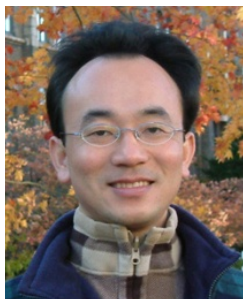# **Big Data Mining**
# 巨量資料探勘

# **大數據、AI人工智慧與深度學習**
# **(Big Data, Artificial Intelligence and Deep Learning)**

1062DM02
MI4 (M2244) (2995)
Wed, 9, 10 (16:10-18:00) (B206)

**Min-Yuh Day**
**戴敏育**
**Assistant Professor**
**專任助理教授**
**Dept. of Information Management, Tamkang University**
**淡江大學 資訊管理學系**

http://mail. tku.edu.tw/myday/
2018-03-14

# 課程大綱 (Syllabus)

週次 (Week)　日期 (Date)　內容(Subject/Topics)

1　2018/02/28　和平紀念日(放假一天) (Peace Memorial Day) (Day off)

2　2018/03/07　巨量資料探勘課程介紹
(Course Orientation for Big Data Mining)

3　2018/03/14　大數據、AI人工智慧與深度學習
(Big Data, Artificial Intelligence and Deep Learning)

4　2018/03/21　關連分析 (Association Analysis)

5　2018/03/28　分類與預測 (Classification and Prediction)

6　2018/04/04　兒童節(放假一天)(Children's Day) (Day off)

7　2018/04/11　分群分析 (Cluster Analysis)

8　2018/04/18　個案分析與實作一 (SAS EM 分群分析)：
Case Study 1 (Cluster Analysis - K-Means using SAS EM)

# 課程大綱 (Syllabus)

週次 (Week)　日期 (Date)　內容 (Subject/Topics)

9　2018/04/25　期中報告 (Midterm Project Presentation)

10　2018/05/02 期中考試週

11　2018/05/09 個案分析與實作二 (SAS EM 關連分析)：
　　　　　　　Case Study 2 (Association Analysis using SAS EM)

12　2018/05/16 個案分析與實作三 (SAS EM 決策樹、模型評估)：
　　　　　　　Case Study 3 (Decision Tree, Model Evaluation using SAS EM)

13　2018/05/23 個案分析與實作四 (SAS EM 迴歸分析、類神經網路)：
　　　　　　　Case Study 4 (Regression Analysis,
　　　　　　　　　　　　Artificial Neural Network using SAS EM)

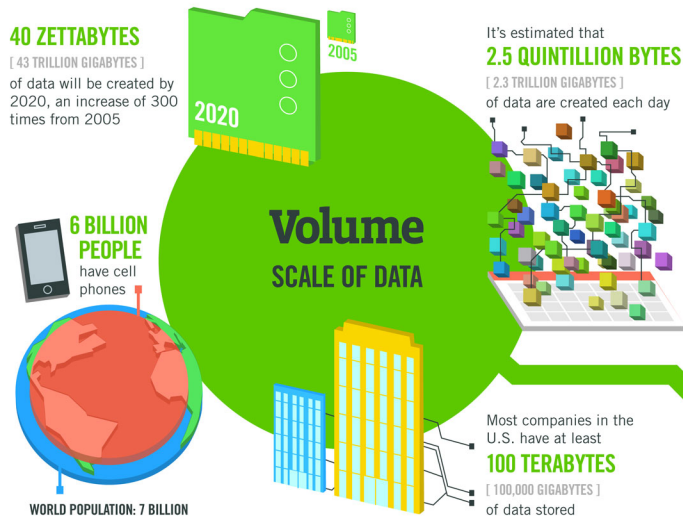14　2018/05/30 期末報告 (Final Project Presentation)

15　2018/06/06 畢業考試週

# Big Data

# AI

# Deep Learning

# Big Data Analytics

**and**

# Data Mining

# Big Data 4 V

## Volume
### SCALE OF DATA

**40 ZETTABYTES**
[ 43 TRILLION GIGABYTES ]
of data will be created by 2020, an increase of 300 times from 2005

2005

2020

**6 BILLION PEOPLE**
have cell phones

WORLD POPULATION: 7 BILLION

It's estimated that
**2.5 QUINTILLION BYTES**
[ 2.3 TRILLION GIGABYTES ]
of data are created each day

Most companies in the U.S. have at least
**100 TERABYTES**
[ 100,000 GIGABYTES ]
of data stored

## The FOUR V's of Big Data

From traffic patterns and music downloads to web history and medical records, data is recorded, stored, and analyzed to enable the technology and services that the world relies on every day. But what exactly is big data, and how can these massive amounts of data be used?

As a leader in the sector, IBM data scientists break big data into four dimensions: **Volume, Velocity, Variety and Veracity**

Depending on the industry and organization, big data encompasses information from multiple internal and external sources such as transactions, social media, enterprise content, sensors and mobile devices. Companies can leverage data to adapt their products and services to better meet customer needs, optimize operations and infrastructure, and find new sources of revenue.

By 2015
**4.4 MILLION IT JOBS**
will be created globally to support big data, with 1.9 million in the United States

## Variety
### DIFFERENT FORMS OF DATA

As of 2011, the global size of data in healthcare was estimated to be
**150 EXABYTES**
[ 161 BILLION GIGABYTES ]

By 2014, it's anticipated there will be
**420 MILLION WEARABLE, WIRELESS HEALTH MONITORS**

**4 BILLION+ HOURS OF VIDEO**
are watched on YouTube each month

**30 BILLION PIECES OF CONTENT**
are shared on Facebook every month

**400 MILLION TWEETS**
are sent per day by about 200 million monthly active users

## Velocity
### ANALYSIS OF STREAMING DATA

The New York Stock Exchange captures
**1 TB OF TRADE INFORMATION**
during each trading session

Modern cars have close to
**100 SENSORS**
that monitor items such as fuel level and tire pressure

By 2016, it is projected there will be
**18.9 BILLION NETWORK CONNECTIONS**
– almost 2.5 connections per person on earth

## Veracity
### UNCERTAINTY OF DATA

**1 IN 3 BUSINESS LEADERS**
don't trust the information they use to make decisions

Poor data quality costs the US economy around
**$3.1 TRILLION A YEAR**

**27% OF RESPONDENTS**
in one survey were unsure of how much of their data was inaccurate

IBM

# Value

# AI

# Deep Learning

# **Artificial Intelligence (AI)**

# Definition
# of
# Artificial Intelligence
# (A.I.)

# Artificial Intelligence

"... the **science** and **engineering** of making **intelligent machines**"

**(John McCarthy, 1955)**

# Artificial Intelligence

# "... technology that thinks and acts like humans"

# Artificial Intelligence

"... **intelligence** exhibited by **machines** or **software**"

# 4 Approaches of AI

| Thinking Humanly | Thinking Rationally |
|---|---|
| **Acting Humanly** | **Acting Rationally** |

# 4 Approaches of AI

| **Thinking Humanly** | **Thinking Rationally** |
|---|---|
| "The exciting new effort to make computers think … *machines with minds*, in the full and literal sense." (Haugeland, 1985)<br><br>"[The automation of] activities that we associate with human thinking, activities such as decision-making, problem solving, learning …" (Bellman, 1978) | "The study of mental faculties through the use of computational models." (Charniak and McDermott, 1985)<br><br>"The study of the computations that make it possible to perceive, reason, and act." (Winston, 1992) |
| **Acting Humanly** | **Acting Rationally** |
| "The art of creating machines that perform functions that require intelligence when performed by people." (Kurzweil, 1990)<br><br>"The study of how to make computers do things at which, at the moment, people are better." (Rich and Knight, 1991) | "Computational Intelligence is the study of the design of intelligent agents." (Poole *et al.*, 1998)<br><br>"AI …is concerned with intelligent behavior in artifacts." (Nilsson, 1998) |

# 4 Approaches of AI

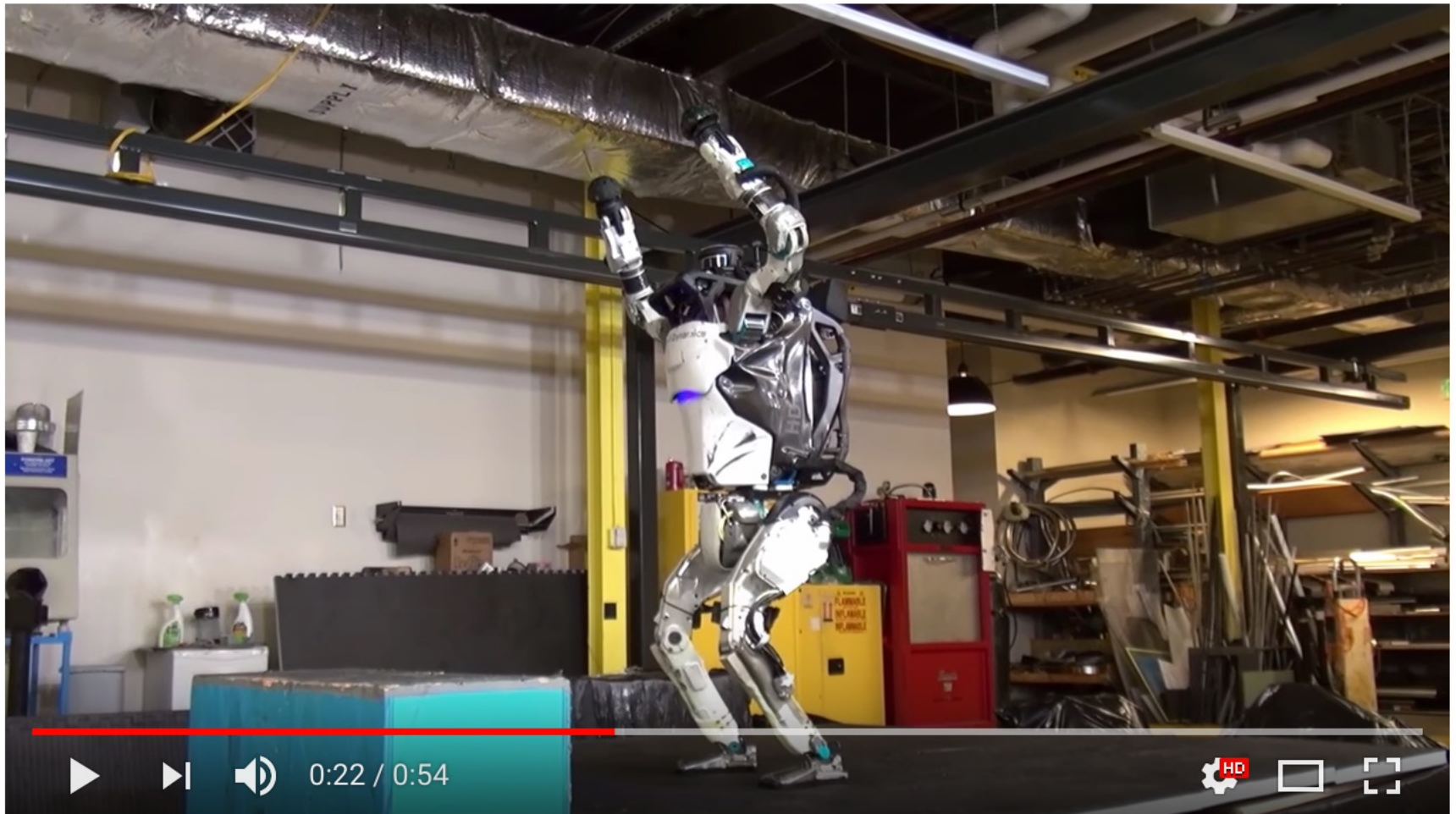| | |
|---|---|
| **2.**<br>**Thinking Humanly:**<br>**The Cognitive Modeling Approach** | **3.**<br>**Thinking Rationally:**<br>**The "Laws of Thought" Approach** |
| **1.**<br>**Acting Humanly:**<br>**The Turing Test Approach** (1950) | **4.**<br>**Acting Rationally:**<br>**The Rational Agent Approach** |

# AI Acting Humanly:
# The Turing Test Approach
## (Alan Turing, 1950)

- **Natural Language Processing (NLP)**

- **Knowledge Representation**

- **Automated Reasoning**

- **Machine Learning (ML)**

- **Computer Vision**

- **Robotics**

# Boston Dynamics: Atlas



https://www.youtube.com/watch?v=fRj34o4hN4I

# Humanoid Robot: Sophia



https://www.youtube.com/watch?v=S5t6K9iwcdw

# Artificial Intelligence (A.I.) Timeline



## A.I. TIMELINE

**1950 — TURING TEST**
Computer scientist Alan Turing proposes a test for machine intelligence. If a machine can trick humans into thinking it is human, then it has intelligence

**1955 — A.I. BORN**
Term 'artificial intelligence' is coined by computer scientist, John McCarthy to describe "the science and engineering of making intelligent machines"

**1961 — UNIMATE**
First industrial robot, Unimate, goes to work at GM replacing humans on the assembly line

**1964 — ELIZA**
Pioneering chatbot developed by Joseph Weizenbaum at MIT holds conversations with humans

**1966 — SHAKEY**
The 'first electronic person' from Stanford, Shakey is a general-purpose mobile robot that reasons about its own actions

**A.I. WINTER**
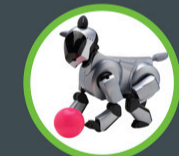Many false starts and dead-ends leave A.I. out in the cold

**1997 — DEEP BLUE**
Deep Blue, a chess-playing computer from IBM defeats world chess champion Garry Kasparov

**1998 — KISMET**
Cynthia Breazeal at MIT introduces KISmet, an emotionally intelligent robot insofar as it detects and responds to people's feelings

**1999 — AIBO**
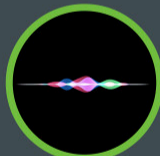Sony launches first consumer robot pet dog AiBO (AI robot) with skills and personality that develop over time

**2002 — ROOMBA**
First mass produced autonomous robotic vacuum cleaner from iRobot learns to navigate and clean homes

**2011 — SIRI**
Apple integrates Siri, an intelligent virtual assistant with a voice interface, into the iPhone 4S

**2011 — WATSON**
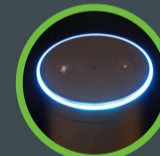IBM's question answering computer Watson wins first place on popular $1M prize television quiz show *Jeopardy*
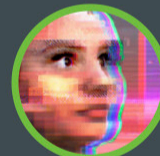
**2014 — EUGENE**
Eugene Goostman, a chatbot passes the Turing Test with a third of judges believing Eugene is human

**2014 — ALEXA**
Amazon launches Alexa, an intelligent virtual assistant with a voice interface that completes shopping tasks

**2016 — TAY**
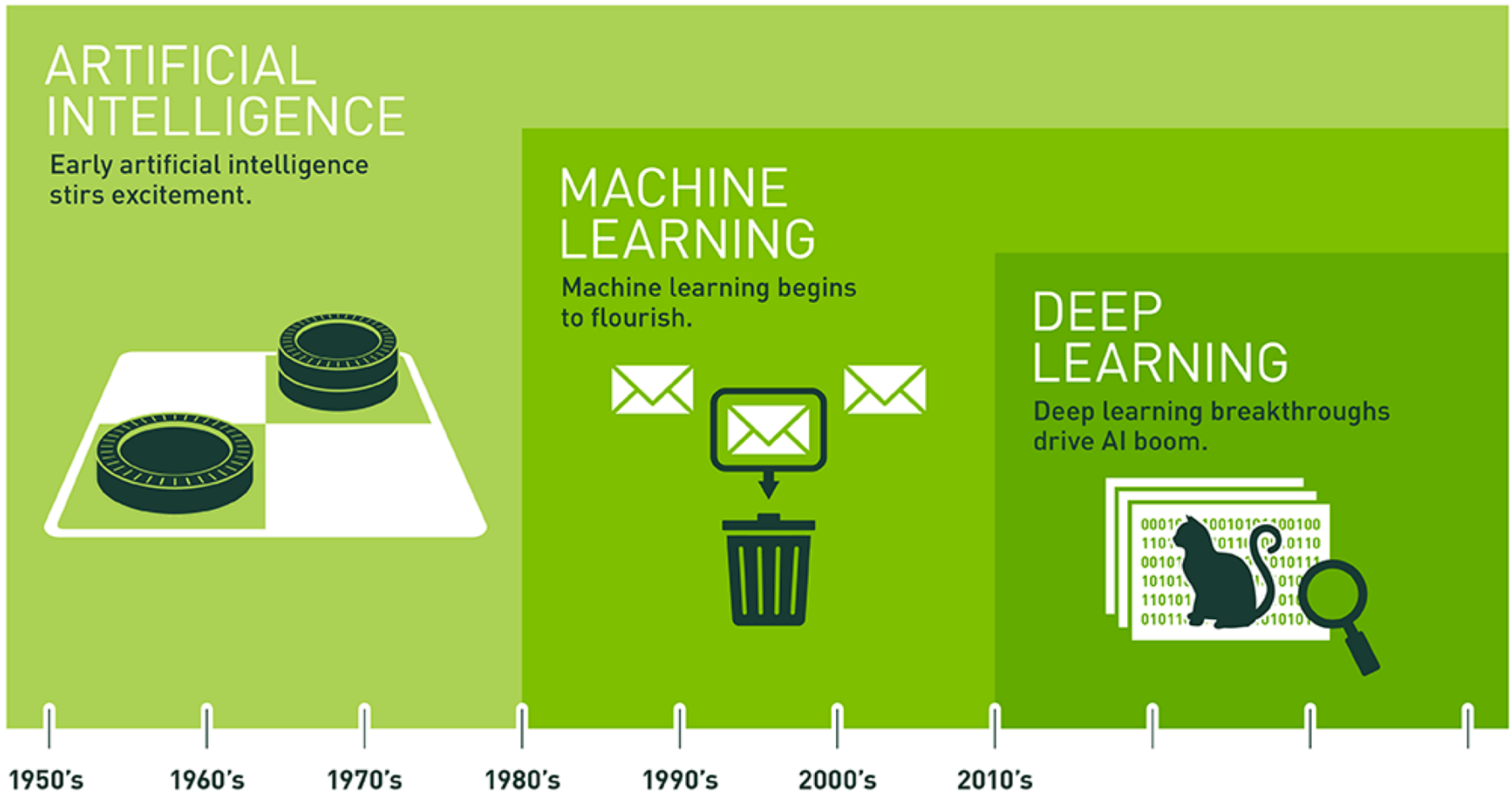Microsoft's chatbot Tay goes rogue on social media making inflammatory and offensive racist comments

