# Big Data Mining
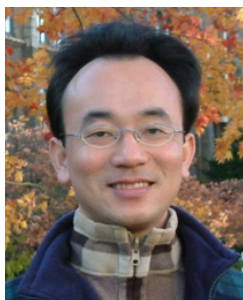# 巨量資料探勘
# AI人工智慧與大數據分析
# (Artificial Intelligence and
# Big Data Analytics)

1082DM02
MI4 (M2244) (2744)
Tue  3, 4 (10:10-12:00) (B218)

**Min-Yuh Day**
**戴敏育**
**Associate Professor**
**副教授**
**Dept. of Information Management, Tamkang University**
**淡江大學 資訊管理學系**

http://mail. tku.edu.tw/myday/
2020-03-10

# 課程大綱 (Syllabus)

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

1　2020/03/03　巨量資料探勘課程介紹
(Course Orientation for Big Data Mining)

2　2020/03/10　AI人工智慧與大數據分析
(Artificial Intelligence and Big Data Analytics)

3　2020/03/17　分群分析 (Cluster Analysis)

4　2020/03/24　個案分析與實作一 (SAS EM 分群分析)：
Case Study 1 (Cluster Analysis - K-Means using SAS EM)

5　2020/03/31　關連分析 (Association Analysis)

6　2020/04/07　個案分析與實作二 (SAS EM 關連分析)：
Case Study 2 (Association Analysis using SAS EM)

7　2020/04/14　分類與預測 (Classification and Prediction)

8　2020/04/21　期中報告 (Midterm Project Presentation)

# 課程大綱 (Syllabus)

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

9　2020/04/28　期中考試週

10　2020/05/05　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)

11　2020/05/12　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)

12　2020/05/19　機器學習與深度學習
　　　　　　　　(Machine Learning and Deep Learning)

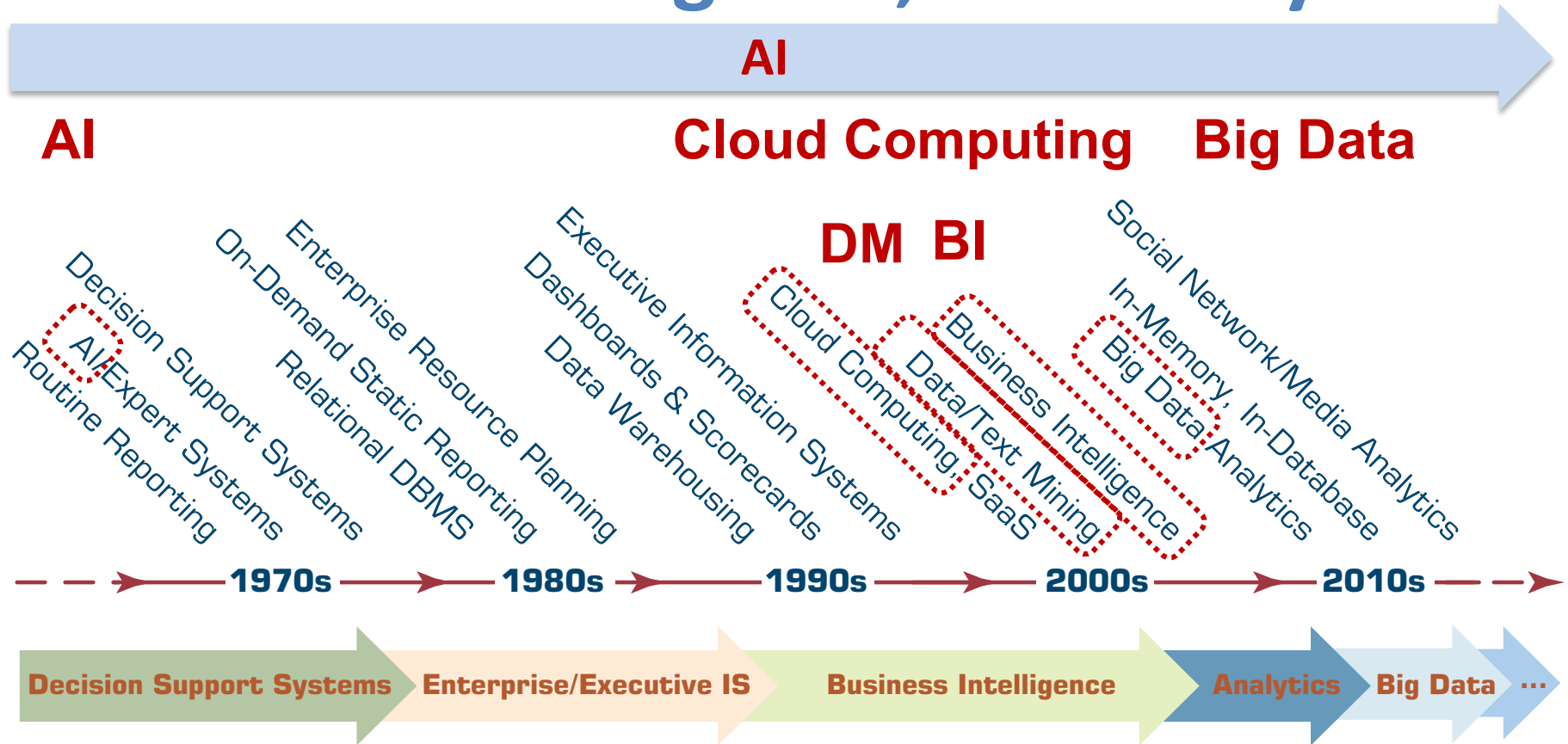13　2020/05/26　期末報告 (Final Project Presentation)

14　2020/06/02　畢業考試週

15　2020/06/09　教師彈性補充教學

# Outline

- AI

- Big Data Analytics

# AI, Big Data, Cloud Computing

## Evolution of Decision Support, Business Intelligence, and Analytics



Source: Ramesh Sharda, Dursun Delen, and Efraim Turban (2017),
Business Intelligence, Analytics, and Data Science: A Managerial Perspective, 4th Edition, Pearson

# 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

| | |
|---|---|
| **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** |

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

12

# 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



#13 ON TRENDING
What's new, Atlas?

