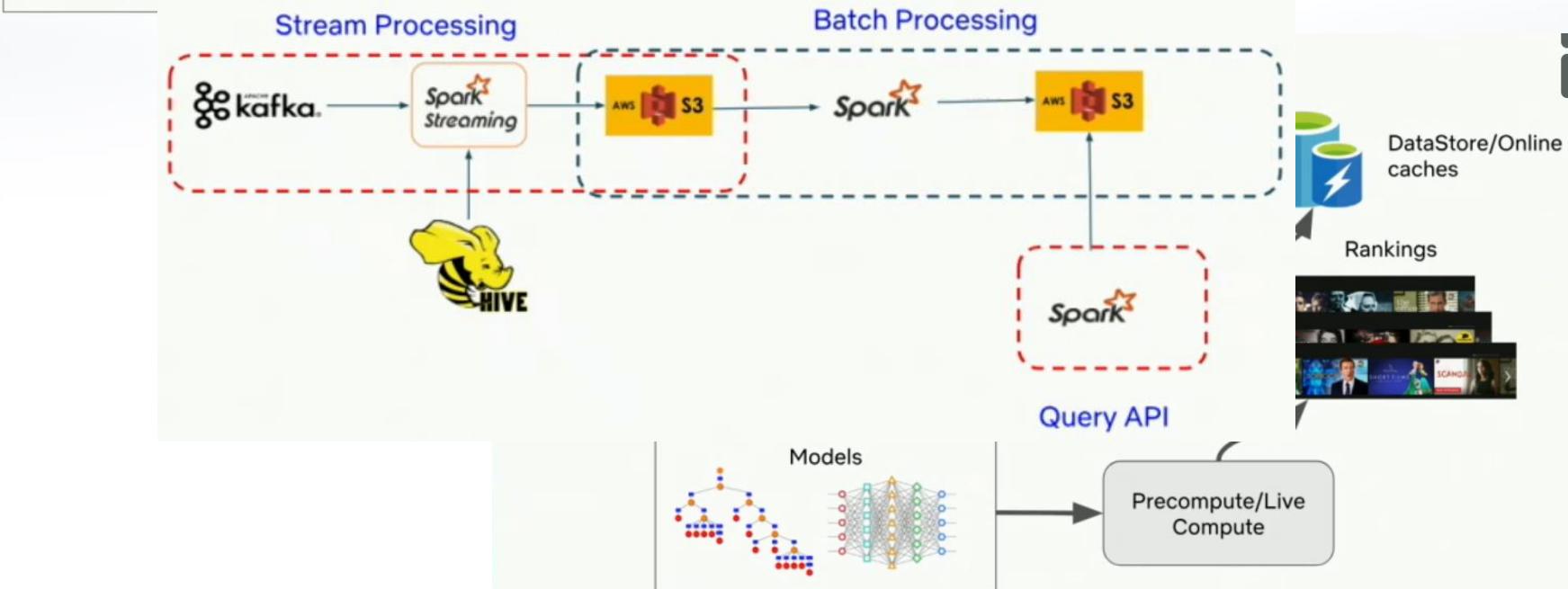
## 2. Big Data Architecture

12

### Traditional Recommendations



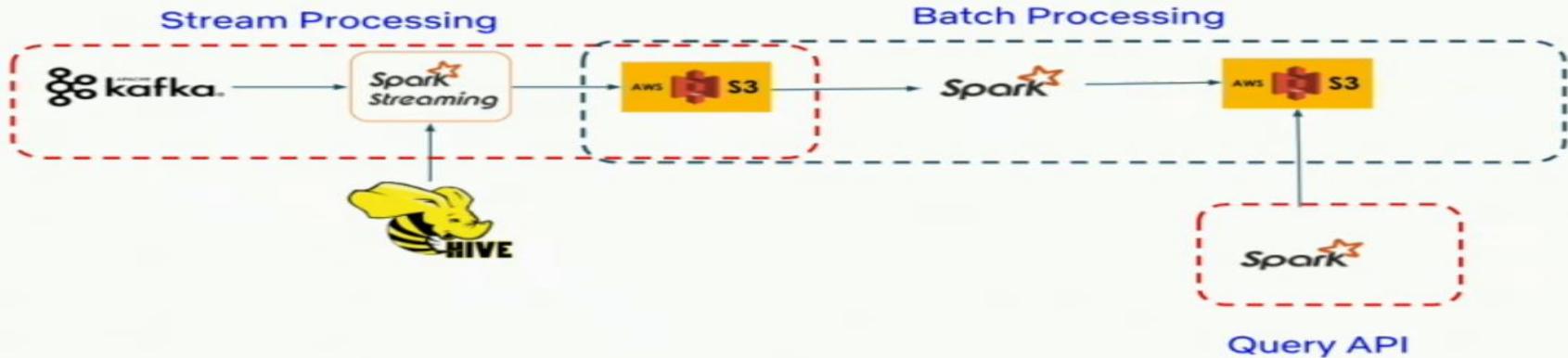
## Attribution Infrastructure



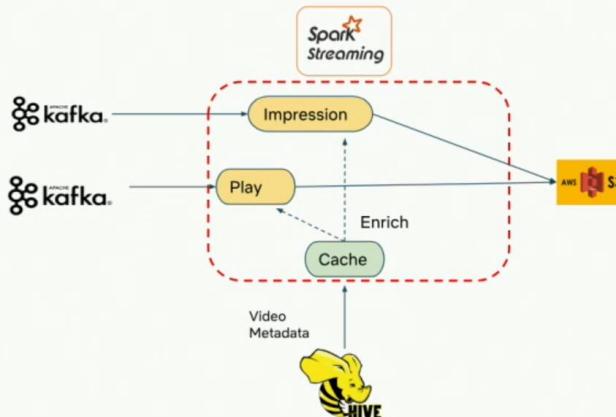