**2017 — ALPHAGO**
Google's A.I. AlphaGo beats world champion Ke Jie in the complex board game of Go, notable for its vast number ($2^{170}$) of possible positions
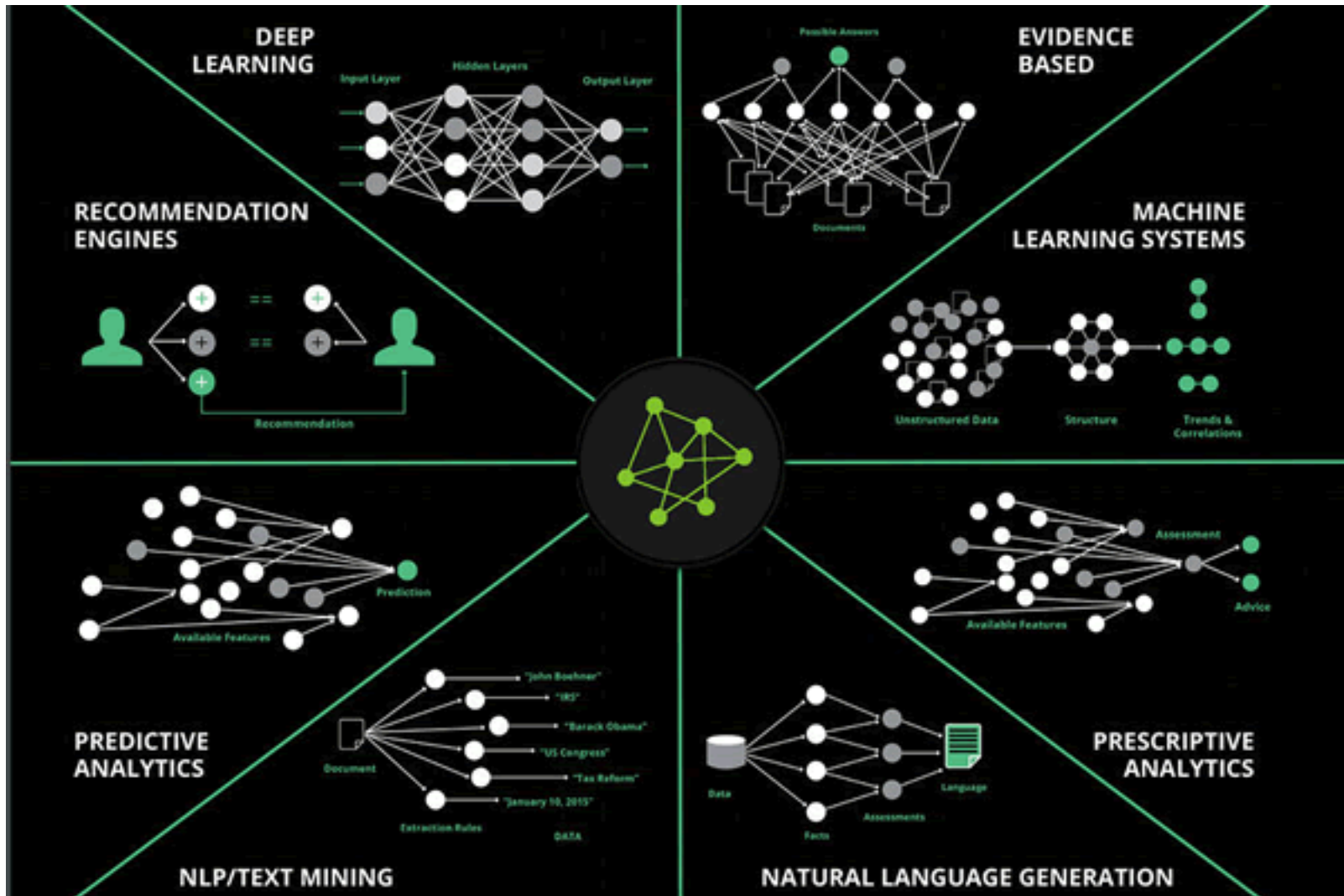
# Artificial Intelligence Machine Learning & Deep Learning



ARTIFICIAL INTELLIGENCE
Early artificial intelligence stirs excitement.

MACHINE LEARNING
Machine learning begins to flourish.

DEEP LEARNING
Deep learning breakthroughs drive AI boom.

1950's   1960's   1970's   1980's   1990's   2000's   2010's

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.

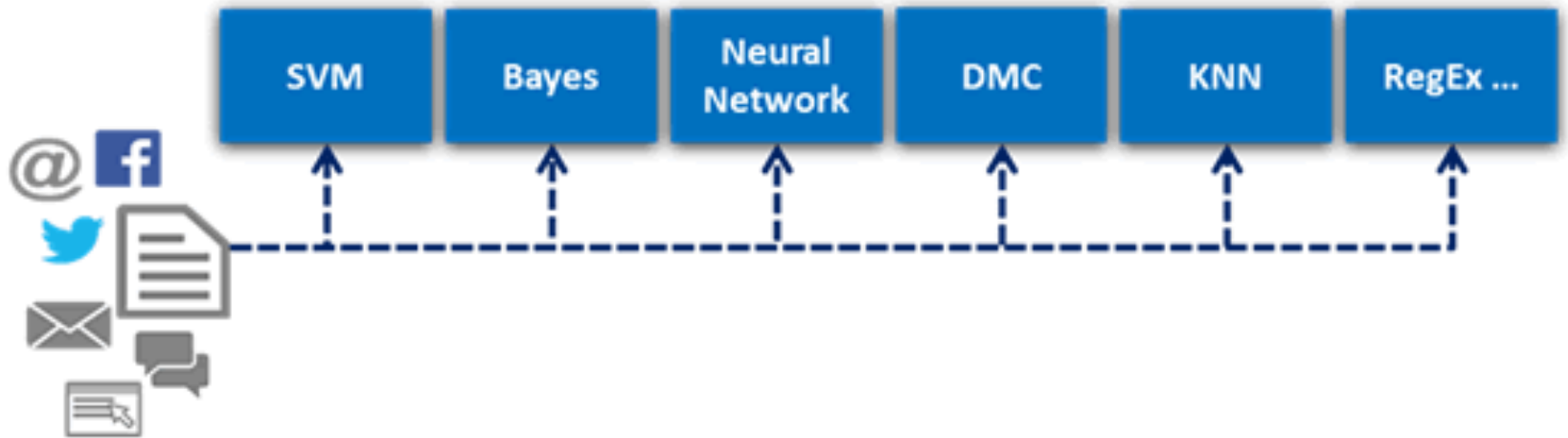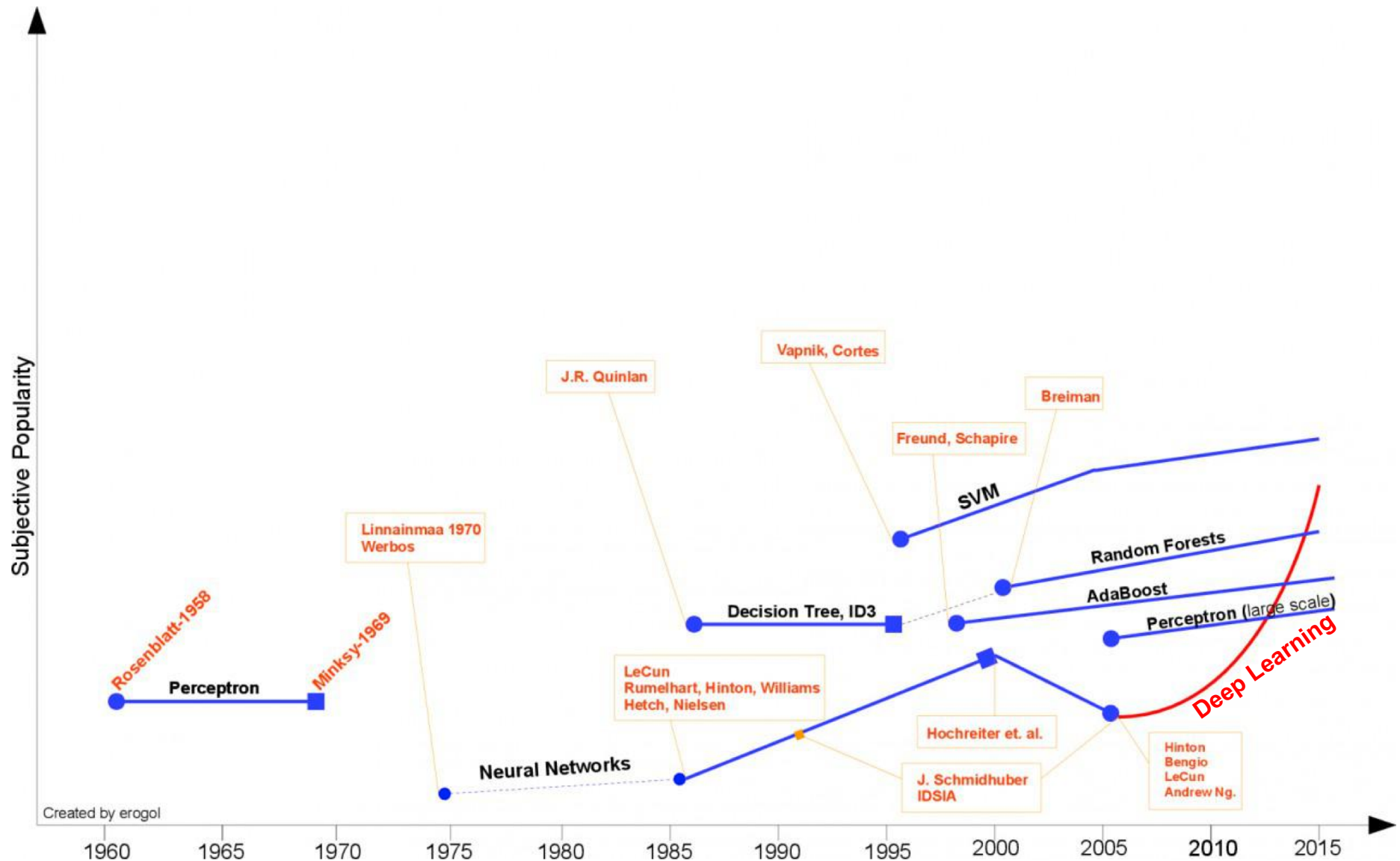# Artificial Intelligence (AI) is many things