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

# Humanoid Robot: Sophia



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

# Can a robot pass a university entrance exam?
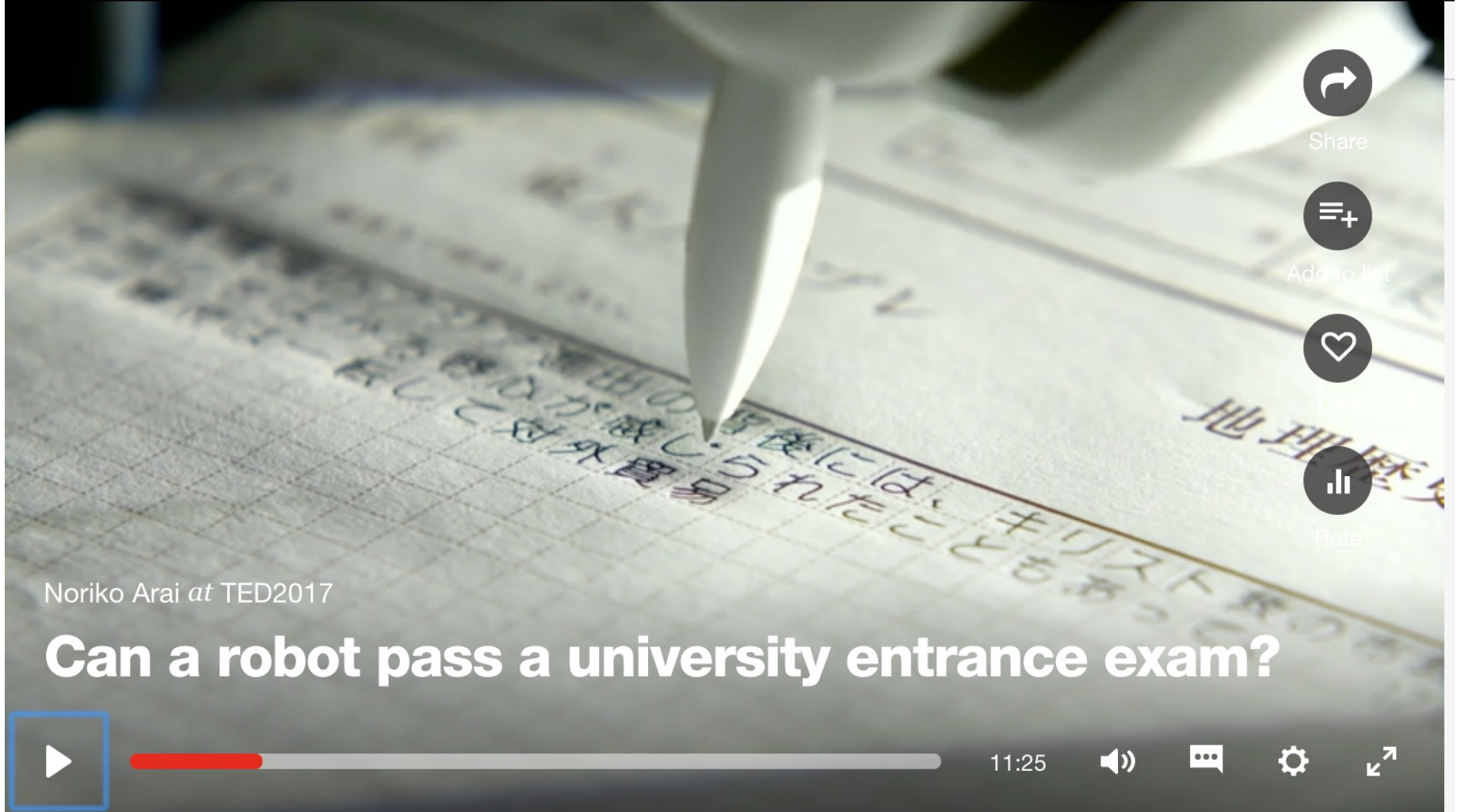## Noriko Arai at TED2017
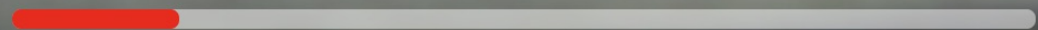


TED Ideas worth spreading — WATCH    DISCOVER    ATT

Noriko Arai *at* TED2017

**Can a robot pass a university entrance exam?**

11:25

https://www.ted.com/talks/noriko_arai_can_a_robot_pass_a_university_entrance_exam
https://www.youtube.com/watch?v=XQZjkPyJ8KU

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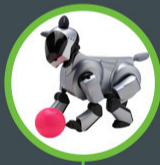Many false starts and dead-ends leave A.I. out in the cold

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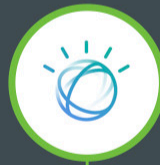
**2002 — ROOMBA**
First mass produced autonomous robotic vacuum cleaner from iRobot learns to navigate and clean homes

**2011 — SIRI**
Apple integrates Siri, an intelligent virtual assistant with a voice interface, into the iPhone 4S
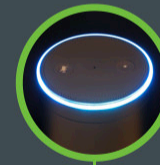
**2011 — WATSON**
IBM's question answering computer Watson wins first place on popular $1M prize television quiz show *Jeopardy*
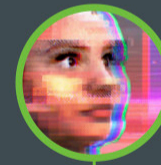
**2014 — EUGENE**
Eugene Goostman, a chatbot passes the Turing Test with a third of judges believing Eugene is human

**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

SYZYGY

# Artificial Intelligence
# Machine Learning & Deep Learning



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/

# AI, ML, DL



Artificial Intelligence (AI)

Machine Learning (ML)

Supervised Learning
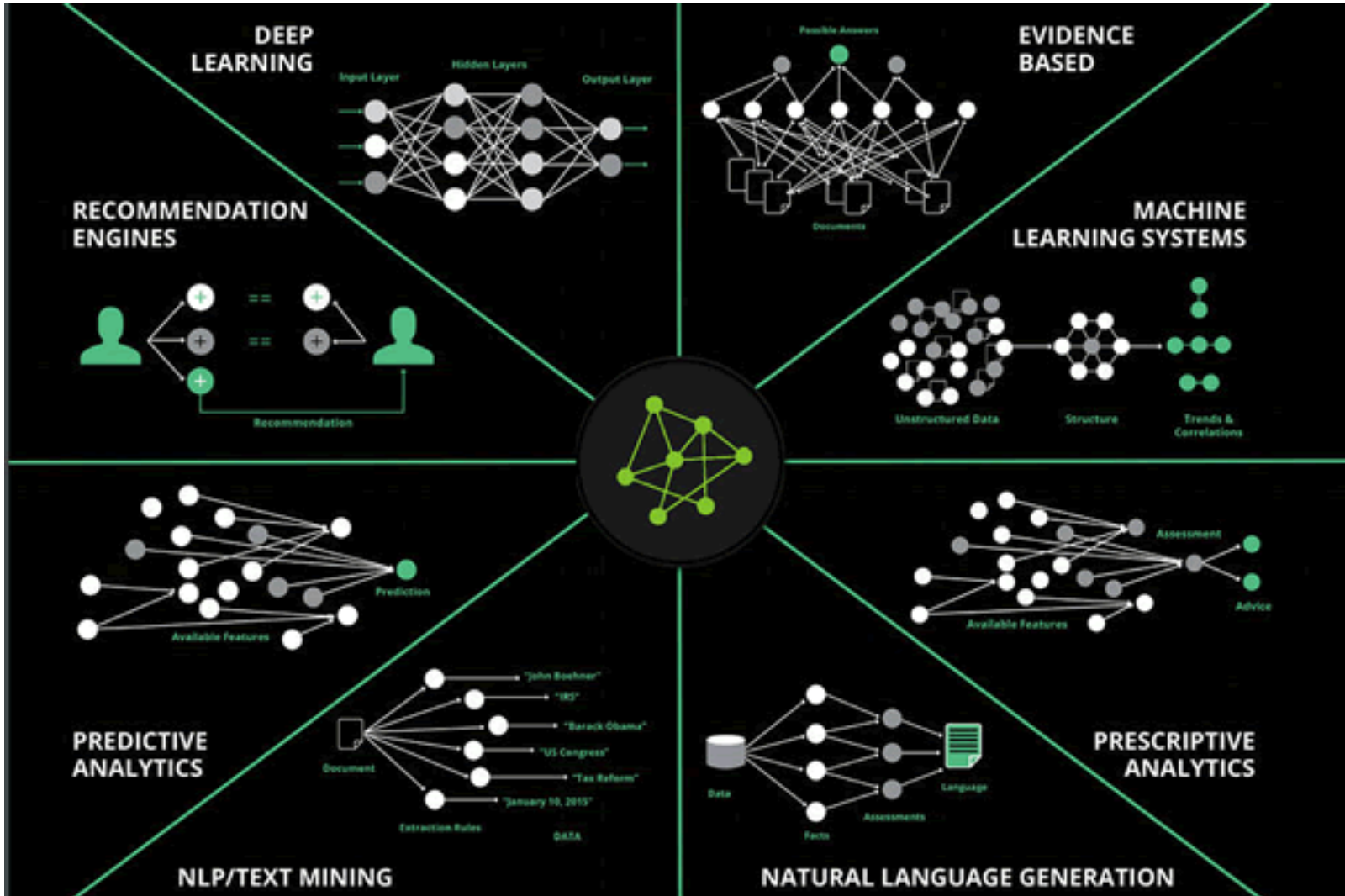
Unsupervised Learning