## 2. Big Data Architecture

12

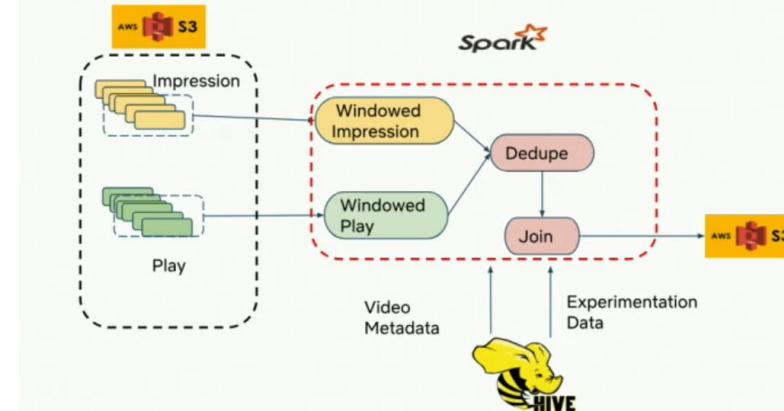
### Attribution Infrastructure



#### Stream Processing - Zoomed in

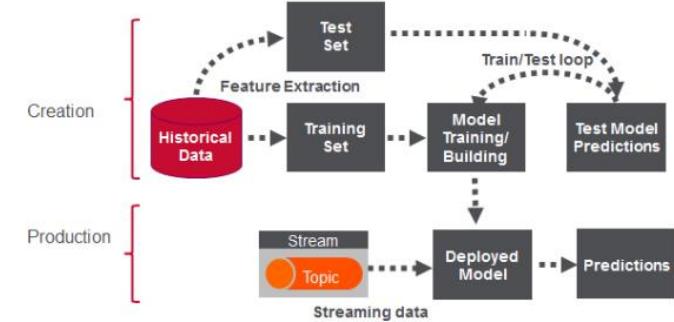
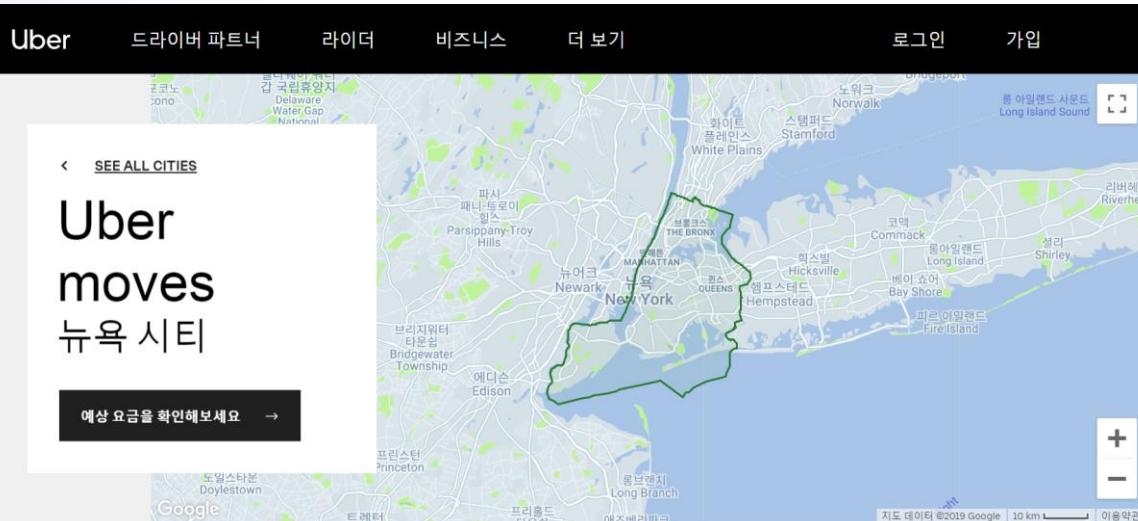