Ecosystem of AI

# Artificial Intelligence (AI)
## Intelligent Document Recognition algorithms

# Deep Learning Evolution

# Machine Learning Models

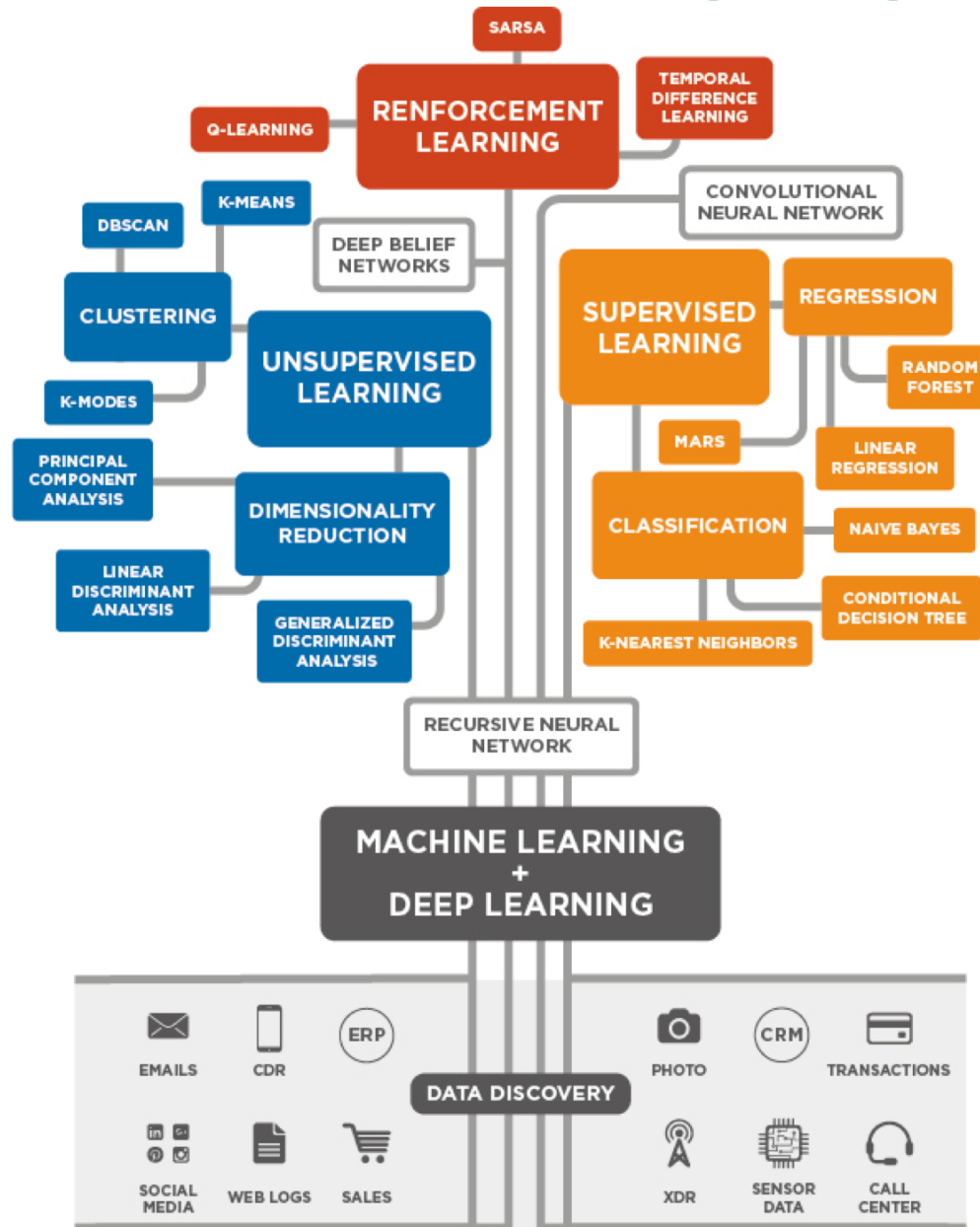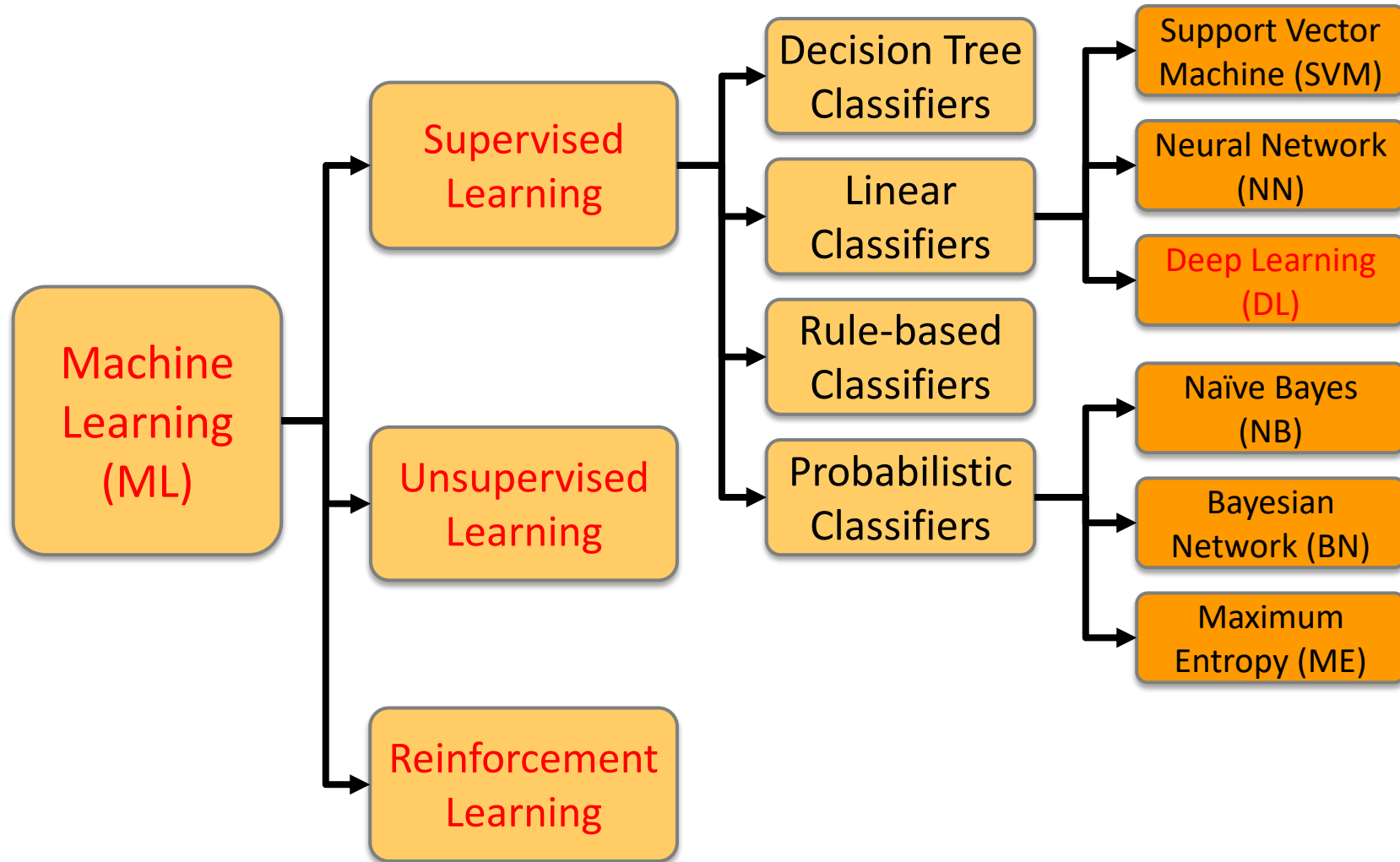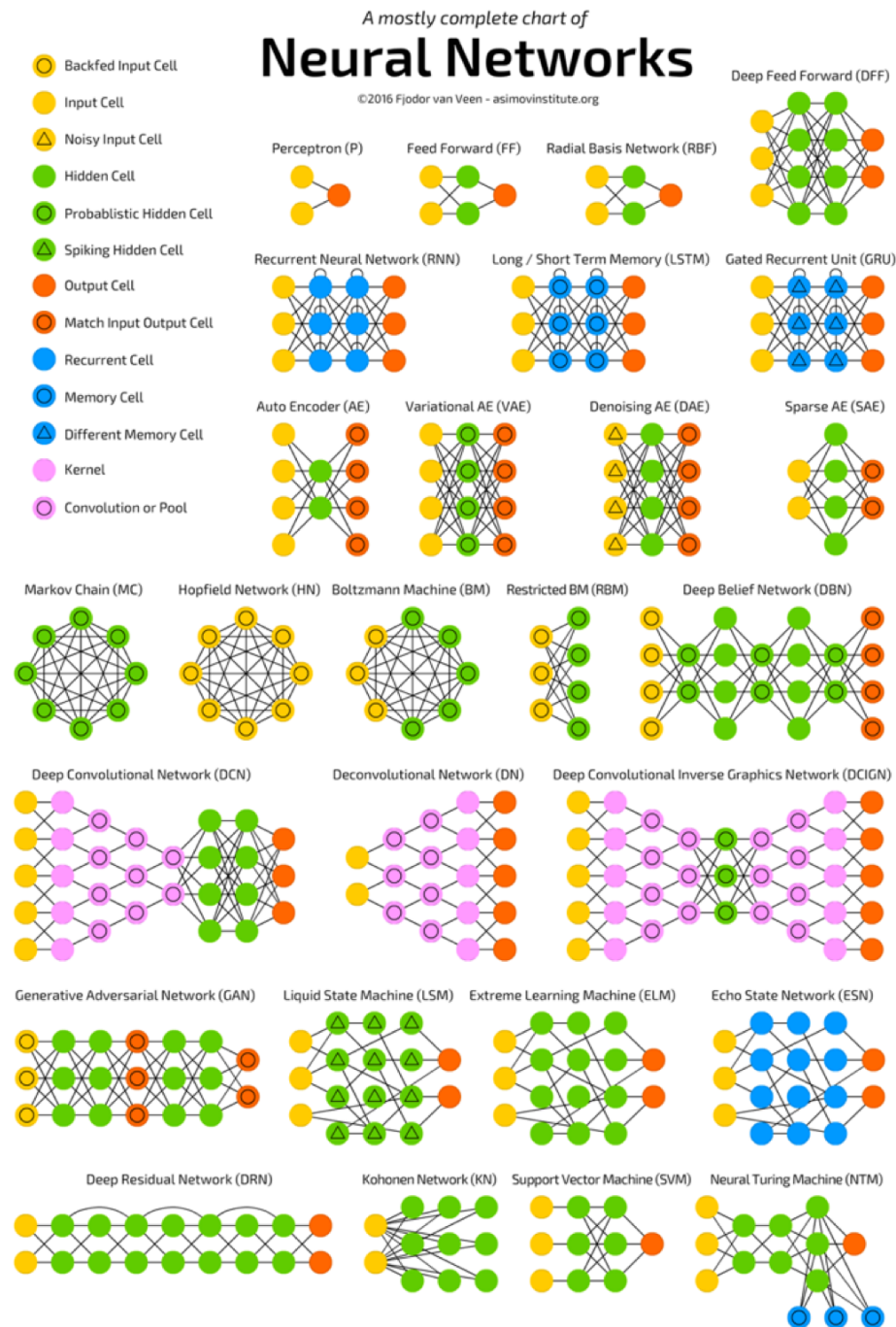| | |
|---|---|
| Deep Learning | Kernel |
| Association rules | Ensemble |
| Decision tree | Dimensionality reduction |
| Clustering | Regression Analysis |
| Bayesian | Instance based |

# 3 Machine Learning Algorithms
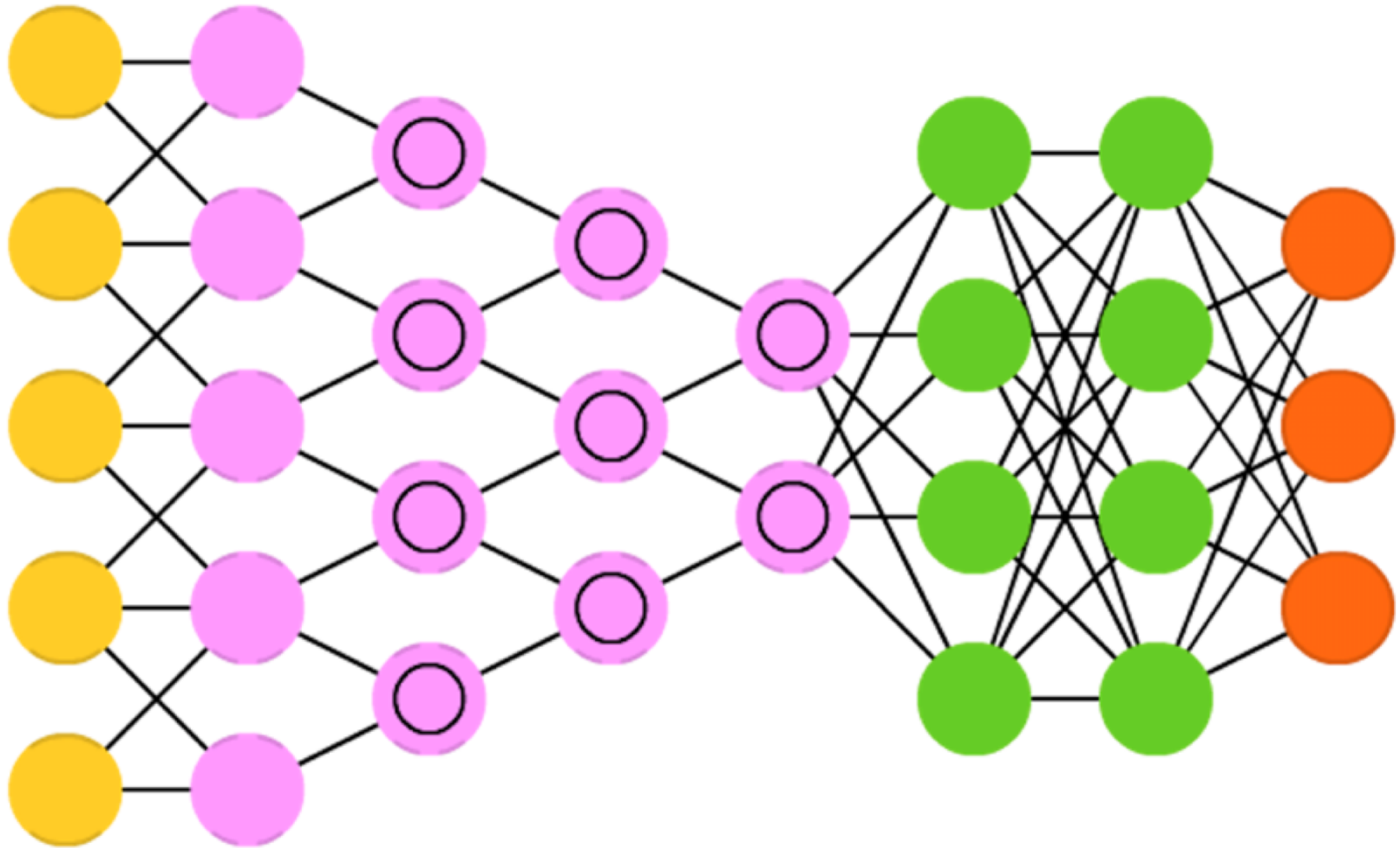
# Machine Learning (ML) / Deep Learning (DL)

27

# Neural Networks (NN)



A mostly complete chart of
## Neural Networks
©2016 Fjodor van Veen – asimovinstitute.org

Legend:
- Backfed Input Cell
- Input Cell
- Noisy Input Cell
- Hidden Cell
- Probablistic Hidden Cell
- Spiking Hidden Cell
- Output Cell
- Match Input Output Cell
- Recurrent Cell
- Memory Cell
- Different Memory Cell
- Kernel
- Convolution or Pool

Perceptron (P) | Feed Forward (FF) | Radial Basis Network (RBF) | Deep Feed Forward (DFF)

Recurrent Neural Network (RNN) | Long / Short Term Memory (LSTM) | Gated Recurrent Unit (GRU)

Auto Encoder (AE) | Variational AE (VAE) | Denoising AE (DAE) | Sparse AE (SAE)

Markov Chain (MC) | Hopfield Network (HN) | Boltzmann Machine (BM) | Restricted BM (RBM) | Deep Belief Network (DBN)

Deep Convolutional Network (DCN) | Deconvolutional Network (DN) | Deep Convolutional Inverse Graphics Network (DCIGN)

Generative Adversarial Network (GAN) | Liquid State Machine (LSM) | Extreme Learning Machine (ELM) | Echo State Network (ESN)

Deep Residual Network (DRN) | Kohonen Network (KN) | Support Vector Machine (SVM) | Neural Turing Machine (NTM)
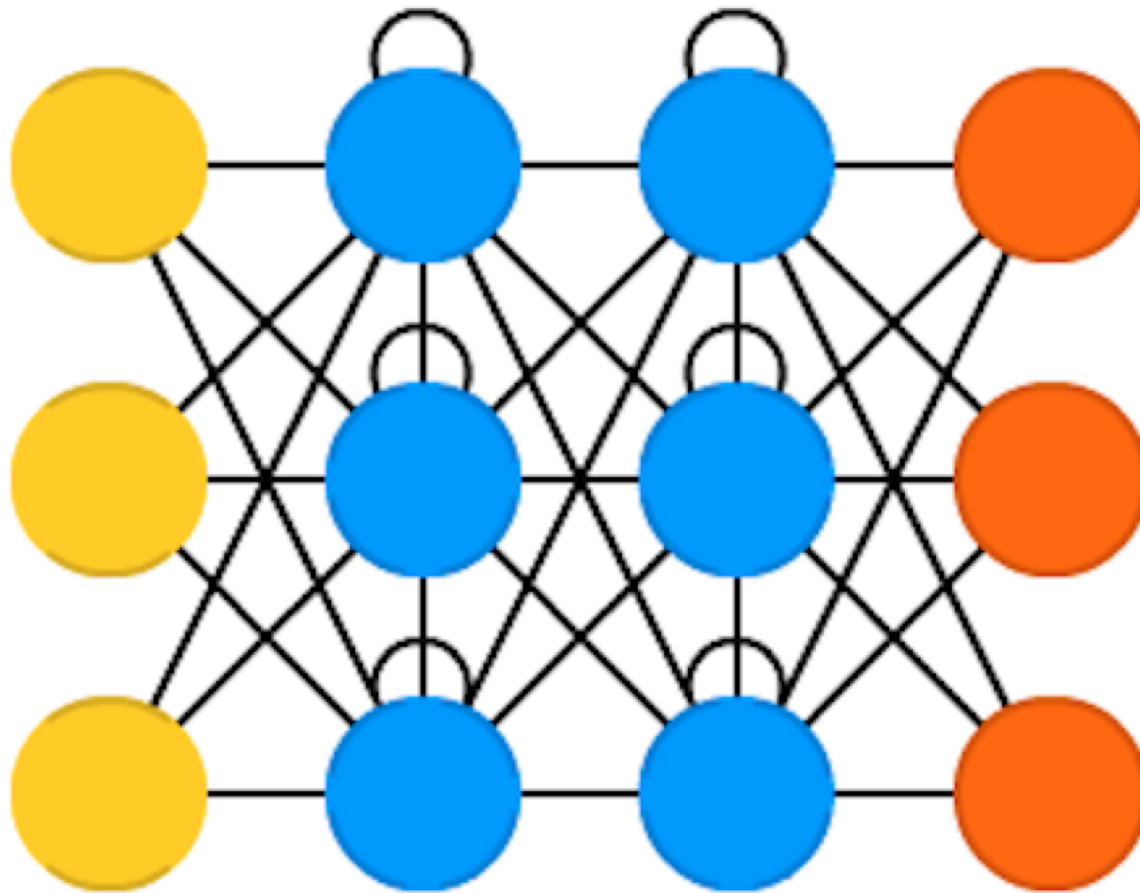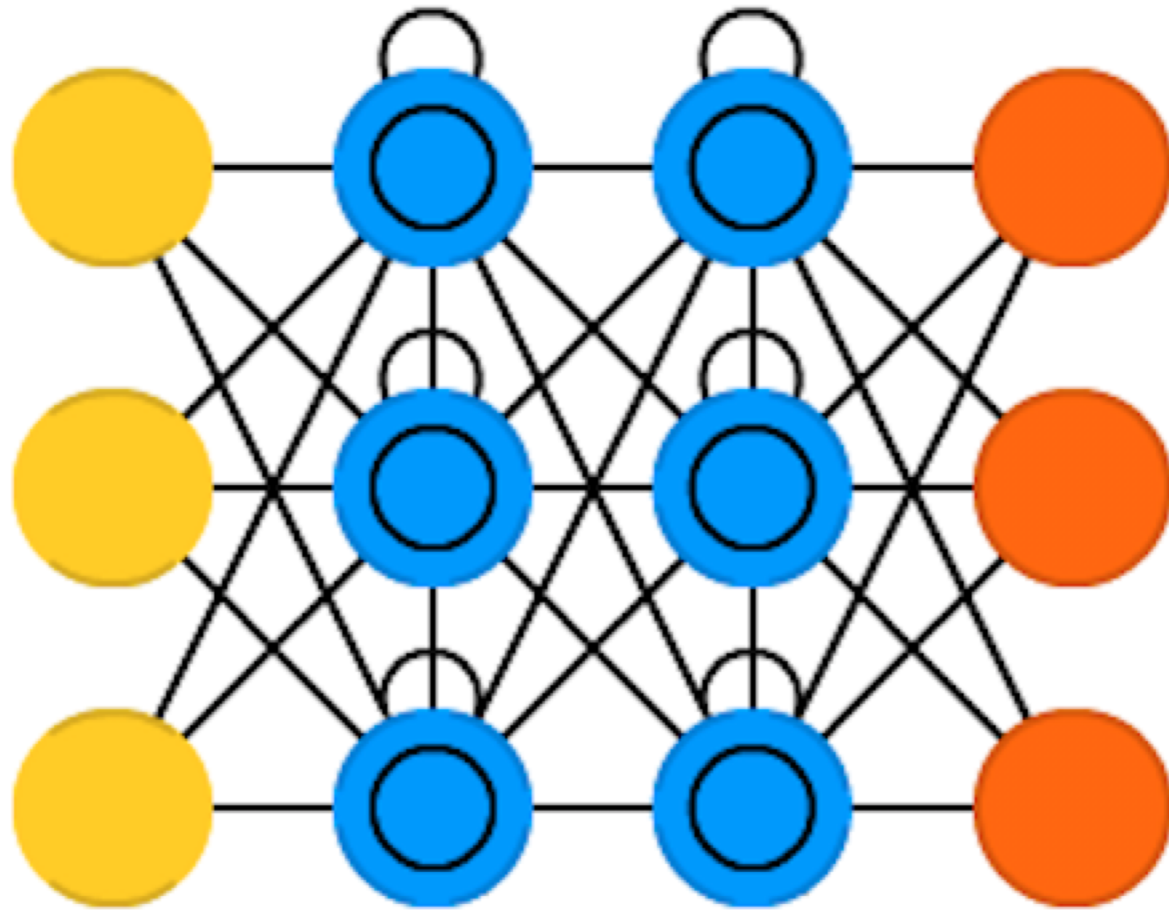
# Convolutional Neural Networks