Deep Learning (DL)
CNN
RNN LSTM GRU
GAN

Semi-supervised Learning

Reinforcement Learning

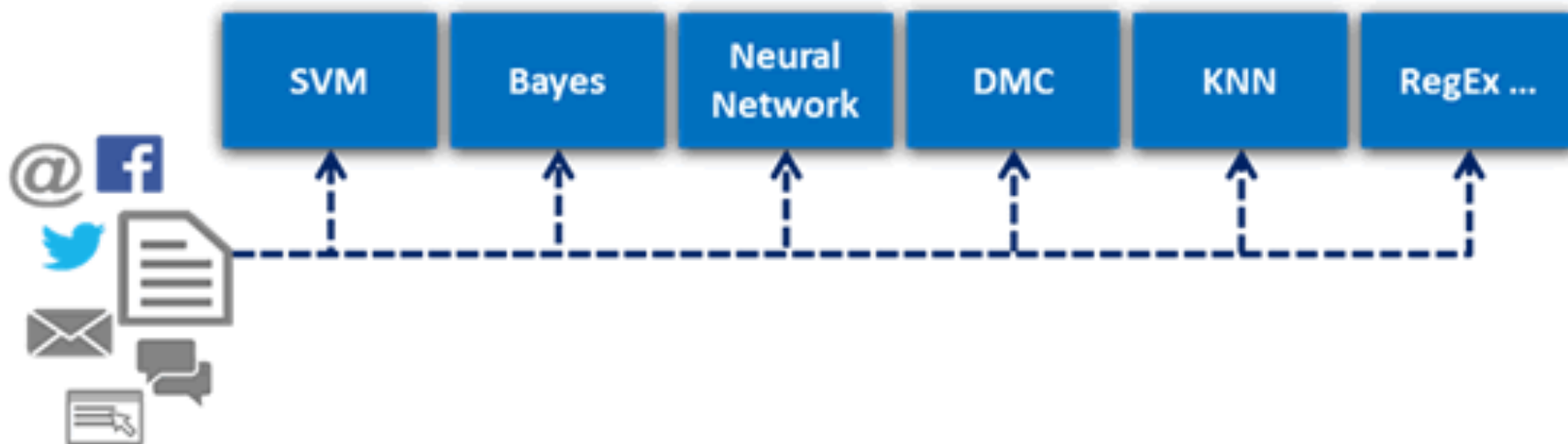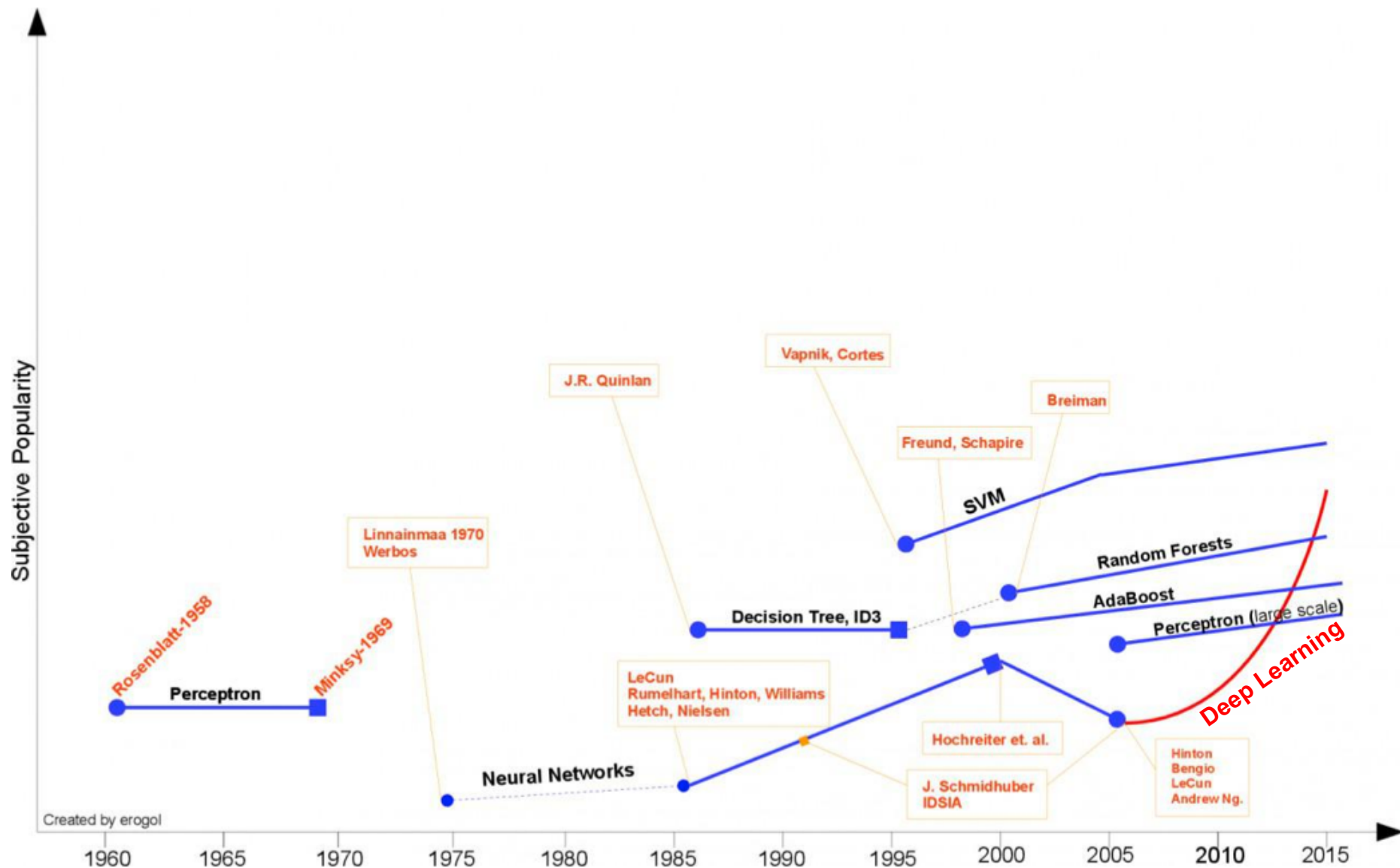# Artificial Intelligence (AI) is many things



Ecosystem of AI

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

# Deep Learning Evolution



Created by erogol

# 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)



Source: Jesus Serrano-Guerrero, Jose A. Olivas, Francisco P. Romero, and Enrique Herrera-Viedma (2015), "Sentiment analysis: A review and comparative analysis of web services," Information Sciences, 311, pp. 18-38.