#### Batch Processing - Zoomed in



## 2. Big Data Architecture

12

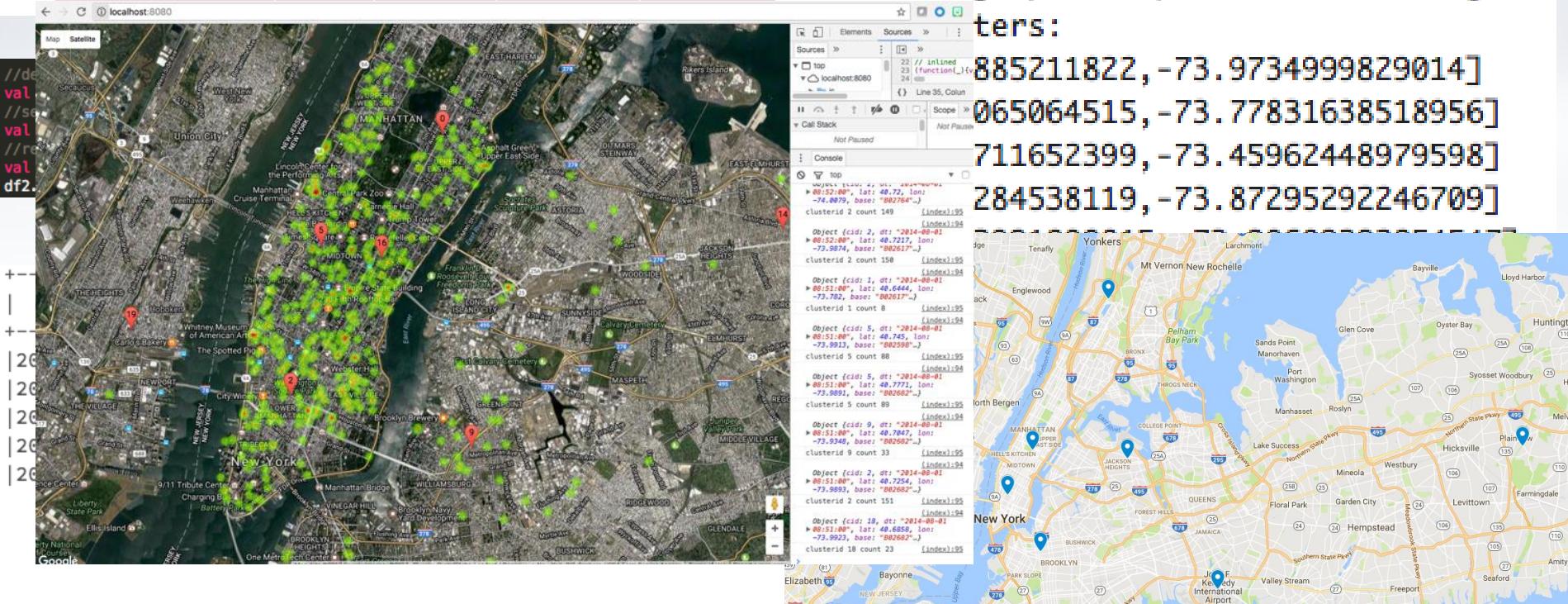
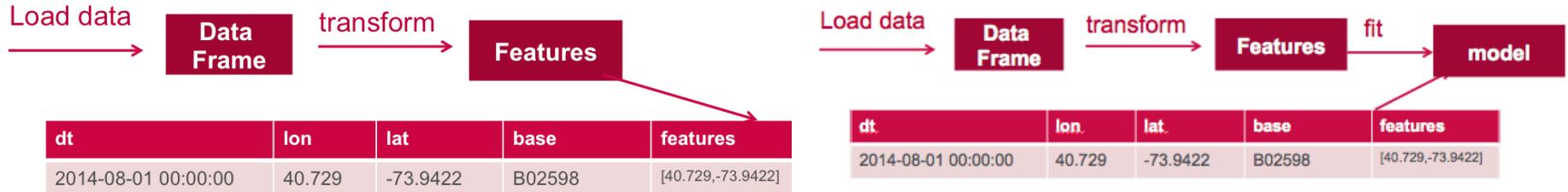


```
import org.apache.spark._  
import org.apache.spark.sql.SQLContext  
import org.apache.spark.sql.types._  
import org.apache.spark.sql.functions._  
import org.apache.spark.sql._  
import org.apache.spark.ml.feature.VectorAssembler  
import org.apache.spark.ml.clustering.KMeans
```

```
val sqlContext = new SQLContext(sc)  
import sqlContext.implicits._  
import sqlContext._  
  
val schema = StructType(Array(  
  StructField("dt", TimestampType, true),  
  StructField("lat", DoubleType, true),  
  StructField("lon", DoubleType, true),  
  StructField("base", StringType, true)  
))
```