## (CNN or Deep Convolutional Neural Networks, DCNN)

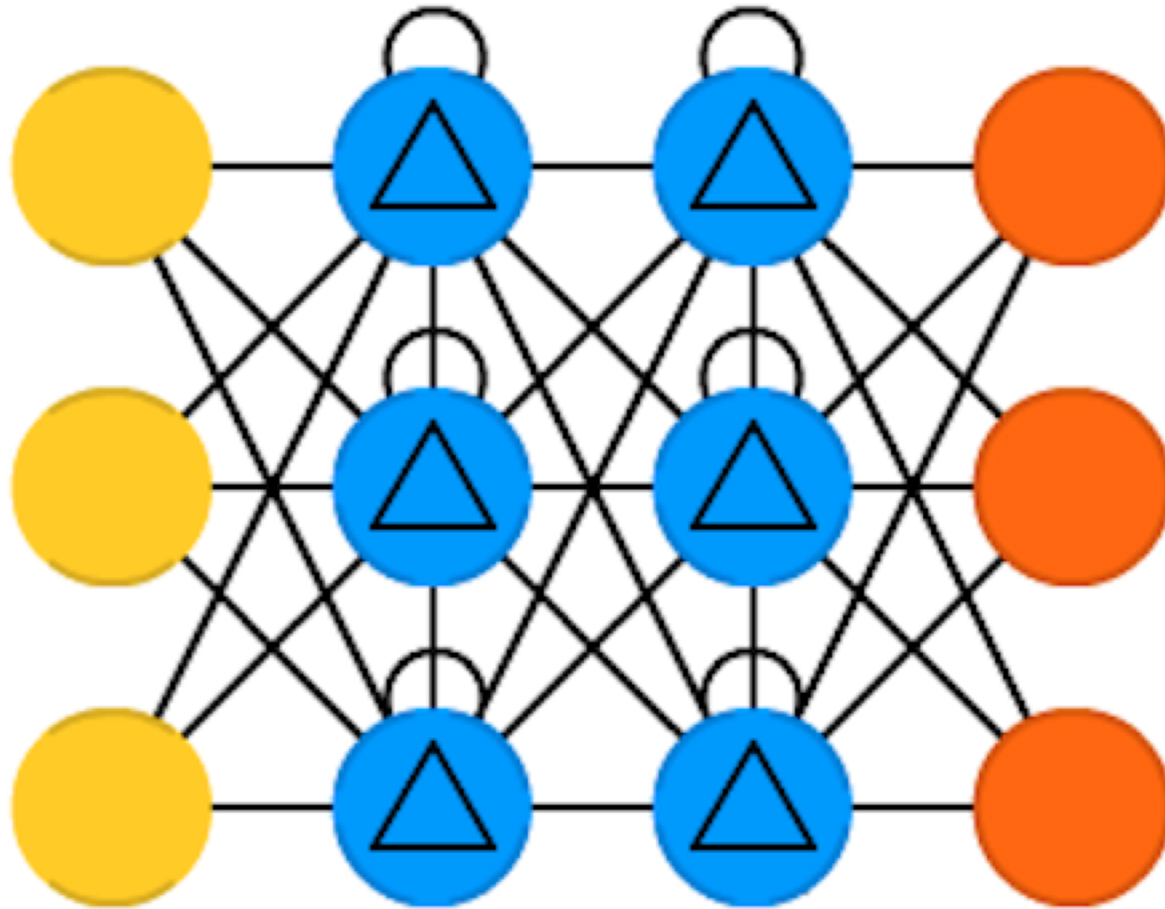# Recurrent Neural Networks (RNN)

# Long / Short Term Memory (LSTM)
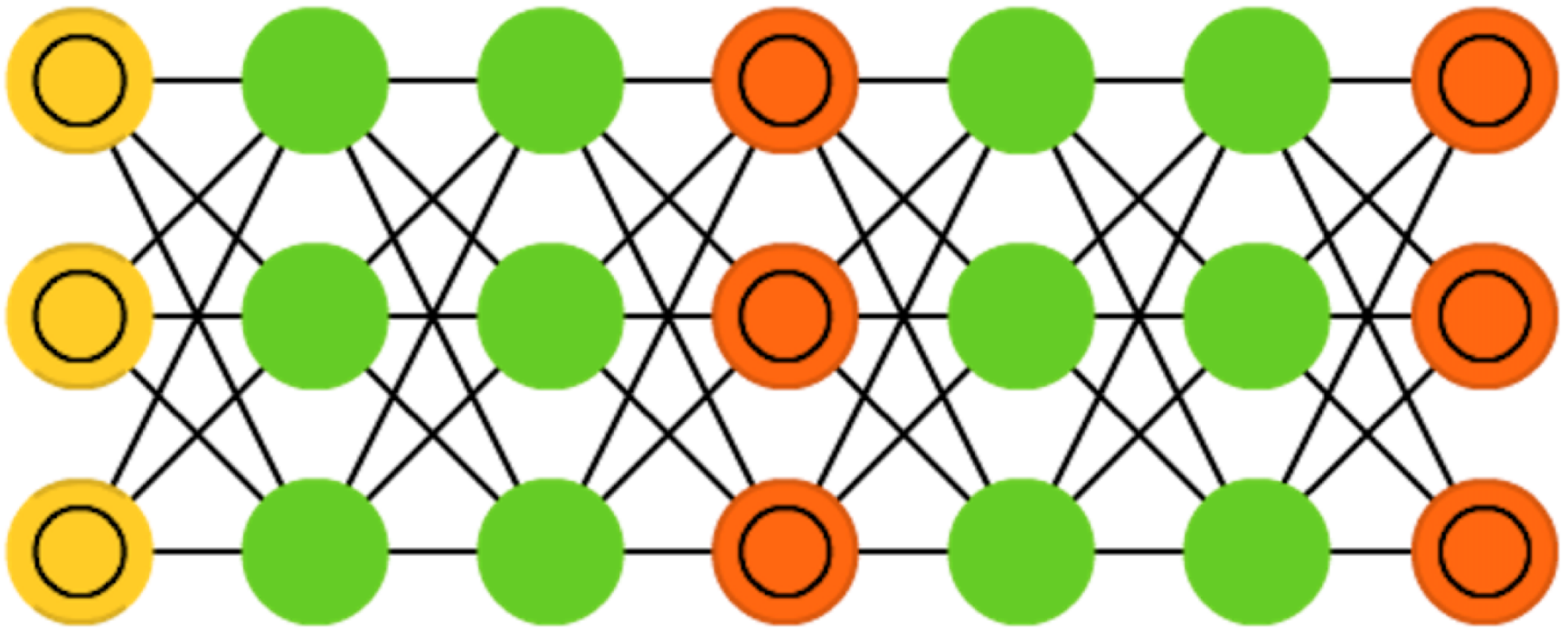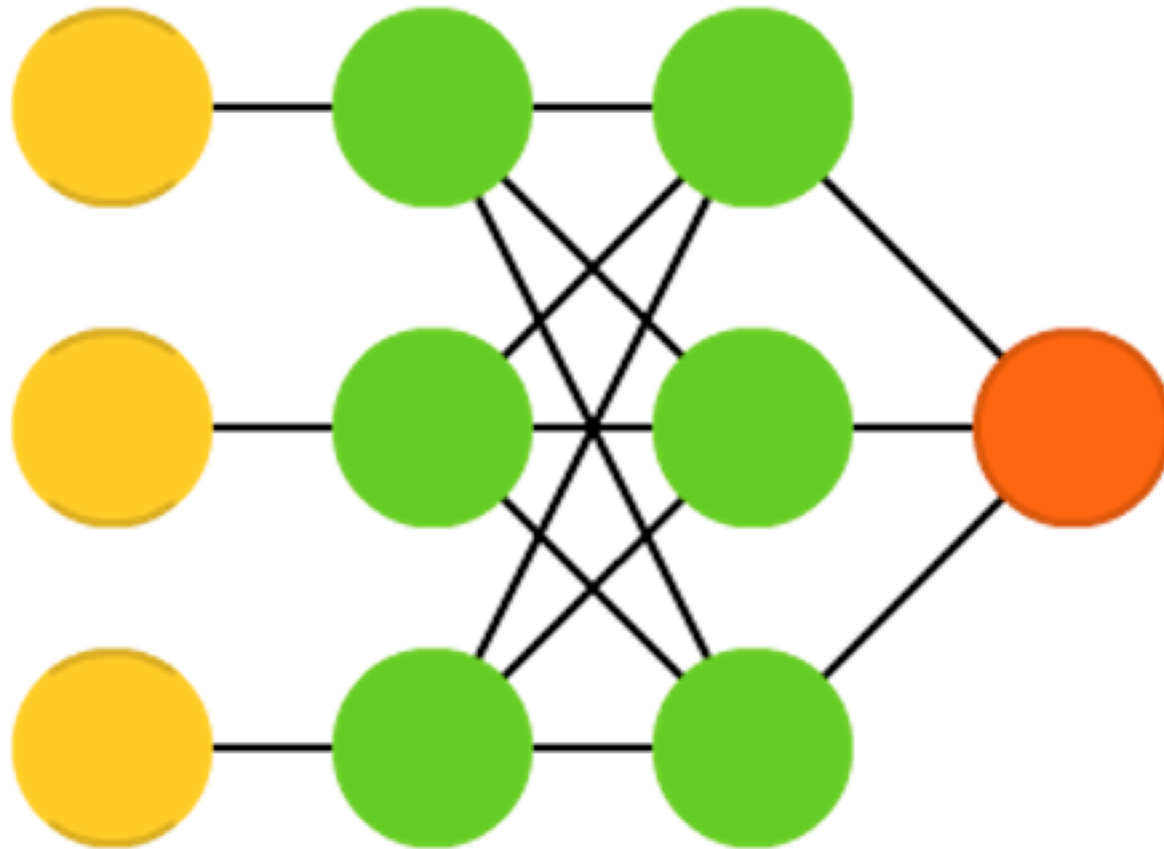


Hochreiter, Sepp, and Jürgen Schmidhuber. "Long short-term memory." Neural computation 9.8 (1997): 1735-1780.

# Gated Recurrent Units (GRU)

Chung, Junyoung, et al. "Empirical evaluation of gated recurrent neural networks on sequence modeling." arXiv preprint arXiv:1412.3555 (2014).
Source: http://www.asimovinstitute.org/neural-network-zoo/

# Generative Adversarial Networks (GAN)



Goodfellow, Ian, et al. "Generative adversarial nets." Advances in Neural Information Processing Systems. 2014.

# Support Vector Machines (SVM)



Cortes, Corinna, and Vladimir Vapnik. "Support-vector networks." Machine learning 20.3 (1995): 273-297.

# Architectures of Big Data Analytics

# Architecture of Big Data Analytics

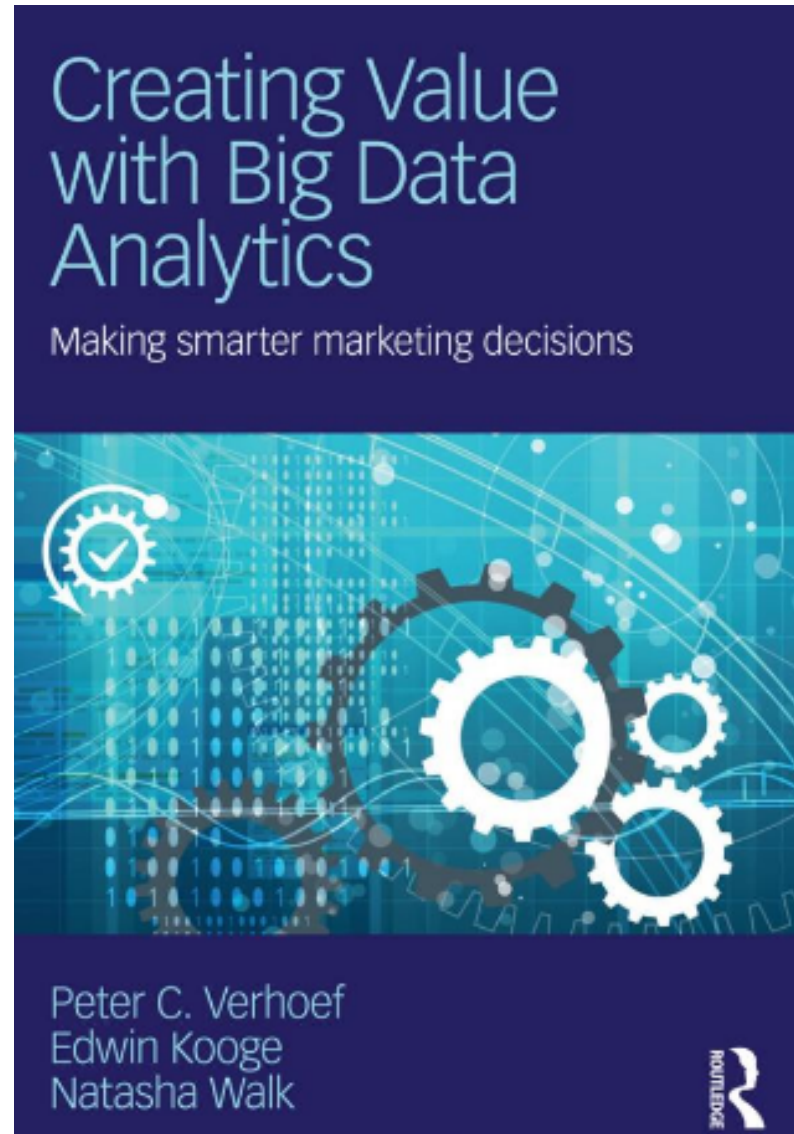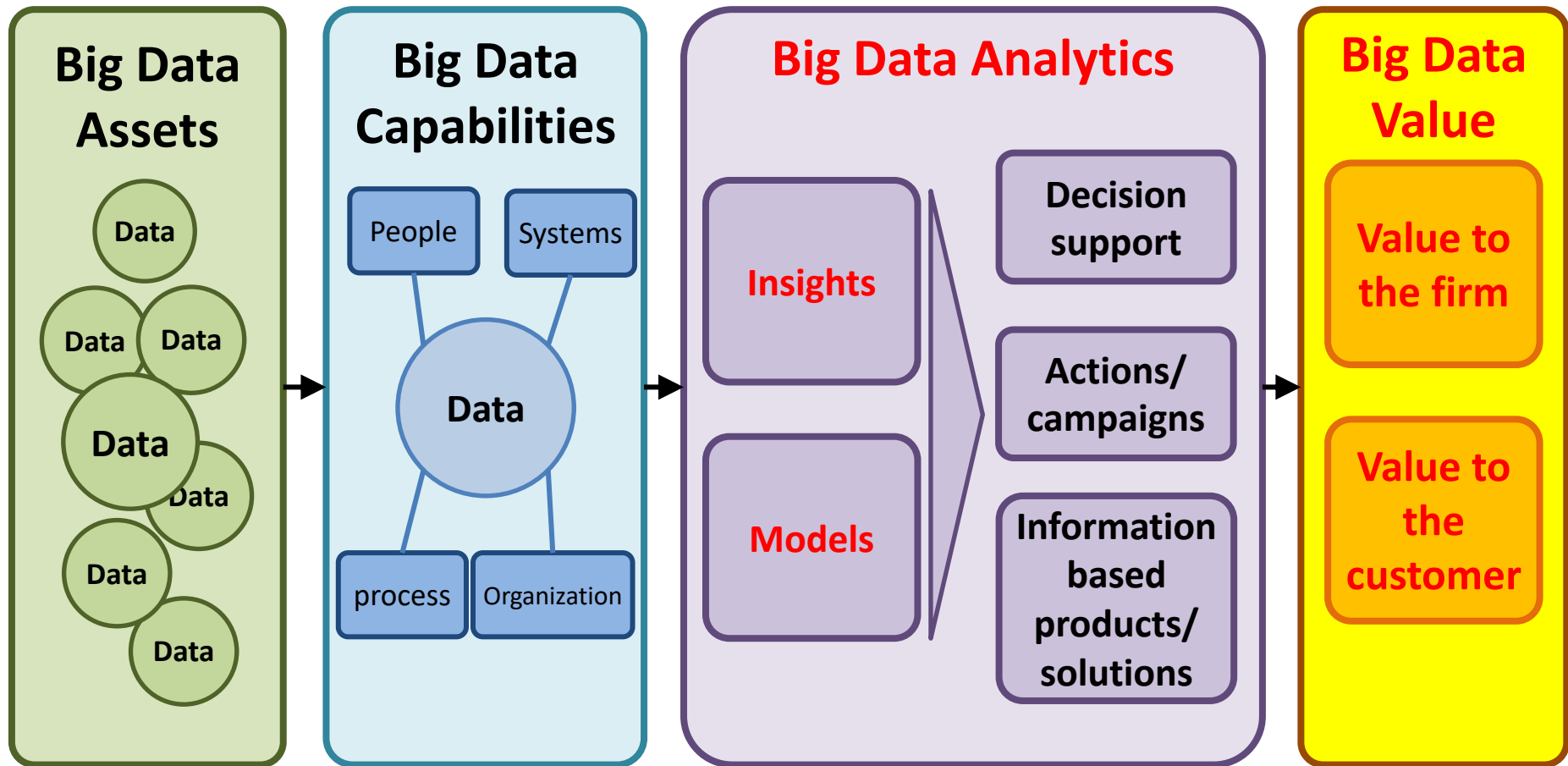**Big Data Sources**