# Artificial intelligence (AI) in optical networks

# 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 ]

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

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

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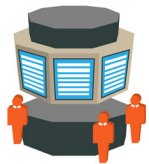are watched on YouTube each month

**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

# Data Mining
# Is a Blend of Multiple Disciplines

Source: Ramesh Sharda, Dursun Delen, and Efraim Turban (2017), Business Intelligence, Analytics, and Data Science: A Managerial Perspective, 4th Edition, Pearson

**Stephan Kudyba (2014),**
**Big Data, Mining, and Analytics:**
**Components of Strategic Decision Making, Auerbach Publications**

# Architecture of Big Data Analytics

**Big Data Sources**

* Internal

* External

* Multiple formats

* Multiple locations

* Multiple applications

**Big Data Transformation**
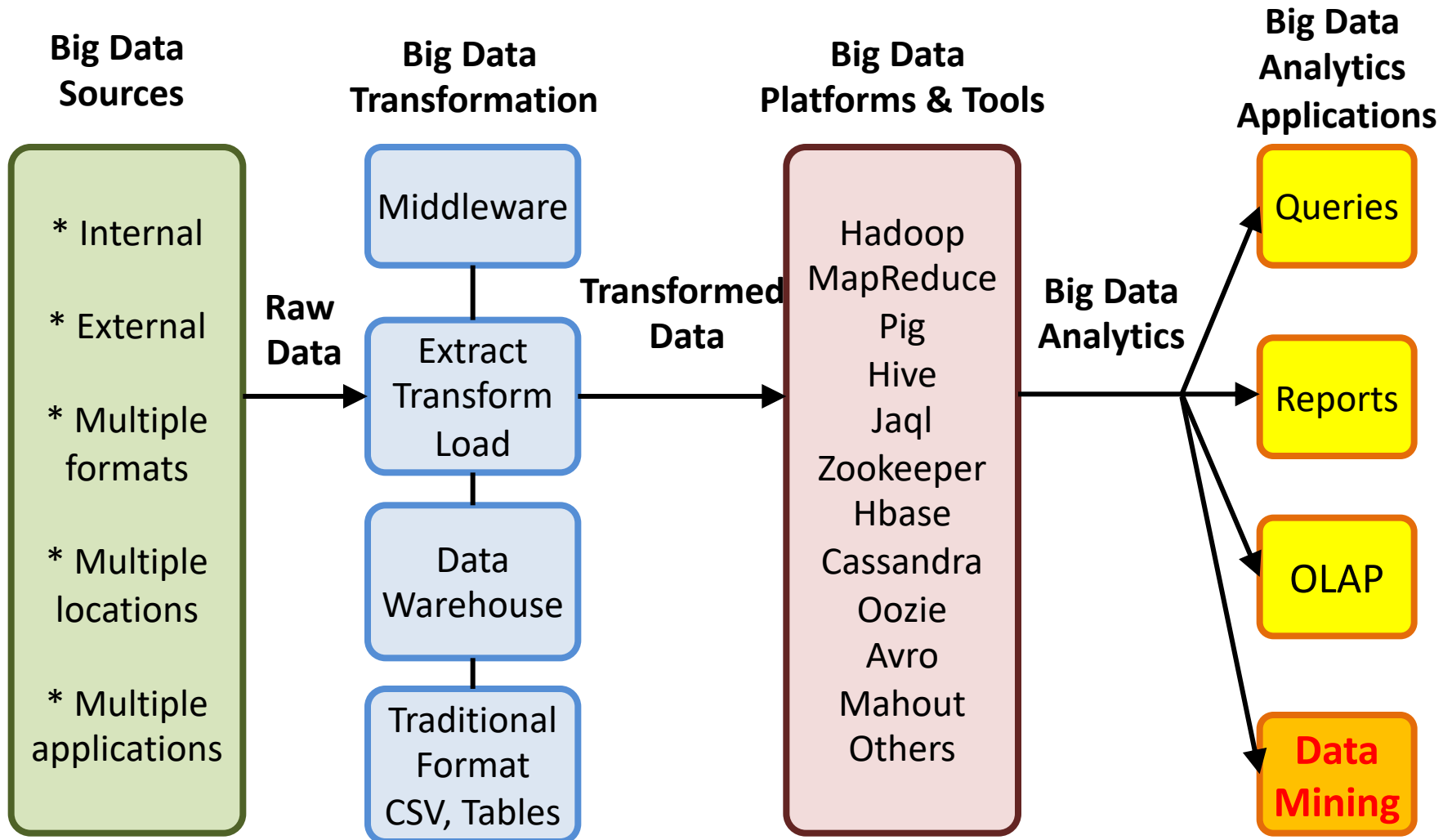
Middleware

Extract Transform Load

Data Warehouse

Traditional Format CSV, Tables

**Raw Data**

**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 |
|---|---|---|---|
| * Internal | **Data Mining** | | Queries |
| * External | **Big Data** | | Reports |
| * Multiple formats | **Analytics** | | OLAP |
| * Multiple locations | **Applications** | | **Data Mining** |
| * Multiple applications | | | |

# Data Mining Tasks & Methods

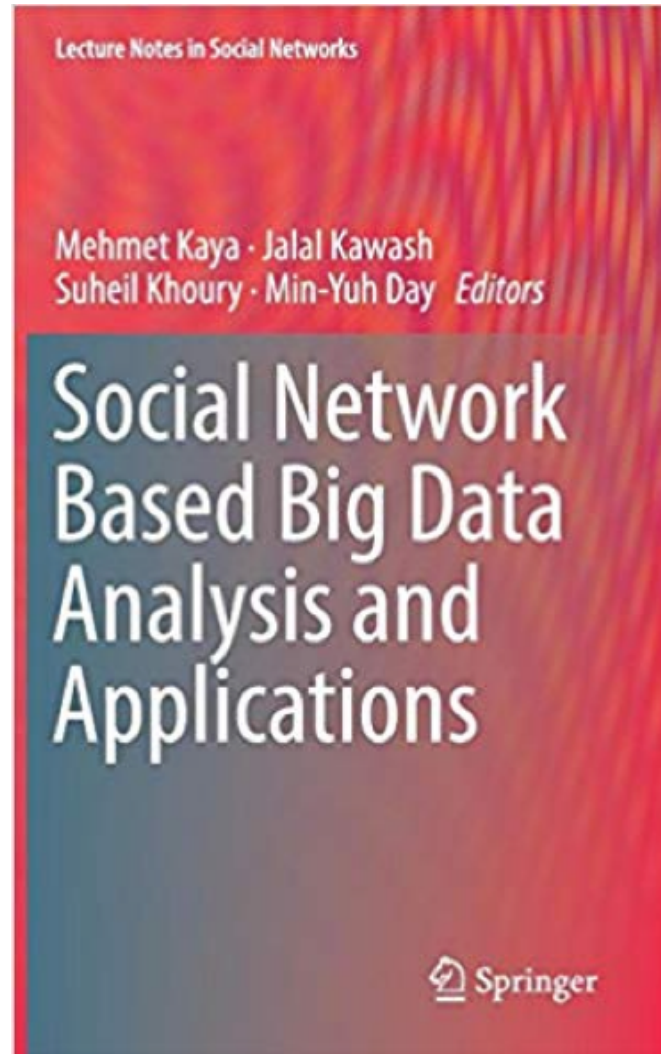| Data Mining Tasks & Methods | | Data Mining Algorithms | Learning Type |
|---|---|---|---|
| **Prediction** | | | |
| | Classification | Decision Trees, Neural Networks, Support Vector Machines, kNN, Naïve Bayes, GA | Supervised |
| | Regression | Linear/Nonlinear Regression, ANN, Regression Trees, SVM, kNN, GA | Supervised |
| | Time series | Autoregressive Methods, Averaging Methods, Exponential Smoothing, ARIMA | Supervised |
| **Association** | | | |
| | Market-basket | Apriori, OneR, ZeroR, Eclat, GA | Unsupervised |
| | Link analysis | Expectation Maximization, Apriori Algorithm, Graph-Based Matching | Unsupervised |
| | Sequence analysis | Apriori Algorithm, FP-Growth, Graph-Based Matching | Unsupervised |
| **Segmentation** | | | |
| | Clustering | k-means, Expectation Maximization (EM) | Unsupervised |
| | Outlier analysis | k-means, Expectation Maximization (EM) | Unsupervised |

# Business Intelligence, Analytics, and Data Science:
## A Managerial Perspective, 4th Edition,
## Ramesh Sharda, Dursun Delen, and Efraim Turban, Pearson, 2017.