# I. Big Data Architecture 와 AI 시스템의 Position

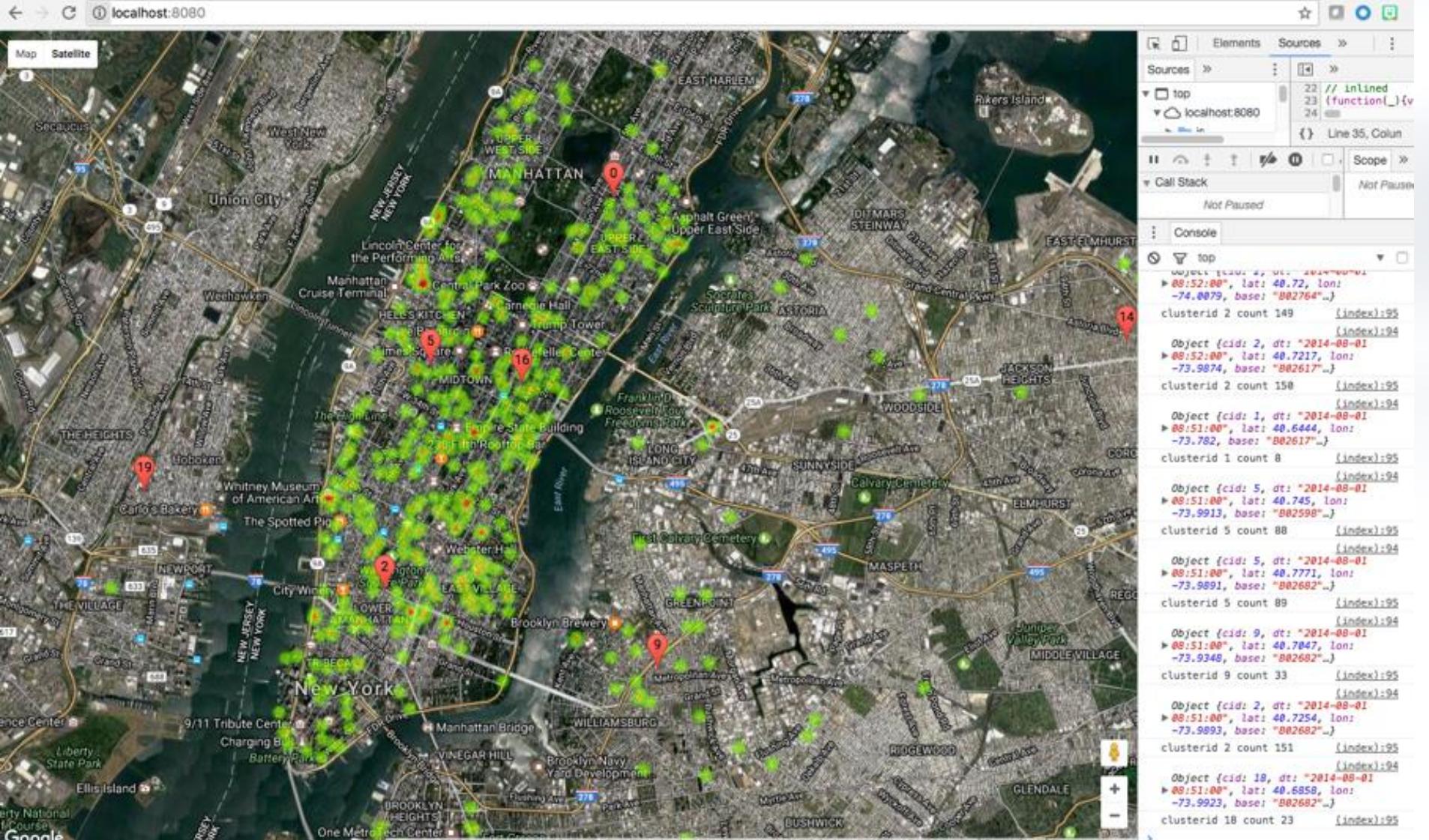
## 2. Big Data Architecture



# I. Big Data Architecture 와 AI 시스템의 Position

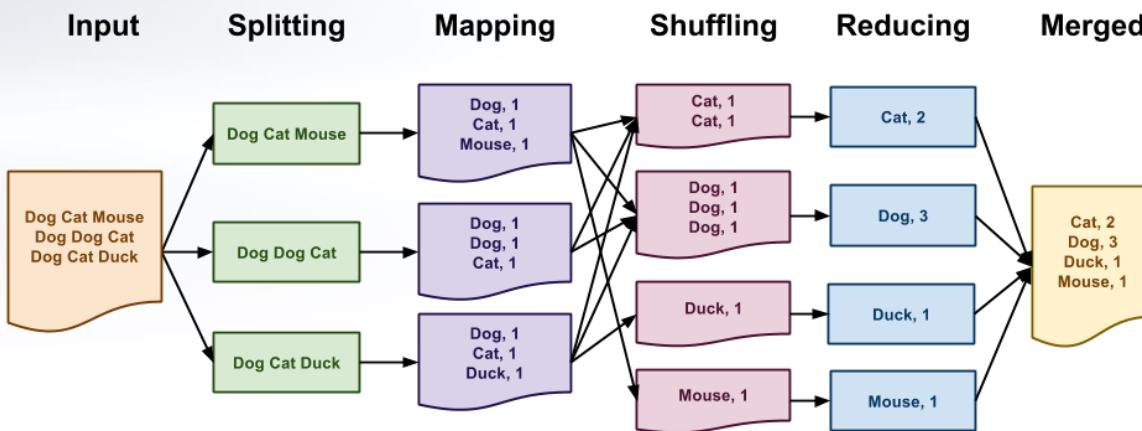
## 2. Big Data Architecture

12



## 2. Big Data Architecture

- Hadoop Map Reduce Technology
  - Mapper
  - Reducer
  - Shuffling



[source] <https://www3.nd.edu/~pbui/teaching/cse.30331.fa16/challenge11.html/>

```

package co.edureka.mapreduce;
import java.io.IOException;
import java.util.StringTokenizer;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;
import org.apache.hadoop.mapreduce.Reducer;
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.input.TextInputFormat;
import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.fs.Path;

public class WordCount {
    public static class Map extends Mapper<LongWritable, Text, Text, IntWritable> {
        public void map(LongWritable key, Text value, Context context) throws IOException {
            StringTokenizer tokenizer = value.toString();
            String line = value.toString();
            while (tokenizer.hasMoreTokens()) {
                value.set(tokenizer.nextToken());
                context.write(value, new IntWritable(1));
            }
        }
    }

    public static class Reduce extends Reducer<Text, IntWritable, Text, IntWritable> {
        public void reduce(Text key, Iterable<IntWritable> values, Context context) throws IOException {
            int sum=0;
            for(IntWritable x: values) {
                sum+=x.get();
            }
            context.write(key, new IntWritable(sum));
        }
    }

    public static void main(String[] args) throws Exception {
        Configuration conf= new Configuration();
        Job job = new Job(conf, "My Word Count Program");
        job.setJarByClass(WordCount.class);
        job.setMapperClass(Map.class);
        job.setReducerClass(Reduce.class);
        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);
        job.setInputFormatClass(TextInputFormat.class);
        job.setOutputFormatClass(TextOutputFormat.class);
        Path outputPath = new Path(args[1]);
        //Configuring the input/output path from the filesystem into the job
        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));