* Internal

* External

* Multiple formats

* Multiple locations

* Multiple applications

**Raw Data** →

**Big Data Transformation**

Middleware

Extract Transform Load

Data Warehouse

Traditional Format CSV, Tables

**Transformed Data** →

**Big Data Platforms & Tools**

Hadoop
MapReduce
Pig
Hive
Jaql
Zookeeper
Hbase
Cassandra
Oozie
Avro
Mahout
Others

**Big Data Analytics** →

**Big Data Analytics Applications**

Queries

Reports

OLAP

**Data Mining**

# Architecture of Big Data Analytics

| Big Data Sources | Big Data Transformation | Big Data Platforms & Tools | Big Data Analytics Applications |
|---|---|---|---|

**Big Data Sources**

* Internal

* External

* Multiple formats

* Multiple locations

* Multiple applications

## Data Mining
## Big Data Analytics Applications

**Big Data Analytics Applications**

Queries

Reports

OLAP

**Data Mining**

# Creating Value with Big Data Analytics:
## Making Smarter Marketing Decisions,
## Peter C. Verhoef and Edwin Kooge, Routledge, 2016

# Big Data Value Creation Model

## Creating Value with Big Data Analytics: Making Smarter Marketing Decisions

# Digital Data Platform for Enterprises Big Data Analytics



**Enterprise Applications**

| Operational Benchmark | Customer focus | Organization Connections | Document Search | Sales Forecast |

Security (Authentication, Authorization, Auditing, Encryption, Protection)

| Variety of Sources | Ingestion layer | Processing Layer | Storage Layer | Analytics Layer | Visualization Apps |
|---|---|---|---|---|---|
| | Data Connectors | Data Mining | Hadoop | Traditional Analytics | |
| | Data Extraction | Data Enrichment | NoSQL | Search Based Analytics | |
| | CDC | Real-time Streaming | RDBMS | Predictive Analytics | |
| | Data Quality | Batch Processing | In-Memory | Ad-hoc Analytics | |

Data Governance and Monitoring (Workflow, lifecycle management, scheduler, manage)

**Digital Data Driven Platform for Enterprises**

# Architecture for Social Big Data Mining

(Hiroshi Ishikawa, 2015)
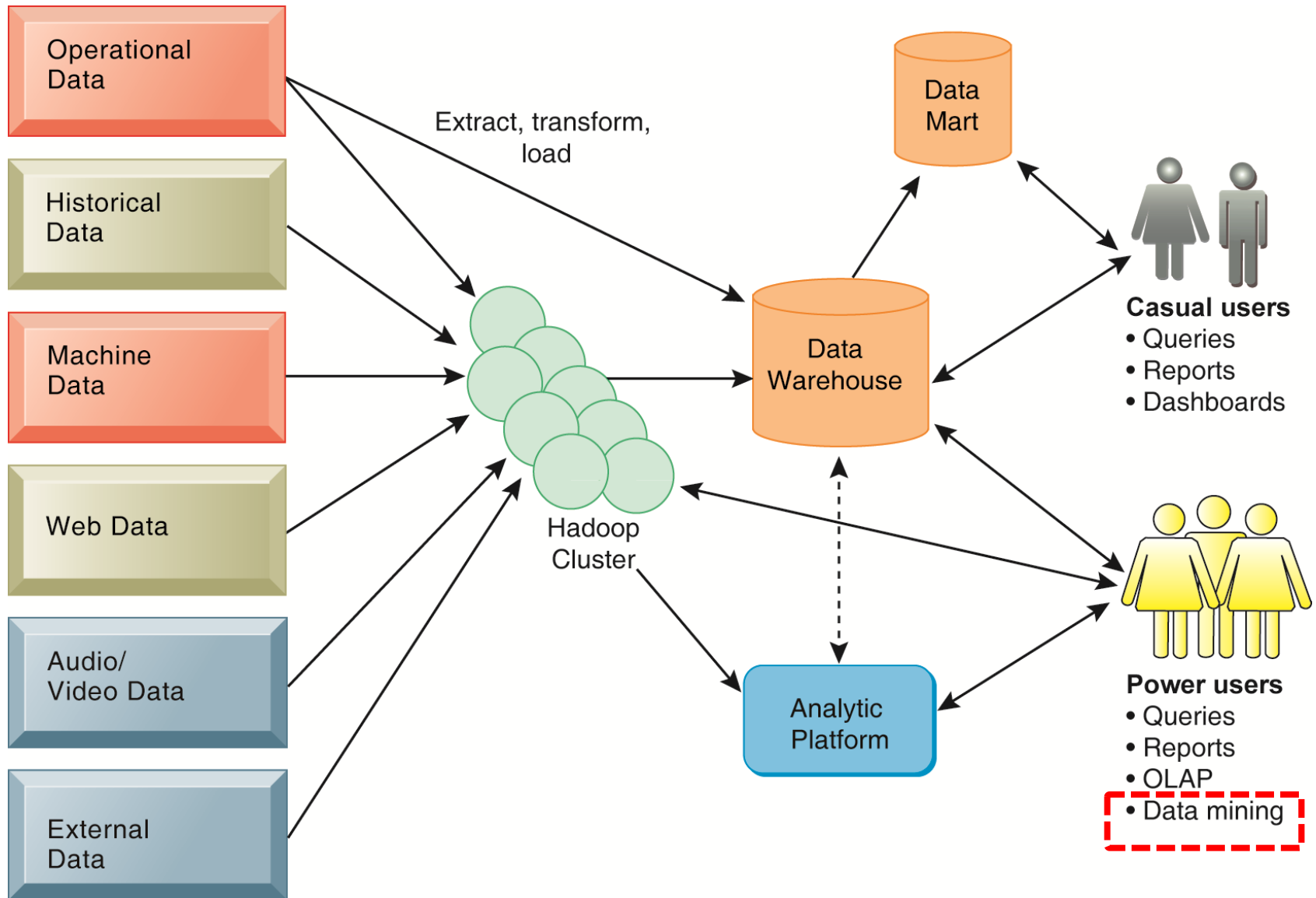
**Enabling Technologies**

- Integrated analysis model

- Natural Language Processing
- Information Extraction
- Anomaly Detection
- Discovery of relationships among heterogeneous data
- Large-scale visualization
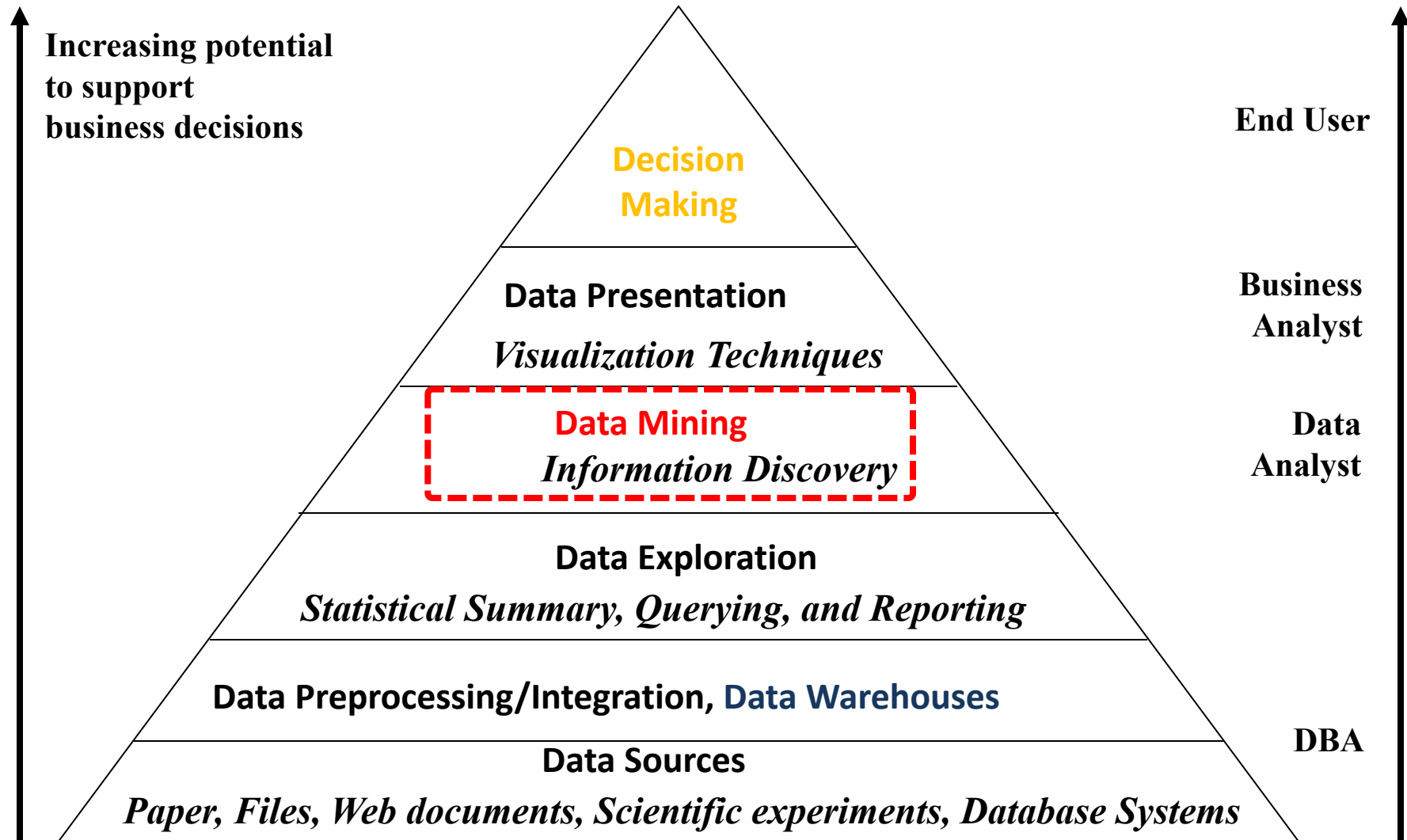
- Parallel distrusted processing

**Analysts**

- Model Construction
- Explanation by Model

- Construction and confirmation of individual hypothesis
- Description and execution of application-specific task

Integrated analysis

**Conceptual Layer**

Data Mining

Multivariate analysis

Application specific task

**Logical Layer**

Software

Hardware

**Social Data**

**Physical Layer**

# Business Intelligence (BI) Infrastructure

# Data Warehouse
# Data Mining and Business Intelligence



Increasing potential to support business decisions

**Decision Making**

End User

**Data Presentation**

*Visualization Techniques*

Business Analyst

**Data Mining**

*Information Discovery*

Data Analyst

**Data Exploration**

*Statistical Summary, Querying, and Reporting*

**Data Preprocessing/Integration, Data Warehouses**

**Data Sources**

*Paper, Files, Web documents, Scientific experiments, Database Systems*

DBA

# The Evolution of BI Capabilities

# Data Science and Business Intelligence

# Data Science and Business Intelligence



**Exploratory**

| Predictive Analytics and Data Mining (Data Science) | |
|---|---|
| Typical Techniques and Data Types | • Optimization, predictive modeling, forecasting, statistical analysis<br>• Structured/unstructured data, many types of sources, very large datasets |
| Common Questions | • What if…?<br>• What's the optimal scenario for our business?<br>• What will happen next? What if these trends continue? Why is this happening? |

## Predictive Analytics and Data Mining (Data Science)

Past · Time · Future

# Predictive Analytics and Data Mining (Data Science)

Structured/unstructured data, many types of sources, very large datasets

Optimization, predictive modeling, forecasting statistical analysis

What if…?
What's the optimal scenario for our business?
What will happen next?
What if these trends countinue?
Why is this happening?

# Data Mining

# AI

# Machine Learning

# Data Mining at the Intersection of Many Disciplines

# Data Mining

## Advanced Data Analysis

### Evolution of Database System Technology

# Evolution of Database System Technology

**Data Collection and Database Creation**

(1960s and earlier)

• Primitive file processing

↓

**Database Management Systems**

(1970s–early 1980s)

• Hierarchical and network database systems
• Relational database systems
• Query languages: SQL, etc.
• Transactions, concurrency control and recovery
• On-line transaction processing (OLTP)

**Advanced Database Systems**

(mid-1980s–present)

• Advanced data models: extended  relational, object-relational, etc.
• Advanced applications: spatial, temporal, multimedia, active, stream and sensor, scientific and engineering, knowledge-based
• XML-based database systems
• Integration with information retrieval
• Data and information integration

**Advanced Data Analysis:**

(late 1980s–present)

• Data warehouse and OLAP
• Data mining and knowledge discovery: generalization, classification, association, clustering
• Advanced data mining applications: stream data mining, bio-data mining, time-series analysis, text mining, Web mining, intrusion detection, etc.
• Data mining applications
• Data mining and society

**New Generation of Information Systems**
(present–future)

51

# Big Data Analysis

- **Too Big,
  too Unstructured,
  too many different source** to be manageable through traditional databases

# Internet Evolution
## Internet of People (IoP): Social Media
## Internet of Things (IoT): Machine to Machine

# Data Mining Technologies

# A Taxonomy for Data Mining Tasks

| Data Mining | Learning Method | Popular Algorithms |
|---|---|---|
| Prediction | Supervised | Classification and Regression Trees, ANN, SVM, Genetic Algorithms |
| Classification | Supervised | Decision trees, ANN/MLP, SVM, Rough sets, Genetic Algorithms |
| Regression | Supervised | Linear/Nonlinear Regression, Regression trees, ANN/MLP, SVM |
| Association | Unsupervised | Apriory, OneR, ZeroR, Eclat |
| Link analysis | Unsupervised | Expectation Maximization, Apriory Algorithm, Graph-based Matching |
| Sequence analysis | Unsupervised | Apriory Algorithm, FP-Growth technique |
| Clustering | Unsupervised | K-means, ANN/SOM |
| Outlier analysis | Unsupervised | K-means, Expectation Maximization (EM) |

# Traditional Analytics

**BI and Analytics**

Operational Data Sources

EDW

Mart

Data Mart

Analytic

Analytic Mart

Unstructured, Semi-structured and Streaming data (i.e. sensor data) handled often outside the Warehouse flow

# Hadoop as a "new data" Store

# Hadoop as an additional input to the EDW

# Hadoop Data Platform As a "staging Layer" as part of a "data Lake"

## – Downstream stores could be Hadoop, data appliances or an RDBMS

# SAS Big data Strategy – SAS areas

# SAS Big data Strategy – SAS areas

# SAS® Within the HADOOP ECOSYSTEM

# SAS enables the entire lifecycle around HADOOP



SAS enableS the entire lifecycle around HADOOP

SAS Visual Analytics
Decision Manager

SAS Scoring Accelerator for Hadoop
SAS Code Accelerator for Hadoop

Decision Manager

IDENTIFY / FORMULATE PROBLEM

DATA PREPARATION

DATA EXPLORATION

TRANSFORM & SELECT

BUILD MODEL

VALIDATE MODEL

DEPLOY MODEL

EVALUATE / MONITOR RESULTS

Done using either the Data Preparation, Data Exploration or Build Model Tools

SAS Visual Analytics
SAS Visual Statistics
SAS In-Memory Statistics for Hadoop

Done using either the Data Preparation, Data Exploration or Build Model Tools

SAS High Performance Analytics Offerings supported by relevant clients like SAS Enterprise Miner, SAS/STAT etc.