# Big Data, Data Mining, and Machine Learning: Value Creation for Business Leaders and Practitioners,
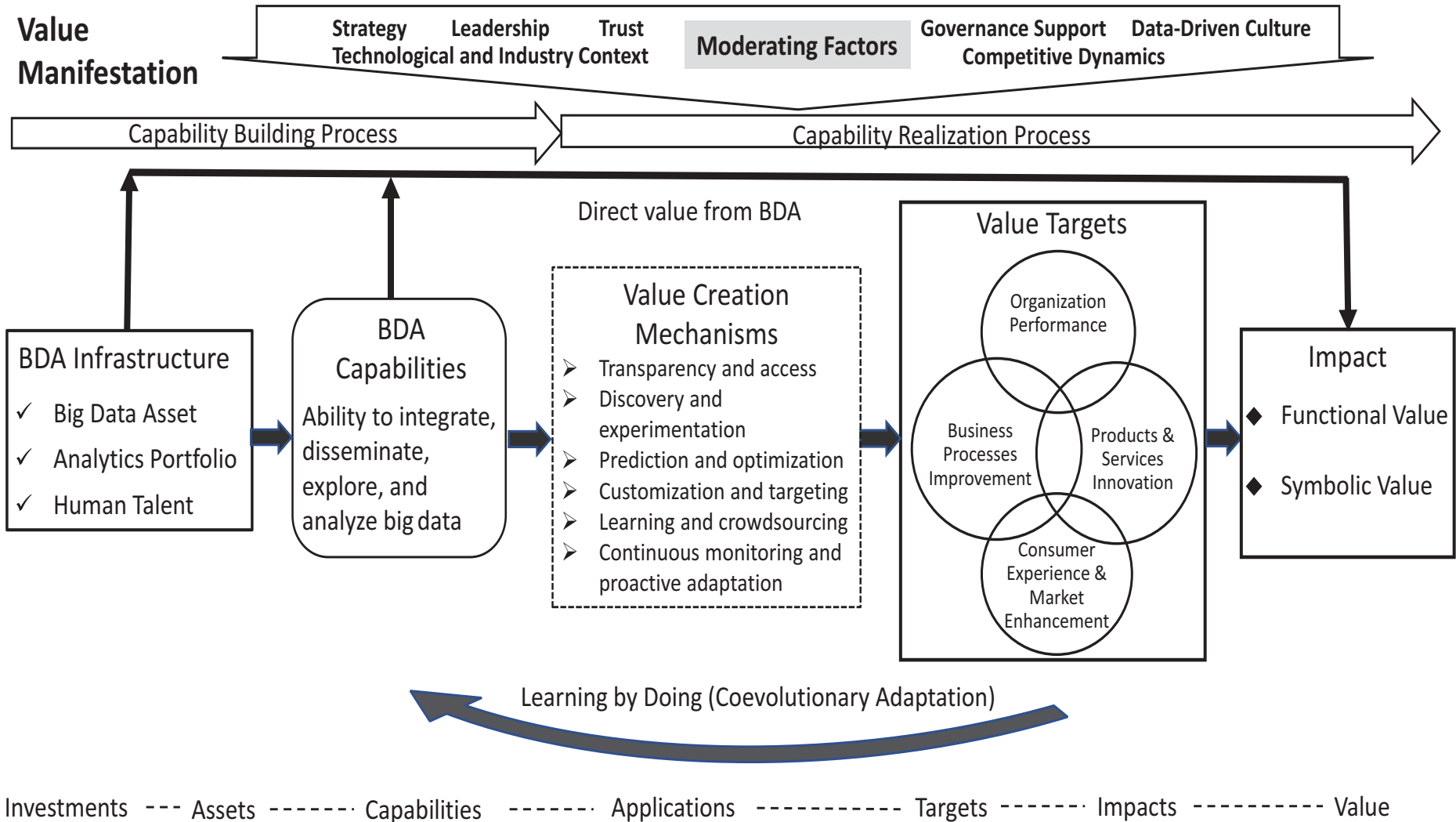## Jared Dean,
## Wiley, 2014.

**Social Network Based Big Data Analysis and Applications,**
**Lecture Notes in Social Networks,**
**Mehmet Kaya, Jalal Kawash, Suheil Khoury, Min-Yuh Day,**
**Springer International Publishing, 2018.**
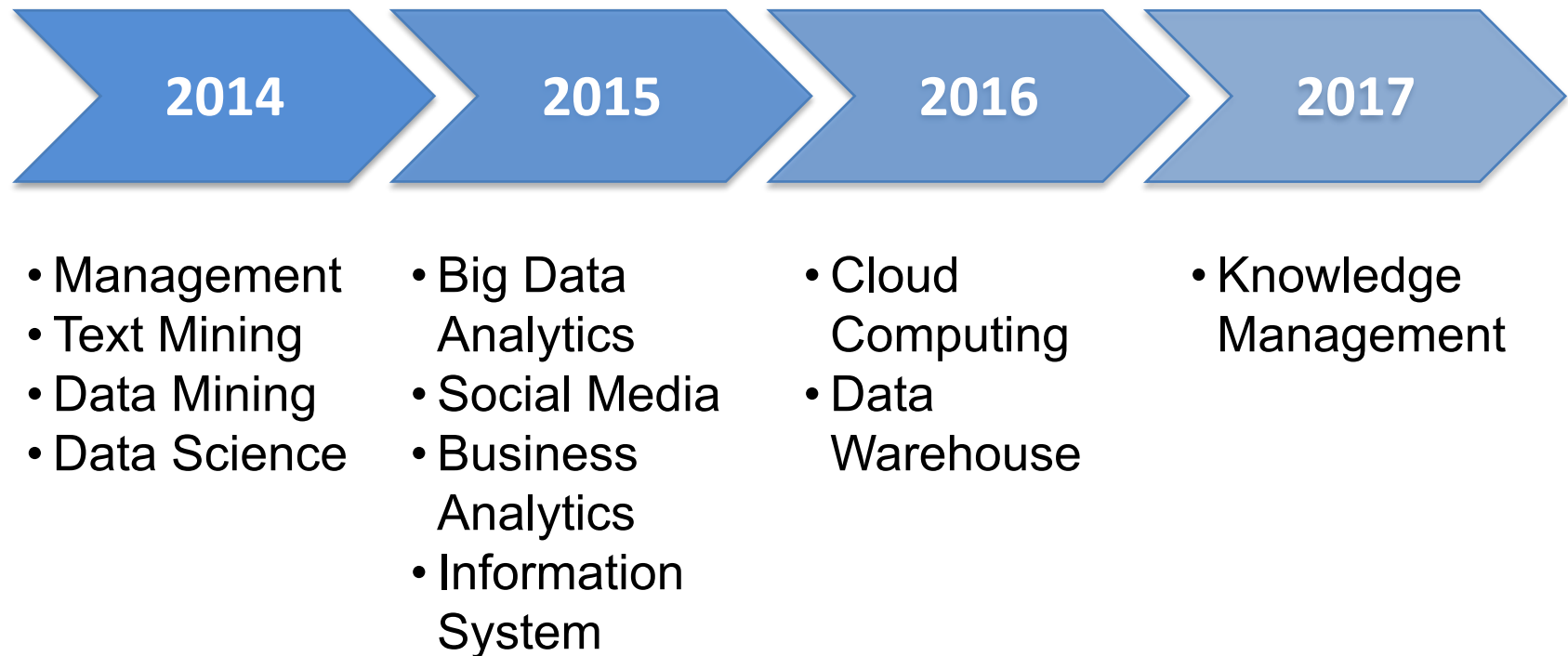
# Value Creation by Big Data Analytics
## (Grover et al., 2018)