        //deleting the output path automatically from hdfs so that we don't have to do
        outputPath.getFileSystem(conf).delete(outputPath);
        //exiting the job only if the flag value becomes false
        System.exit(job.waitForCompletion(true) ? 0 : 1);
    }
}

```

- 대부분의 분석용 data는 이미 MapReduce 연산 처리가 되어 있으므로, 특별한 경우를 제외하고는 분석가가 MapReduce native Coding 작업을 할 필요가 없음.
- 하지만 Data Scientist는 저수준의 MapReduce Operation 구조를 이해하고 있어야, IT 조직과 원활한 소통이 가능

## 2. Big Data Architecture

14

### ▪ Hadoop Map Reduce 가상 체험 실습

- 분산시스템에서의 병렬 처리로 수행되는 Data 연산 과정의 Look & Feel
- AI Service Server 구축을 위해 수없이 반복해야 할 Software 설치/재설치 및 환경 Setting에 익숙해지기
- Windows에 비해 다소 불편(?) 할 수 있는 Linux 중심의 Open Source에 감사하는 Data Scientist 지향적 Mind-Set 전환

설치할 S/W	Link
Virtual Box 6.0.6	<ul style="list-style-type: none"><li>▪ <a href="https://download.virtualbox.org/virtualbox/6.0.6/VirtualBox-6.0.6-130049-Win.exe">https://download.virtualbox.org/virtualbox/6.0.6/VirtualBox-6.0.6-130049-Win.exe</a></li></ul>
Centos 7	<ul style="list-style-type: none"><li>▪ <a href="http://isoredirect.centos.org/centos/7/isos/x86_64/CentOS-7-x86_64-Minimal-1810.iso">http://isoredirect.centos.org/centos/7/isos/x86_64/CentOS-7-x86_64-Minimal-1810.iso</a></li></ul>
Hadoop 3.2.0	<ul style="list-style-type: none"><li>▪ <a href="https://www.apache.org/dyn/closer.cgi/hadoop/common/hadoop-3.2.0/hadoop-3.2.0.tar.gz">https://www.apache.org/dyn/closer.cgi/hadoop/common/hadoop-3.2.0/hadoop-3.2.0.tar.gz</a></li></ul>

AI 관점

- Huge Data 접근의 Entry Point – Data Scientist의 Road Map
  - Peta Data Processing
  - 30억 Rows + 1000개의 Columns에서의 Fact Finding
  - 전세계 Web 정보 전체에서의 New Trend Discovery
  - 1일 0.1초당 측정되는 10000개 Sensor에 나타난 Pattern

Rethink Business

Scale-out !!!

## 2. Big Data Architecture

### Hadoop: Setting up a Single Node Cluster.

- [Purpose](#)
- [Prerequisites](#)
  - [Supported Platforms](#)
  - [Required Software](#)
  - [Installing Software](#)
- [Download](#)
- [Prepare to Start the Hadoop Cluster](#)
- [Standalone Operation](#)
- [Pseudo-Distributed Operation](#)
  - [Configuration](#)
  - [Setup passphraseless ssh](#)
  - [Execution](#)
  - [YARN on a Single Node](#)
- [Fully-Distributed Operation](#)



```
export JAVA_HOME=/usr/java/latest
```

```
$ bin/hadoop
```