# Big Data, Big Analytics:

**Emerging Business Intelligence and Analytic Trends for Today's Businesses**

# Big Data, Prediction vs. Explanation

Source: Agarwal, R., & Dhar, V. (2014). Editorial—Big Data, Data Science, and Analytics: The Opportunity and Challenge for IS Research. Information Systems Research, 25(3), 443-448.

# Big Data:
# The Management Revolution

Source: McAfee, A., & Brynjolfsson, E. (2012). Big data: the management revolution.*Harvard business review*.

# Business Intelligence and Enterprise Analytics

- Predictive analytics

- Data mining

- Business analytics

- Web analytics

- <span style="color:red">Big-data</span> analytics

# Three Types of Business Analytics

- Prescriptive Analytics

- Predictive Analytics

- Descriptive Analytics

# Three Types of Business Analytics



| | | |
|---|---|---|
| Optimization | "What's the best that can happen?" | **Prescriptive Analytics** |
| Randomized Testing | "What if we try this?" | |
| Predictive Modeling / Forecasting | "What will happen next?" | **Predictive Analytics** |
| Statistical Modeling | "Why is this happening?" | |
| Alerts | "What actions are needed?" | **Descriptive Analytics** |
| Query / Drill Down | "What exactly is the problem?" | |
| Ad hoc Reports / Scorecards | "How many, how often, where?" | |
| Standard Report | "What happened?" | |

# Big Data



Mobile Sensors · Social Media · Video Surveillance · Video Rendering · Smart Grids · Geophysical Exploration · Medical Imaging · Gene Sequencing

# Big Data Growth is increasingly unstructured

# Typical Analytic Architecture

# Data Evolution and the Rise of Big Data Sources

# Emerging Big Data Ecosystem

# Key Roles for the New Big Data Ecosystem



| Role |
|------|
| Deep Analytical Talent — Data Scientists Projected U.S. talent gap: 140,000 to 190,000 |
| Data Savvy Professionals — Projected U.S. talent gap: 1.5 million |
| Technology and Data Enablers |

Note: Figures above reflect a projected talent gap in US in 2018, as shown in McKinsey May 2011 article "Big Data: The Next Frontier for Innovation, Competition, and Productivity"

# Profile of a Data Scientist

- **Quantitative**
  - mathematics or statistics
- **Technical**
  - software engineering, machine learning, and programming skills
- **Skeptical mind-set** and **critical thinking**
- **Curious** and **creative**
- **Communicative** and **collaborative**

# Data Scientist Profile

# Big Data Analytics Lifecycle

# Key Roles for a Successful Analytics Project

# Overview of Data Analytics Lifecycle

# Overview of Data Analytics Lifecycle

1. Discovery

2. Data preparation

3. Model planning

4. Model building

5. Communicate results

6. Operationalize

# Key Outputs from a Successful Analytics Project

# Data Mining Process

# Data Mining Process

- A manifestation of best practices

- A systematic way to conduct DM projects

- Different groups has different versions

- Most common standard processes:

  - CRISP-DM
    (Cross-Industry Standard Process for Data Mining)

  - SEMMA
    (Sample, Explore, Modify, Model, and Assess)

  - KDD
    (Knowledge Discovery in Databases)

# Data Mining Process (SOP of DM)

What main methodology are you using for your **analytics**, **data mining**, or **data science** projects ?

# Data Mining Process



| | 2014 poll | 2007 poll |
|---|---|---|
| CRISP-DM (86) | 43% | 42% |
| My own (55) | 27.5% | 19% |
| SEMMA (17) | 8.5% | 13% |
| Other, not domain-specific (16) | 8% | 4% |
| KDD Process (15) | 7.5% | 7.3% |
| My organizations' (7) | 3.5% | 5.3% |
| A domain-specific methodology (4) | 2% | 4.7% |
| None (0) | 0% | 4.7% |

2014 poll    2007 poll

# Data Mining:

## Core Analytics Process

## The KDD Process for Extracting Useful Knowledge from Volumes of Data

Source: Fayyad, U., Piatetsky-Shapiro, G., & Smyth, P. (1996). The KDD Process for Extracting Useful Knowledge from Volumes of Data. Communications of the ACM, 39(11), 27-34.

Fayyad, U., Piatetsky-Shapiro, G., & Smyth, P. (1996).
**The KDD Process for**
**Extracting Useful Knowledge**
**from Volumes of Data.**
Communications of the ACM, 39(11), 27-34.

# Data Mining
## Knowledge Discovery in Databases (KDD) Process
(Fayyad et al., 1996)

# Knowledge Discovery (KDD) Process



**Data mining:**
**core of knowledge discovery process**

Pattern Evaluation

Knowledge

Data Mining

Task-relevant Data

Data Warehouse    Selection

Data Cleaning

Data Integration

Databases

Source: Han & Kamber (2006)

# Data Mining Process: CRISP-DM

# Data Mining Process: CRISP-DM

Step 1: Business Understanding

Step 2: Data Understanding

Step 3: Data Preparation (!)

Accounts for ~85% of total project time

Step 4: Model Building

Step 5: Testing and Evaluation

Step 6: Deployment

- The process is highly repetitive and experimental (DM: art versus science?)

# Data Preparation – A Critical DM Task



**Real-world Data**
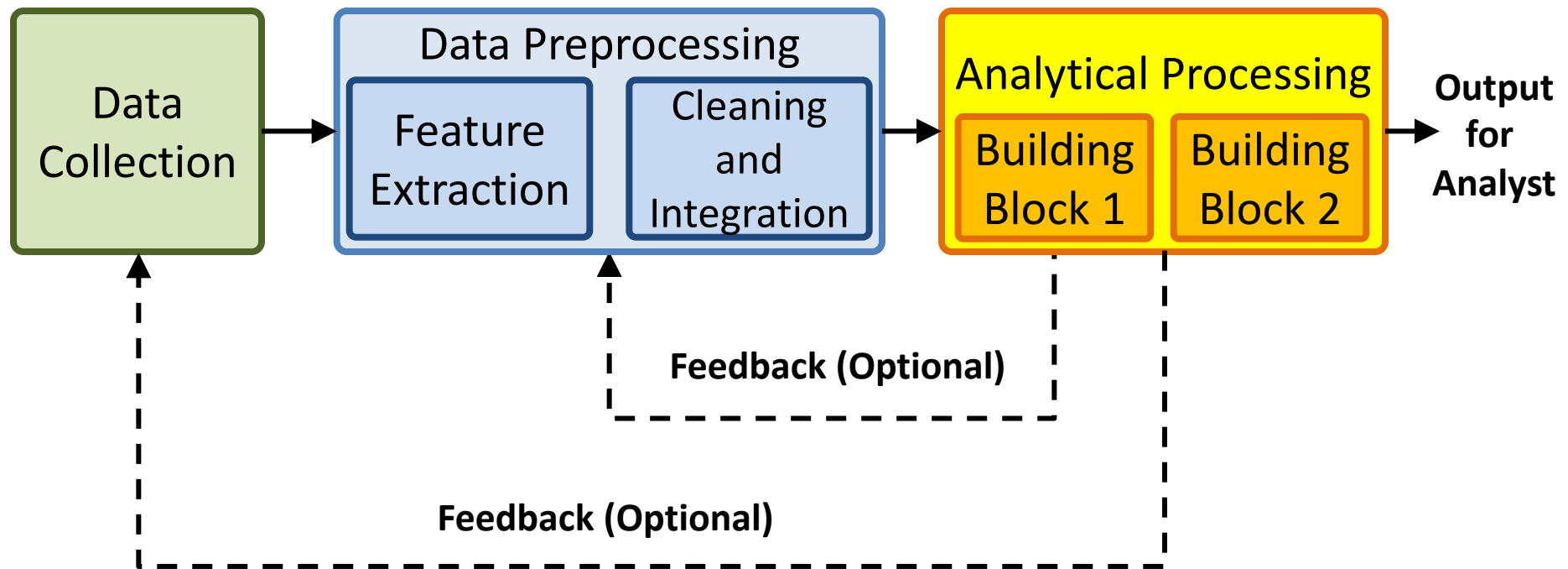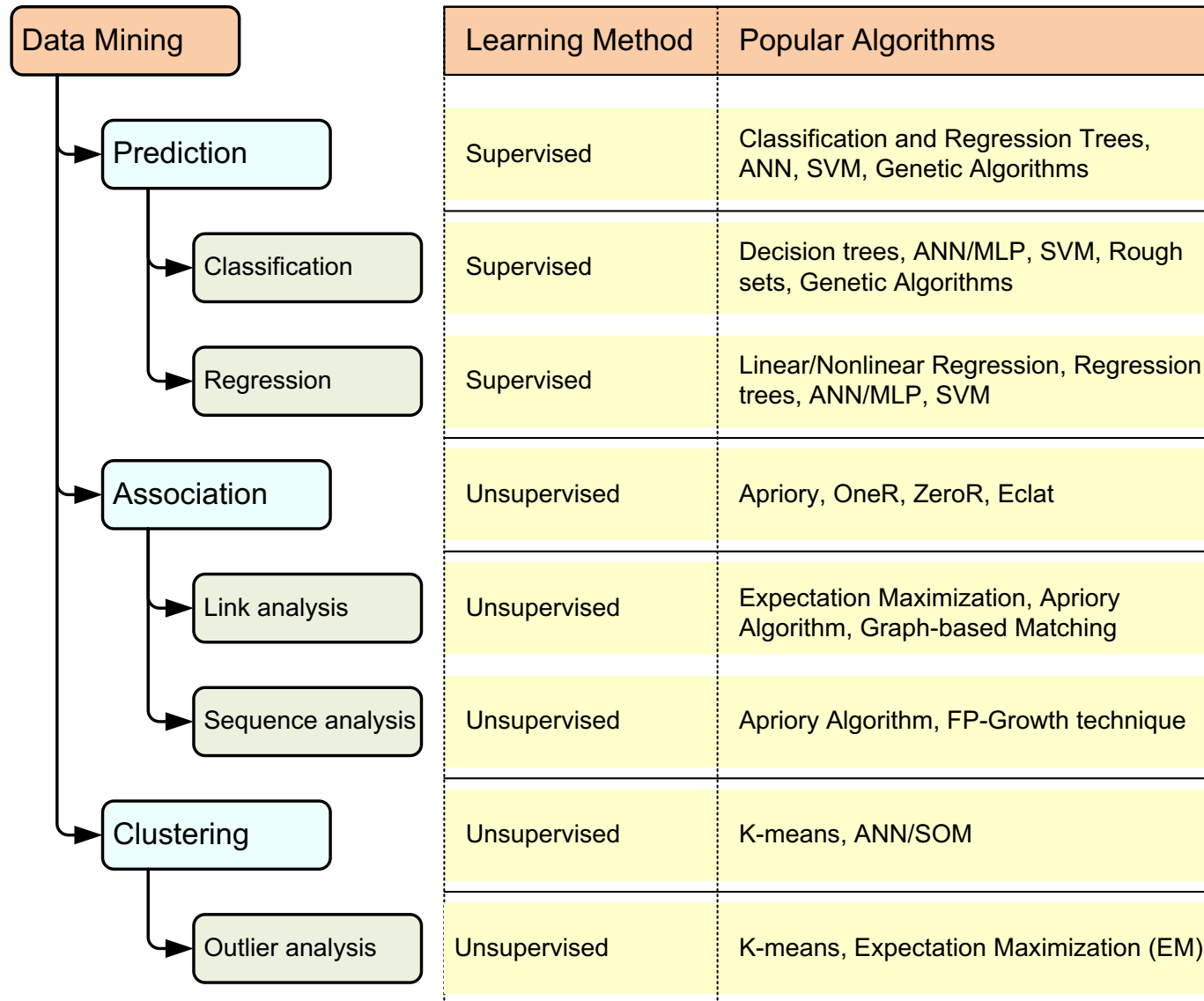
**Data Consolidation**
- Collect data
- Select data
- Integrate data

**Data Cleaning**
- Impute missing values
- Reduce noise in data
- Eliminate inconsistencies

**Data Transformation**
- Normalize data
- Discretize/aggregate data
- Construct new attributes

**Data Reduction**
- Reduce number of variables
- Reduce number of cases
- Balance skewed data

**Well-formed Data**

# Data Mining Process: SEMMA



**Sample** (Generate a representative sample of the data)

**Explore** (Visualization and basic description of the data)

**Modify** (Select variables, transform variable representations)

**Model** (Use variety of statistical and machine learning models )

**Assess** (Evaluate the accuracy and usefulness of the models)

SEMMA

# Data Mining Processing Pipeline
## (Charu Aggarwal, 2015)

# A Taxonomy for Data Mining Tasks

| | Learning Method | Popular Algorithms |
|---|---|---|
| Data Mining | | |
| Prediction | Supervised | Classification and Regression Trees, ANN, SVM, Genetic Algorithms |
| Classification | Supervised | Decision trees, ANN/MLP, SVM, Rough sets, Genetic Algorithms |
| Regression | Supervised | Linear/Nonlinear Regression, Regression trees, ANN/MLP, SVM |
| Association | Unsupervised | Apriory, OneR, ZeroR, Eclat |
| Link analysis | Unsupervised | Expectation Maximization, Apriory Algorithm, Graph-based Matching |
| Sequence analysis | Unsupervised | Apriory Algorithm, FP-Growth technique |
| Clustering | Unsupervised | K-means, ANN/SOM |
| Outlier analysis | Unsupervised | K-means, Expectation Maximization (EM) |

# **Fundamental Big Data**: **MapReduce Paradigm, Hadoop and Spark Ecosystem**

National Security · Cyber security · Maritime security · Smarter Transport · ...

**VISUAL ANALYTICS**

**DYNAMIC & INTERACTIVE**
Dashboard  Graph
Map

**ENHANCE**
Understanding  Investigation
User Experience

**BIG ANALYTICS**

**QUERY & FILTER**
Complex queries
$R^2I^2$

**DETECT**
Anomalies
Communities
Typologies

**PREDICT**
Tending
Real-time
Prediction

**DECIDE**
Simulation
Optimization

**BIG DATA – Batch**

**BIG DATA – Real Time**

**DATA**

Complex by nature

Complex by structure

# MapReduce Paradigm

# MapReduce Paradigm



Map

Reduce

Big Data

Map0　Map1　Map2　Map3

Reduce0　Reduce1　Reduce2　Reduce3

MapReduce Data

Output Data

# MapReduce Word Count

**Input**

Dog Love Cat
Bird Love Bird
Dog Bird Cat

# MapReduce Word Count

**Input** → **Output**

Dog Love Cat
Bird Love Bird
Dog Bird Cat

Bird, 3
Cat, 2
Dog, 2
Love, 2

# MapReduce Word Count

Input  Split  Map  Shuffle  Reduce  Output

Dog Love Cat
Bird Love Bird
Dog Bird Cat

Dog Love Cat

Bird Love Bird

Dog Bird Cat

Dog, 1
Love, 1
Cat, 1

Bird, 1
Love, 1
Bird, 1

Dog, 1
Bird, 1
Cat, 1

Bird, (1, 1, 1)

Cat, (1, 1)

Dog, (1, 1)

Love, (1, 1)

Bird, 3

Cat, 2

Dog, 2

Love, 2

Bird, 3
Cat, 2
Dog, 2
Love, 2

# Hadoop Ecosystem

The **Apache™ Hadoop®** project develops **open-source software** for reliable, scalable, **distributed computing**.