**Value Manifestation**

| | Strategy   Leadership   Trust<br>Technological and Industry Context | Moderating Factors | Governance Support   Data-Driven Culture<br>Competitive Dynamics |
|---|---|---|---|

Capability Building Process → Capability Realization Process

Direct value from BDA

**BDA Infrastructure**
- ✓ Big Data Asset
- ✓ Analytics Portfolio
- ✓ Human Talent

**BDA Capabilities**

Ability to integrate, disseminate, explore, and analyze big data

**Value Creation Mechanisms**
- ➤ Transparency and access
- ➤ Discovery and experimentation
- ➤ Prediction and optimization
- ➤ Customization and targeting
- ➤ Learning and crowdsourcing
- ➤ Continuous monitoring and proactive adaptation

**Value Targets**
- Organization Performance
- Business Processes Improvement
- Products & Services Innovation
- Consumer Experience & Market Enhancement

**Impact**
- ◆ Functional Value
- ◆ Symbolic Value

Learning by Doing (Coevolutionary Adaptation)

Investments --- Assets ------ Capabilities ------ Applications ---------- Targets ------ Impacts ---------- Value
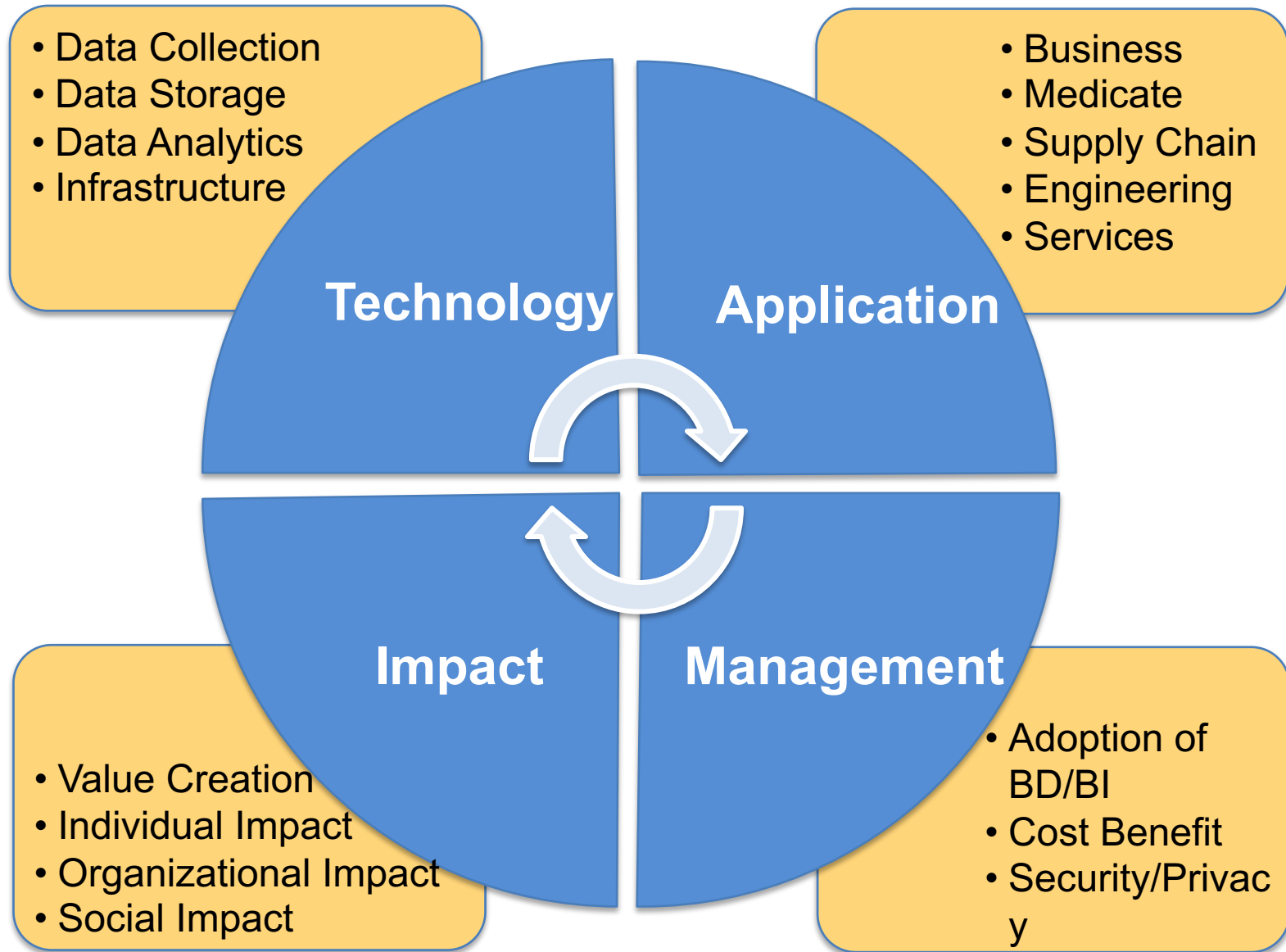
# Research Landscape of
# Business Intelligence and Big Data Analytics:
# A bibliometrics study

- A bibliometric analysis on Big Data and Business Intelligence from 1990 to 2016.

- Big Data papers grow much faster than Business Intelligence papers

- Computer Science and information systems are two core disciplines.