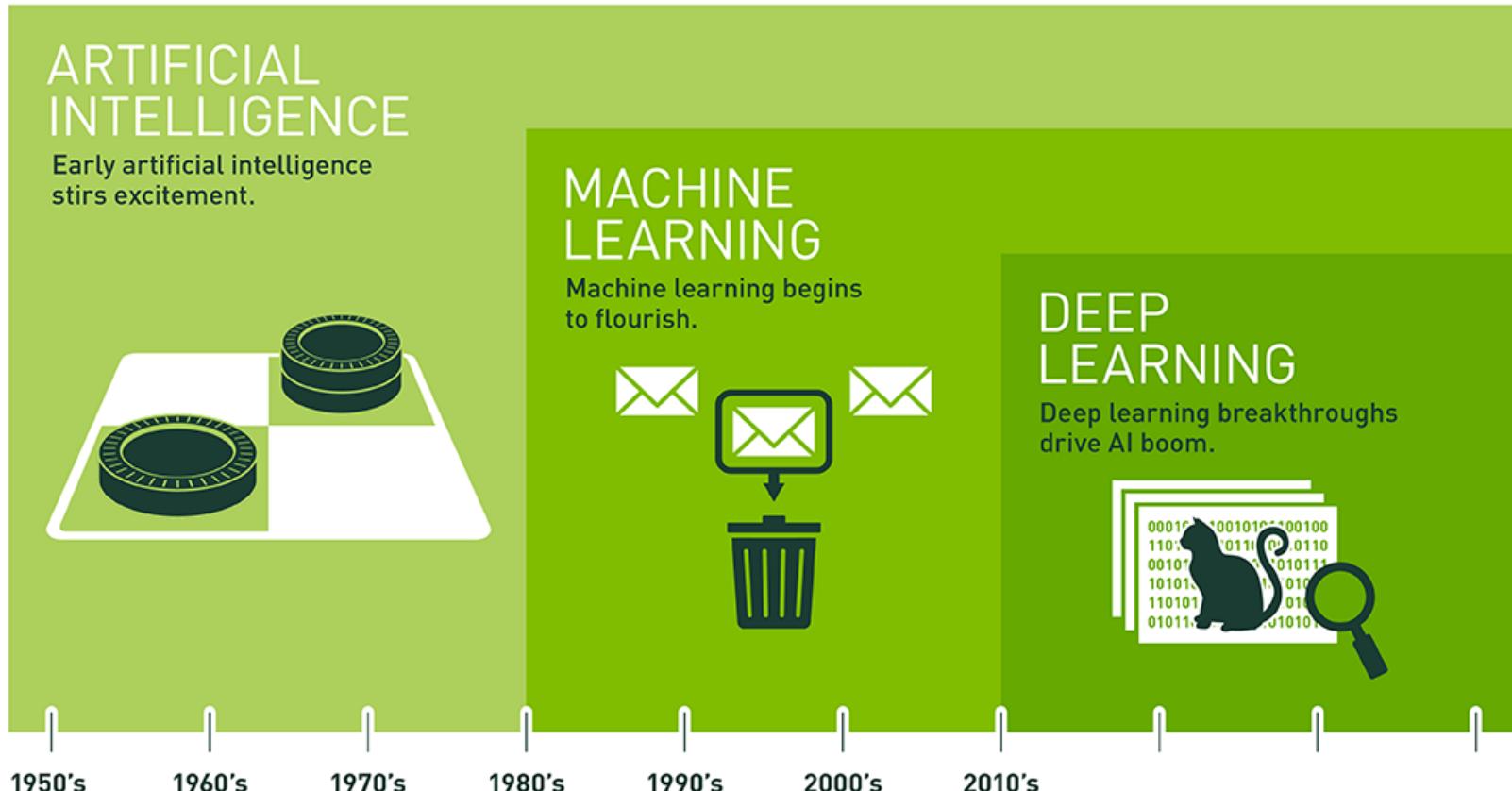
```
$ mkdir input
$ cp etc/hadoop/*.xml input
$ bin/hadoop jar share/hadoop/mapreduce/
  hadoop-mapreduce-examples-3.2.0.jar g
  rep input output 'dfs[a-z.]+'
$ cat output/*
```

[source] <https://hadoop.apache.org/docs/r3.2.0/hadoop-project-dist/hadoop-common/SingleCluster.html/>

## 3. AI, Machine Learning and Deep Learning

16

- AI, ML, DL



Since an early flush of optimism in the 1950s, smaller subsets of artificial intelligence – first machine learning, then deep learning, a subset of machine learning – have created ever larger disruptions.

[source] <https://blogs.nvidia.com/blog/2016/07/29/whats-difference-artificial-intelligence-machine-learning-deep-learning-ai/>

## 3. AI, Machine Learning and Deep Learning

Artificial  
Intelligence  
정의

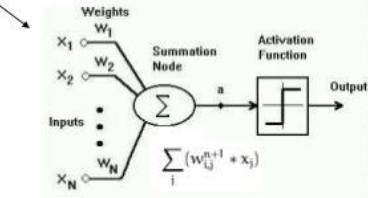
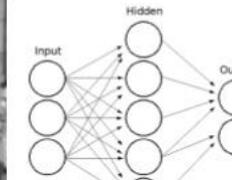
Wikipedia

- ▶ Machine Intelligence
- ▶ Intelligent Agent
- ▶ Learning & Problem Solving



Artificial Intelligence

1957: Frank Rosenblatt's Perceptron, the first artificial neural network



- 1) In computer science, artificial intelligence (AI), sometimes called **machine intelligence**, is intelligence demonstrated by machines, in contrast to the natural intelligence displayed by humans and animals.
- 2) Computer science defines AI research as the study of "**intelligent agents**": any device that perceives its environment and takes actions that maximize its chance of successfully achieving its goals.
- 3) Colloquially, the term "artificial intelligence" is used to describe machines that mimic "cognitive" functions that humans associate with other human minds, such as "**learning**" and "**problem solving**".