**MapReduce** **Processing**

+

**HDFS** **Storage**

# Big Data with Hadoop Architecture

# Big Data with Hadoop Architecture
## Logical Architecture
### Processing: MapReduce

# Big Data with Hadoop Architecture
## Logical Architecture
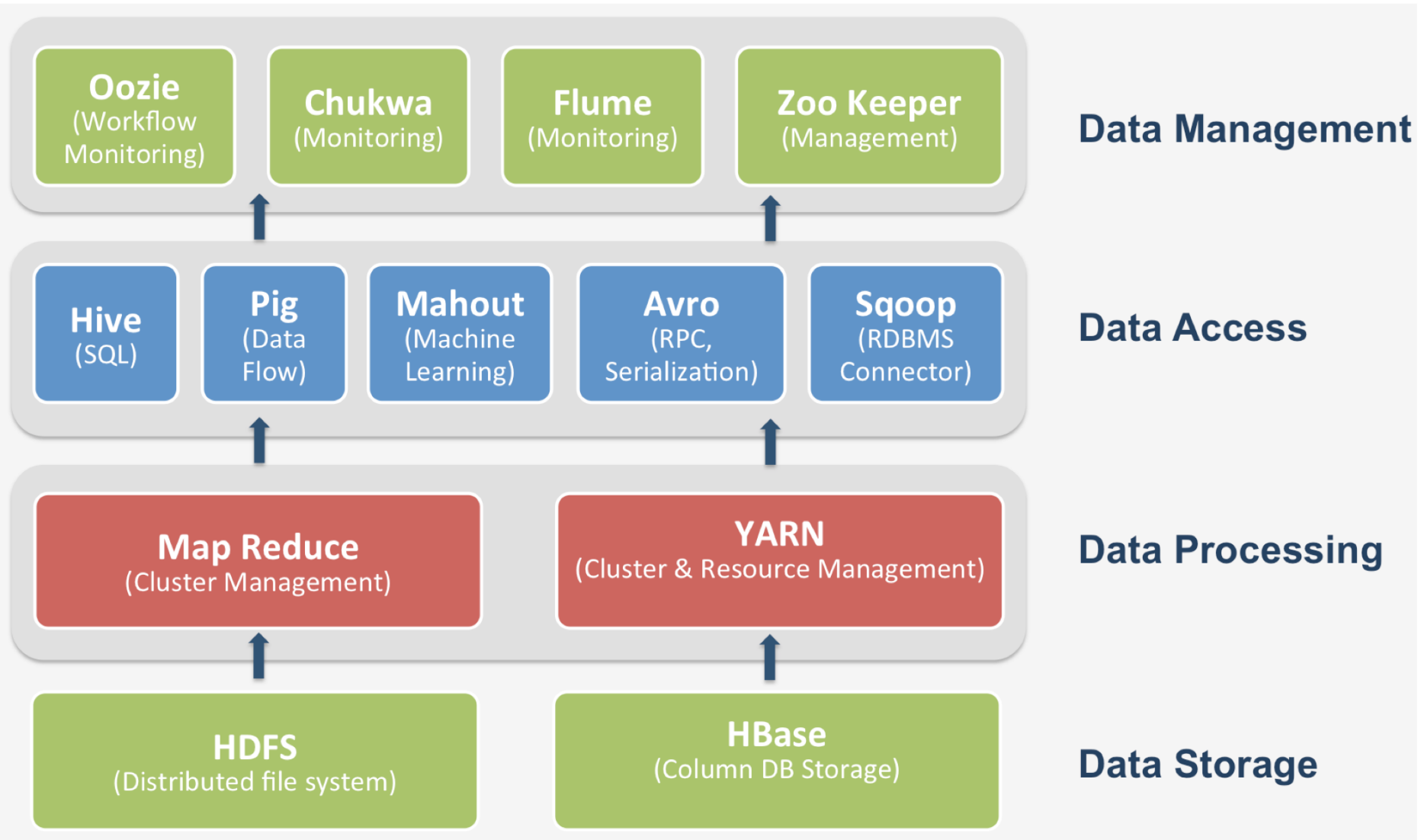### Storage: HDFS

# Big Data with Hadoop Architecture Process Flow
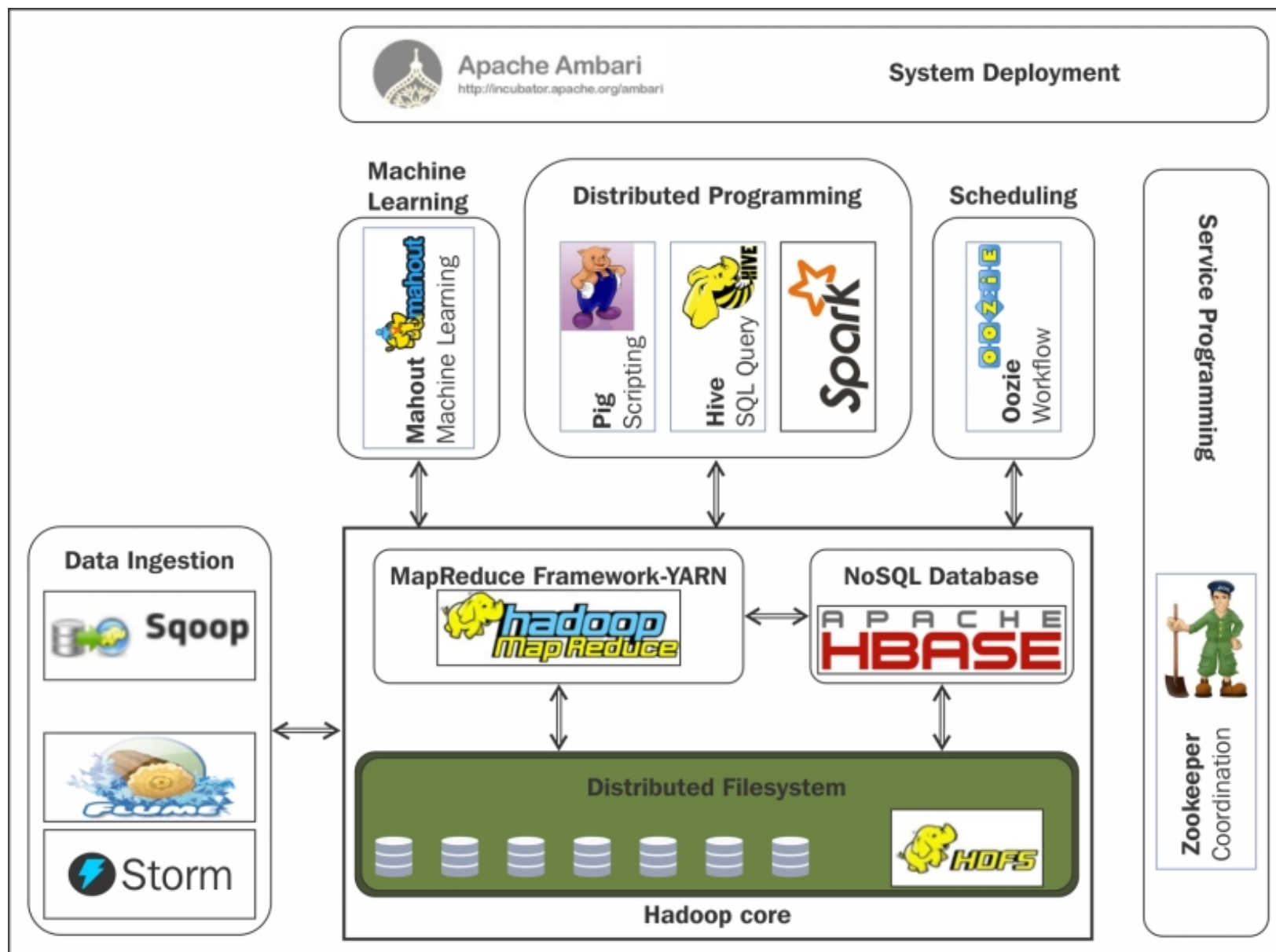
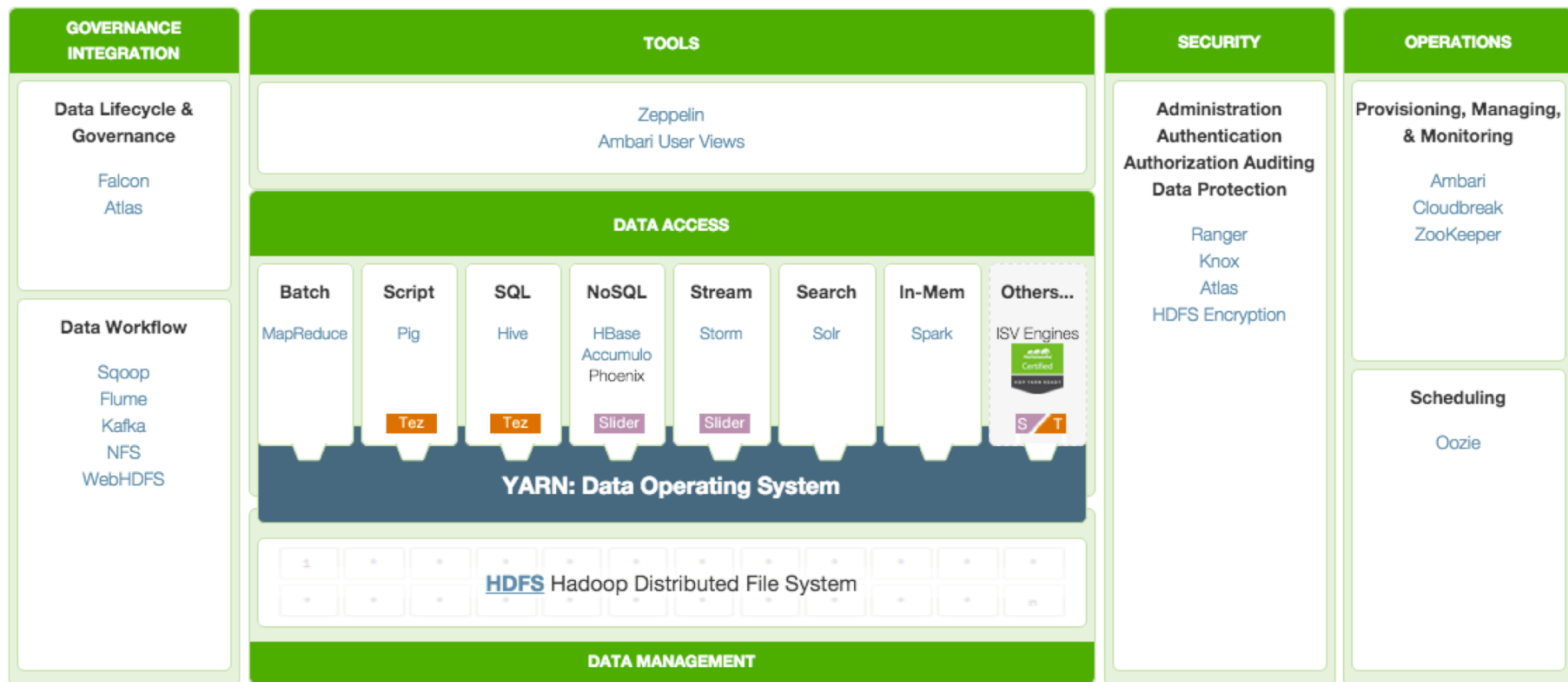# Big Data with Hadoop Architecture
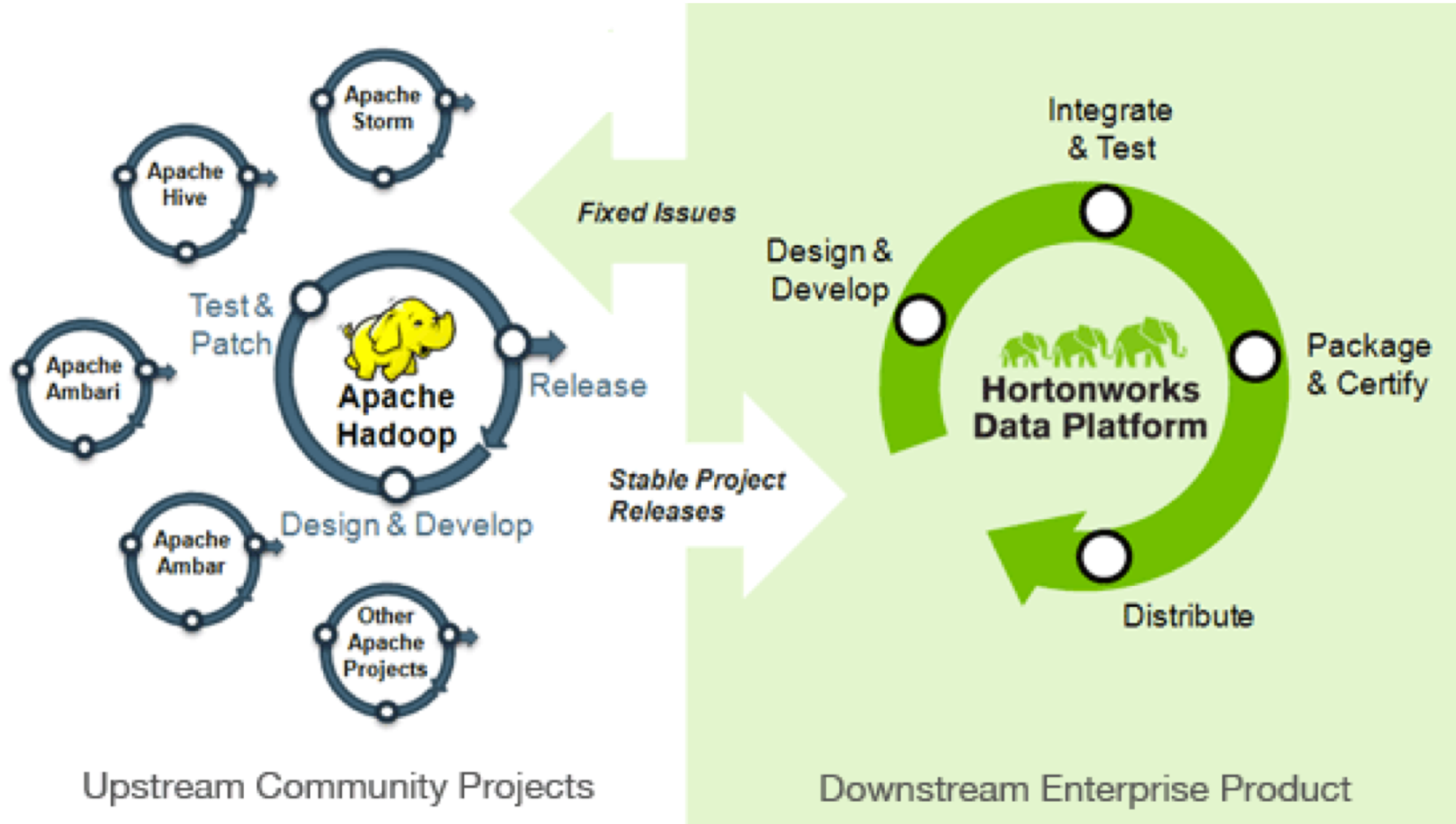## Hadoop Cluster

# Hadoop Ecosystem

| | | | | |
|---|---|---|---|---|
| **Oozie** (Workflow Monitoring) | **Chukwa** (Monitoring) | **Flume** (Monitoring) | **Zoo Keeper** (Management) | **Data Management** |

| | | | | | |
|---|---|---|---|---|---|
| **Hive** (SQL) | **Pig** (Data Flow) | **Mahout** (Machine Learning) | **Avro** (RPC, Serialization) | **Sqoop** (RDBMS Connector) | **Data Access** |

| | | |
|---|---|---|
| **Map Reduce** (Cluster Management) | **YARN** (Cluster & Resource Management) | **Data Processing** |

| | | |
|---|---|---|
| **HDFS** (Distributed file system) | **HBase** (Column DB Storage) | **Data Storage** |