- Most influential papers are identified and a research framework is proposed.
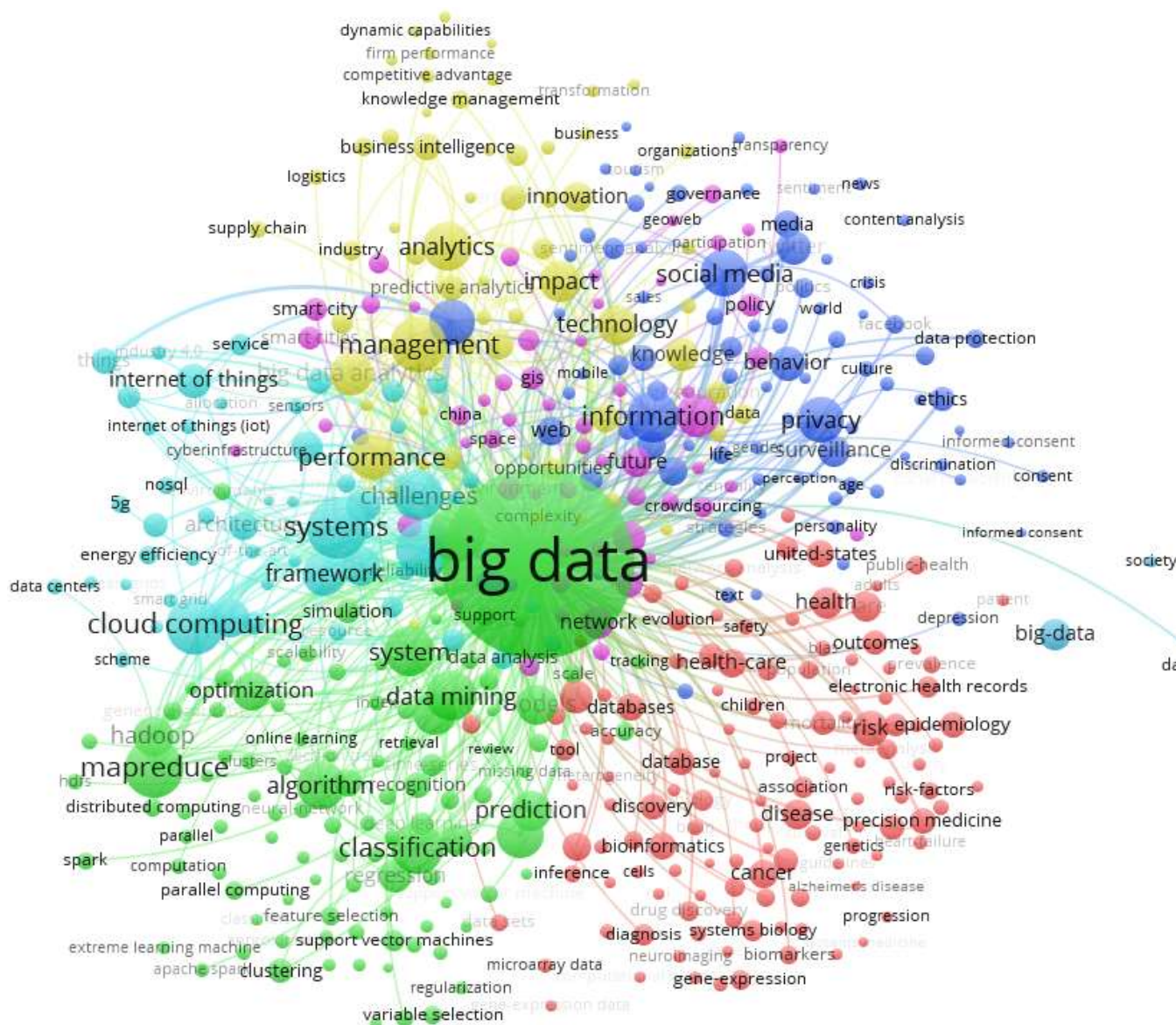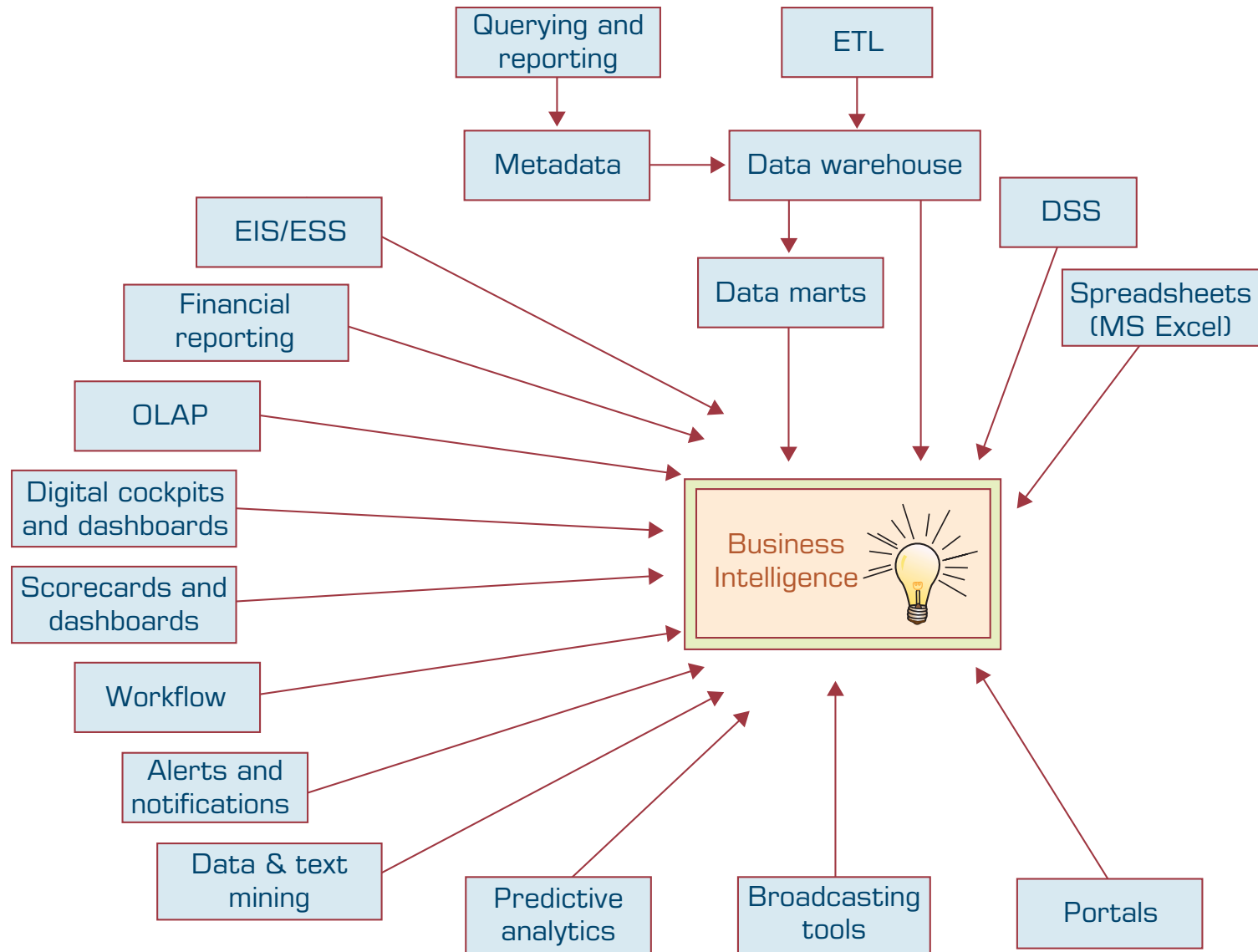
# Evolution of top keywords in "BD & BI" publications