## 3. AI, Machine Learning and Deep Learning

Machine  
Learning  
정의

Wikipedia

- ▶ Machine : Computer
  - ▶ Learning : the ability to learn without being explicitly programmed
- [source] IBM 1959

"[Machine Learning is the] field of study that gives computers the ability to learn without being explicitly programmed." – Arthur Samuel (1959)



- 1) Evolved from the study of **pattern recognition** and **computational learning theory** in **artificial intelligence**  
(<http://www.britannica.com/EBchecked/topic/1116194/machine-learning>)
- 2) Machine learning is closely related to (and often overlaps with) computational **statistics**, which also focuses on **prediction-making** through the use of computers.
- 3) It has strong ties to **mathematical optimization**, which delivers methods, theory and application domains to the field.
- 4) Machine learning is sometimes conflated with **data mining**
- 5) Machine learning can also be **unsupervised** and be used to learn and establish baseline behavioral profiles for various entities and then used to find meaningful anomalies.

## 3. AI, Machine Learning and Deep Learning

19

Deep  
Learning  
정의

- ▶ Broader Family of Neural Network Methods
- ▶ Learning : Supervised
- ▶ Deep vs Shallow

Wikipedia

- 1) Deep learning (also known as deep structured learning or hierarchical learning) is part of a broader family of neural network methods based on convolutional neural networks (CNN)s. Learning can be supervised, semi-supervised or unsupervised.
- 2) Deep learning architectures such as deep neural networks, deep belief networks recurrent neural networks and convolutional neural networks have been applied to fields including computer vision, speech recognition, natural language processing, audio recognition, social network filtering, machine translation, bioinformatics, drug design, medical image analysis, material inspection and board game programs, where they have produced results comparable to and in some cases superior to human experts.