# Hadoop Ecosystem

# HDP (Hortonworks Data Platform)
## A Complete Enterprise Hadoop Data Platform

# Apache Hadoop
# Hortonworks Data Platform



Upstream Community Projects

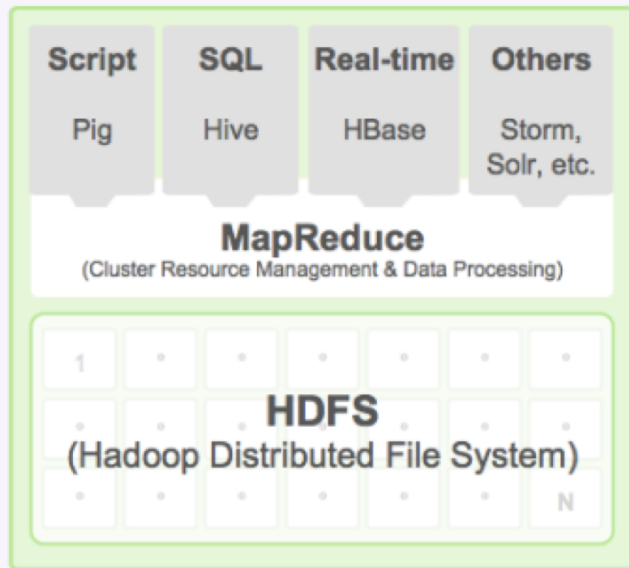Downstream Enterprise Product

# Hadoop and Data Analytics Tools
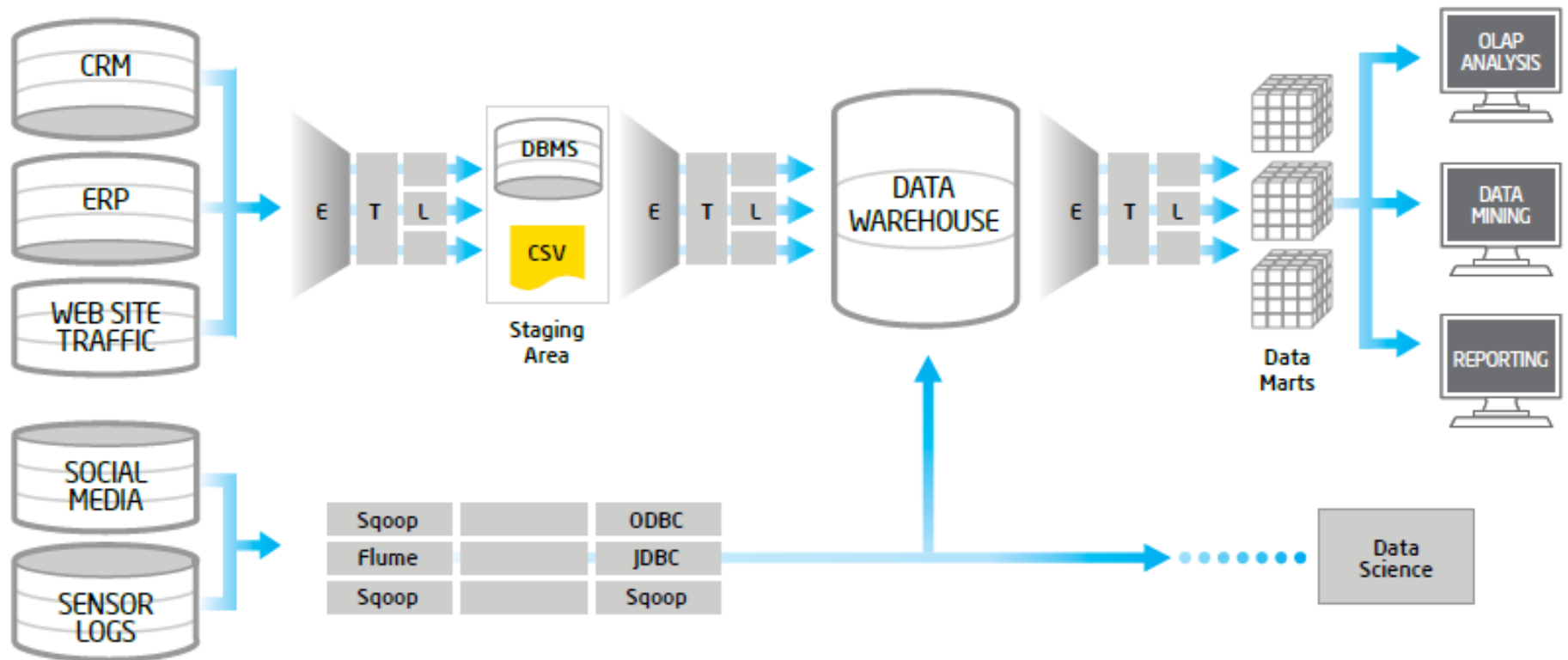
# Hadoop 1 → Hadoop 2

# Big Data Solution

# Traditional ETL Architecture

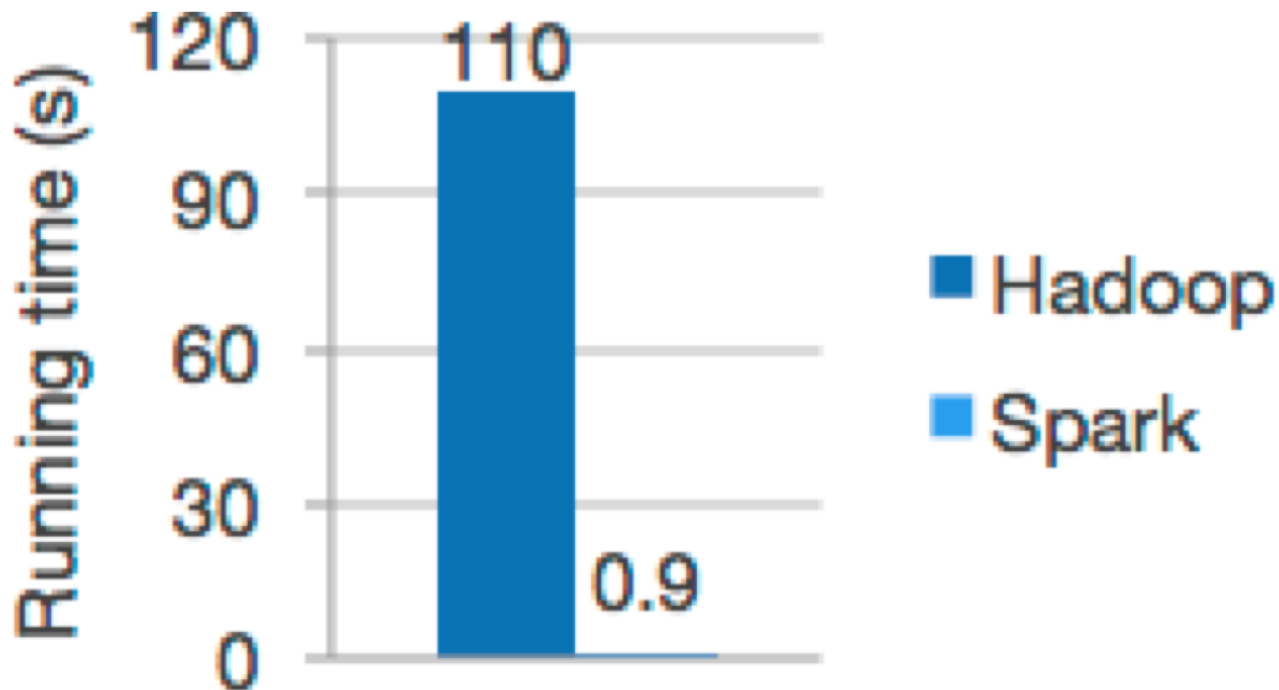# Offload ETL with Hadoop (Big Data Architecture)

# Spark Ecosystem

*Lightning-fast cluster computing*

# Apache Spark

## is a fast and general engine for

## large-scale data processing.

# Logistic regression in Hadoop and Spark



Run programs up to 100x faster than Hadoop MapReduce in memory, or 10x faster on disk.

# Ease of Use

- Write applications quickly in Java, Scala, Python, R.

# Word count in Spark's Python API

```python
text_file = spark.textFile("hdfs://...")


text_file.flatMap(lambda line: line.split())
    .map(lambda word: (word, 1))
    .reduceByKey(lambda a, b: a+b)
```

# Spark and Hadoop

# Spark Ecosystem

| Spark SQL | Spark Streaming | MLlib (machine learning) | GraphX (graph) |
|-----------|-----------------|--------------------------|----------------|

## Apache Spark

# Spark Ecosystem



Spark SQL + DataFrames | Streaming | MLlib *Machine Learning* | GraphX *Graph Computation*

Spark Core API

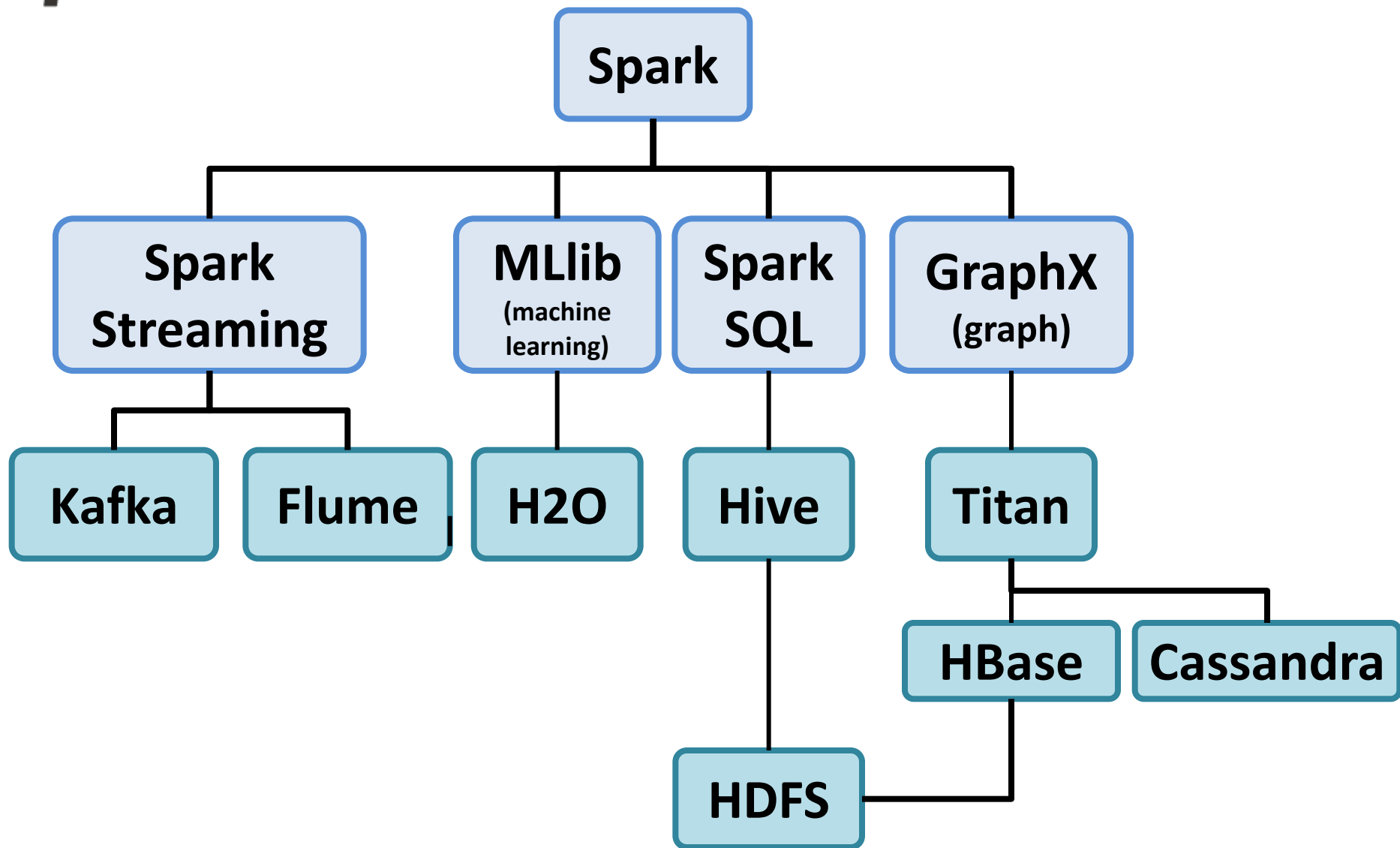R | SQL | Python | Scala | Java

Source: https://databricks.com/spark/about

# Spark Ecosystem

# SMACK Stack

- **Spark**
  - fast and general engine for distributed, large-scale data processing

- **Mesos**
  - cluster resource management system that provides efficient resource isolation and sharing across distributed applications

- **Akka**
  - a toolkit and runtime for building highly concurrent, distributed, and resilient message-driven applications on the JVM
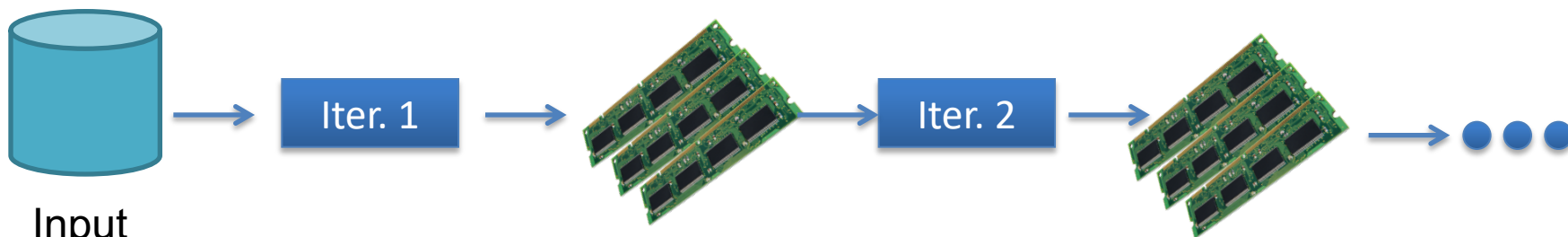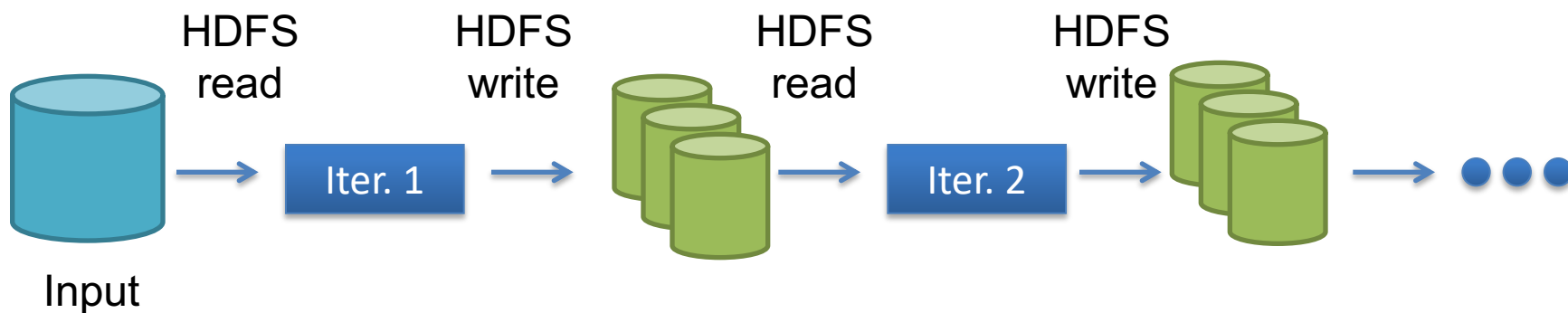
- **Cassandra**
  - distributed, highly available database designed to handle large amounts of data across multiple datacenters

- **Kafka**
  - a high-throughput, low-latency distributed messaging system designed for handling real-time data feeds

# Hadoop vs. Spark



HDFS read · HDFS write · HDFS read · HDFS write

Input → Iter. 1 → Iter. 2 → ●●●

Input → Iter. 1 → Iter. 2 → ●●●

# Summary

- Big Data

- Artificial Intelligence

- Deep Learning

- <span style="color:red">Architectures of Big Data Analytics</span>

- <span style="color:red">Data Mining Process</span>

- Fundamental Big Data:
  MapReduce Paradigm,
  Hadoop and Spark Ecosystem

# References

- EMC Education Services (2015), Data Science and Big Data Analytics: Discovering, Analyzing, Visualizing and Presenting Data, Wiley

- Stuart Russell and Peter Norvig (2016), Artificial Intelligence: A Modern Approach, 3rd Edition, Pearson International

- Peter C. Verhoef and Edwin Kooge (2016), Creating Value with Big Data Analytics: Making Smarter Marketing Decisions, Routledge

- Shiva Achari (2015), Hadoop Essentials - Tackling the Challenges of Big Data with Hadoop, Packt Publishing

- Mike Frampton (2015), Mastering Apache Spark, Packt Publishing

- Deepak Ramanathan (2014),
  SAS Modernization architectures - Big Data Analytics,
  http://www.slideshare.net/deepakramanathan/sas-modernization-architectures-big-data-analytics