## 3. AI, Machine Learning and Deep Learning

20

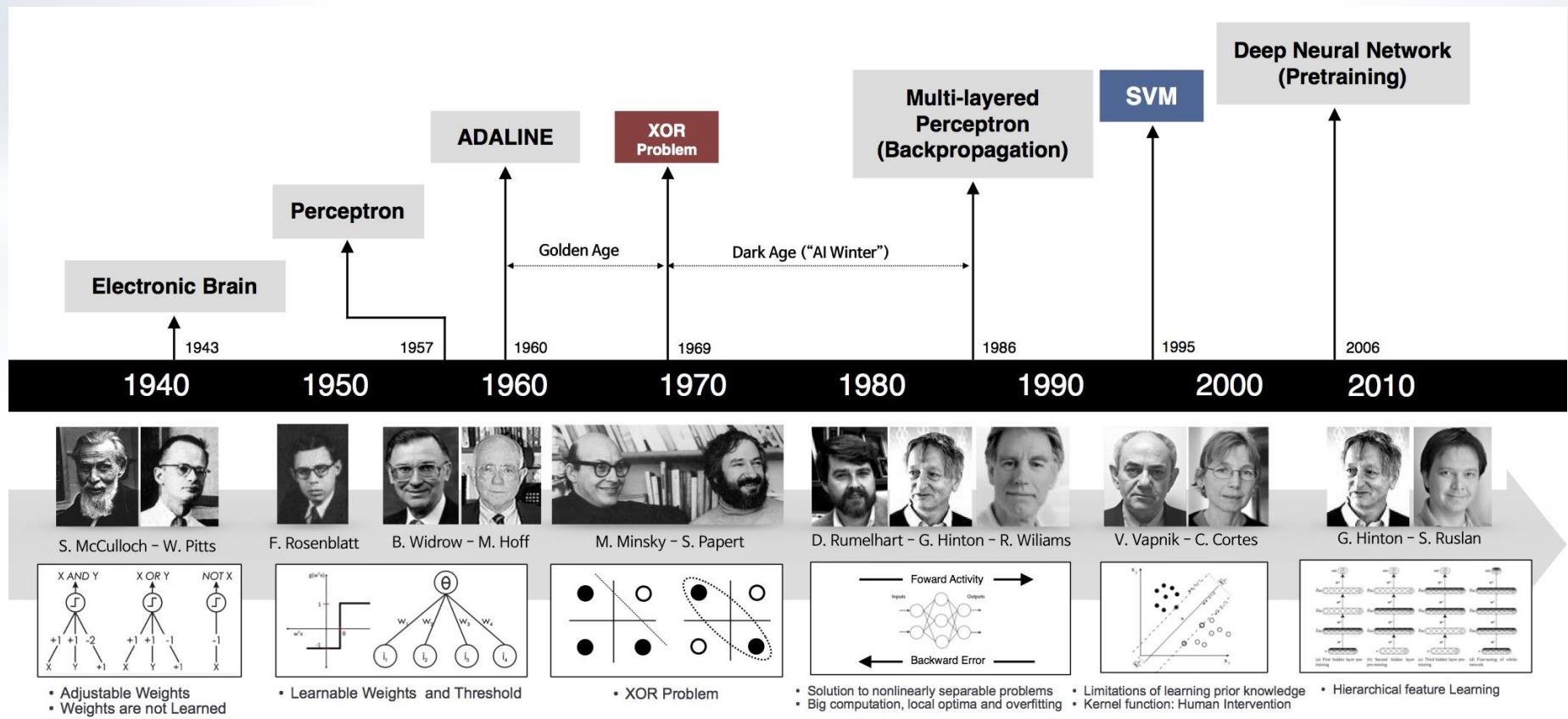
- AI, ML, DL

AI	ML	DL
1950s	1960s	1970s
실세계 시뮬레이션	의사결정	복잡한 문제 해결
▪ 인간처럼 사고하고 행동하는 Device	▪ 의사결정에 필요한 데 이터의 학습(Learning)	▪ 인간의 오감으로 발견하기 어려운 특징의 탐색
AI > ML > AI	ML > AI	ML

## 3. AI, Machine Learning and Deep Learning

### Brief History

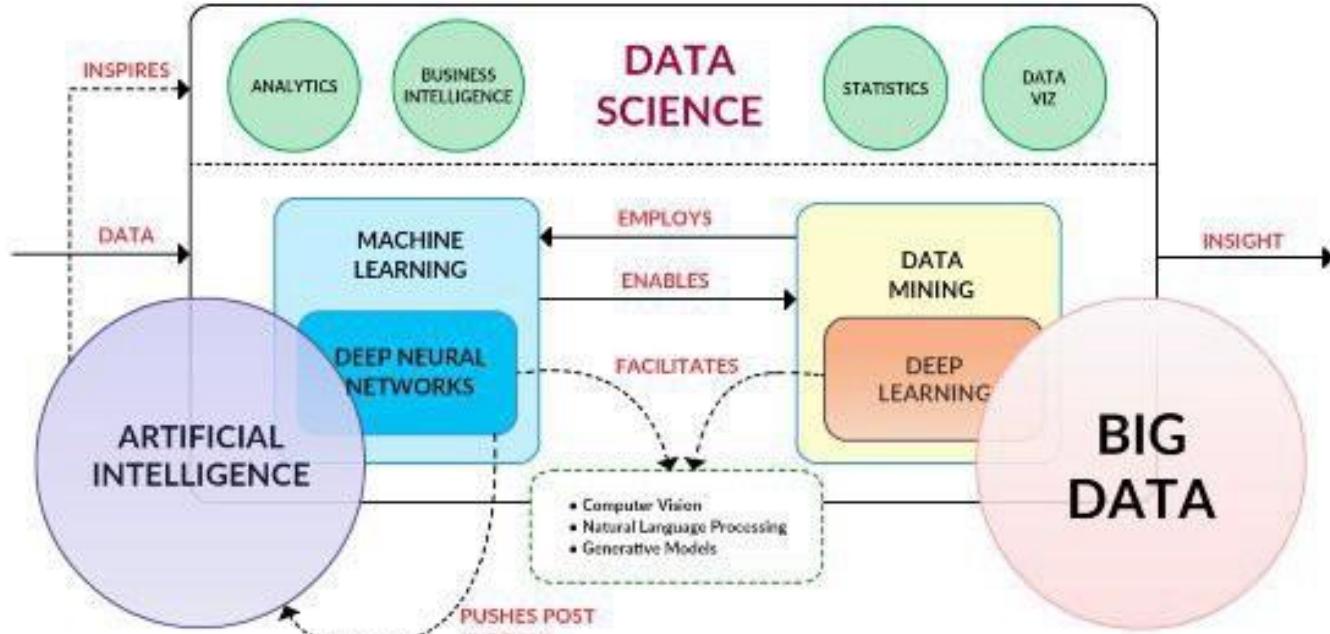
- 1970s – AI winter
- SVM의 등장 – Deep Learning 능가
- 고성능 컴퓨터 + Backpropagation의 실현으로 DL이 Main Stream으로 포지셔닝



[source] <https://medium.com/ibm-data-science-experience/deep-learning-with-data-science-experience-8478cc0f81ac>

## 4. Big Data 시스템에서의 AI의 위치와 역할

22



[source] <https://hackernoon.com/how-big-data-is-empowering-ai-and-machine-learning-4e93a1004c8f>

### ▪ AI의 R & R

- Heuristic으로 접근해서 발견할 수 없는 기회에의 Challenge
- 접근하지 못했던 실시간 Streaming Data(실시간 시장의 흐름, 공정의 상태 변화 등)의 모니터링 실현으로 상시 업무로의 전환
- 경쟁우위(Competitive Advantage)를 점할 수 있는 최초의 제품/Service/Application 창출
- 상시업무에서의 중요도가 극히 낮은 Big Data의 정체성 확인과 빅데이터 유통성의 전사 확산

## 4. Big Data 시스템에서의 AI의 위치와 역할

22

A사 사례

B사 사례