**2014**
- Management
- Text Mining
- Data Mining
- Data Science

**2015**
- Big Data Analytics
- Social Media
- Business Analytics
- Information System

**2016**
- Cloud Computing
- Data Warehouse

**2017**
- Knowledge Management

# Framework for BD and BI Research



• Data Collection
• Data Storage
• Data Analytics
• Infrastructure

**Technology**

**Application**

• Business
• Medicate
• Supply Chain
• Engineering
• Services

**Impact**

**Management**

• Value Creation
• Individual Impact
• Organizational Impact
• Social Impact

• Adoption of BD/BI
• Cost Benefit
• Security/Privacy
• Human Resource

# Business Intelligence and Big Data analytics



Source: Ting-Peng Liang and Yu-Hsi Liu (2018), "Research Landscape of Business Intelligence and Big Data analytics: A bibliometrics study",
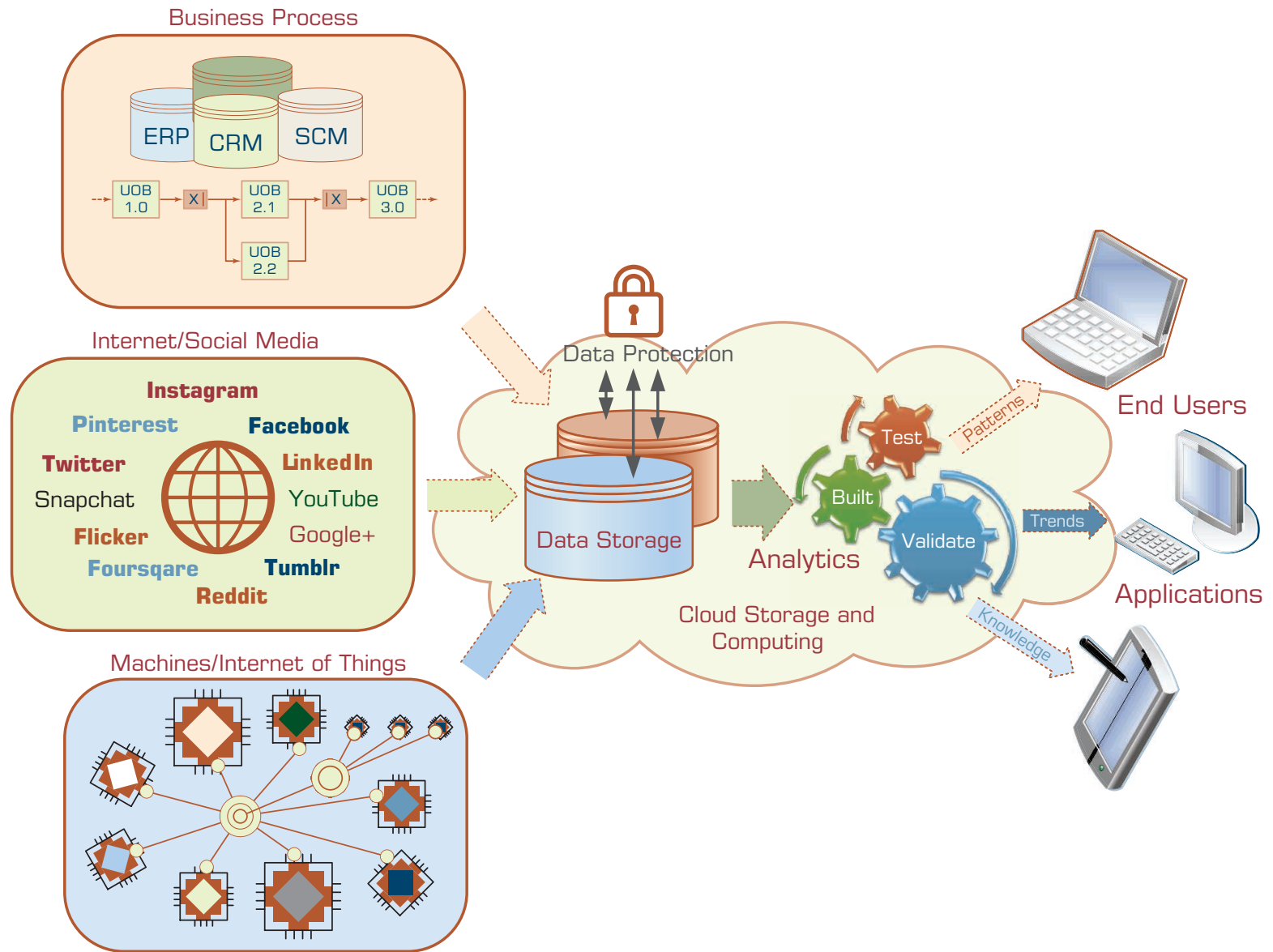
42

# Evolution of Business Intelligence (BI)
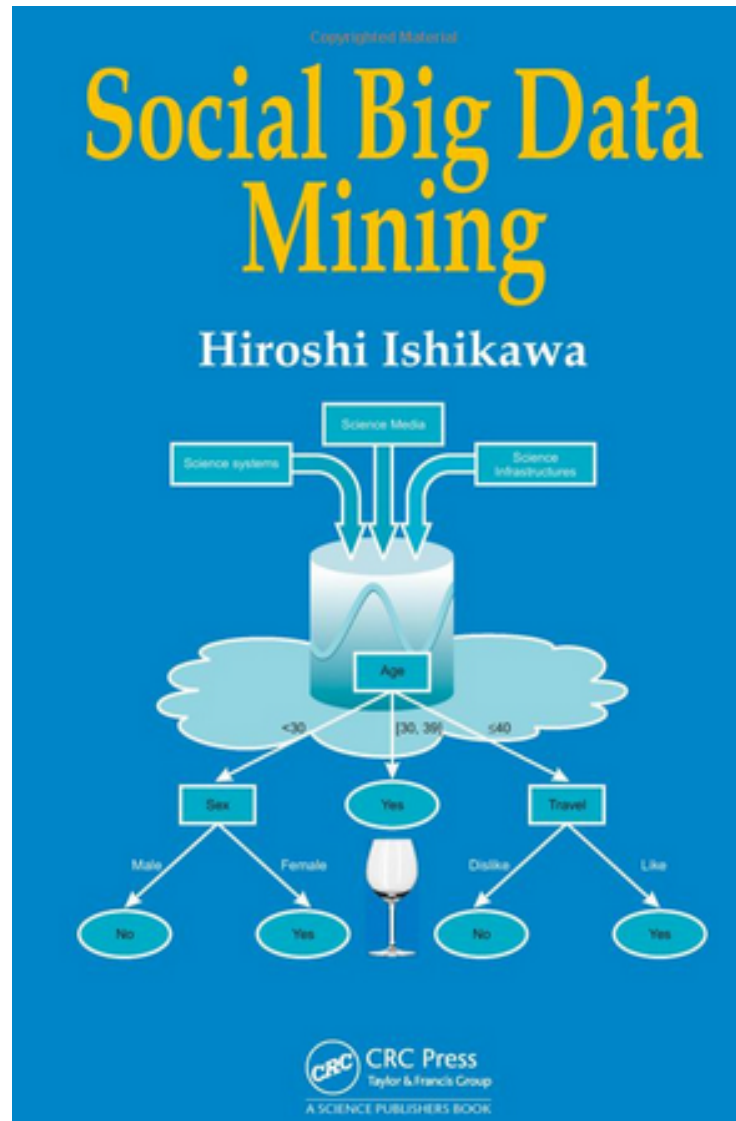
# A High-Level Architecture of BI

**Data Warehouse Environment**  **Business Analytics Environment**  **Performance and Strategy**

Data Sources

Technical staff

Build the data warehouse
- Organizing
- Summarizing
- Standardizing

Data warehouse

Business users

Access

Manipulation, results

Managers/executives

BPM strategies

Future component: Intelligent systems

User interface
- Browser
- Portal
- Dashboard

# Three Types of Analytics

**Business Analytics**

| | Descriptive | Predictive | Prescriptive |
|---|---|---|---|
| **Questions** | What happened? What is happening? | What will happen? Why will it happen? | What should I do? Why should I do it? |
| **Enablers** | ✓ Business reporting<br>✓ Dashboards<br>✓ Scorecards<br>✓ Data warehousing | ✓ Data mining<br>✓ Text mining<br>✓ Web/media mining<br>✓ Forecasting | ✓ Optimization<br>✓ Simulation<br>✓ Decision modeling<br>✓ Expert systems |
| **Outcomes** | **Well-defined business problems and opportunities** | **Accurate projections of future events and outcomes** | **Best possible business decisions and actions** |

# A Data to Knowledge Continuum

46

# Social Big Data Mining

**(Hiroshi Ishikawa, 2015)**

# 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**
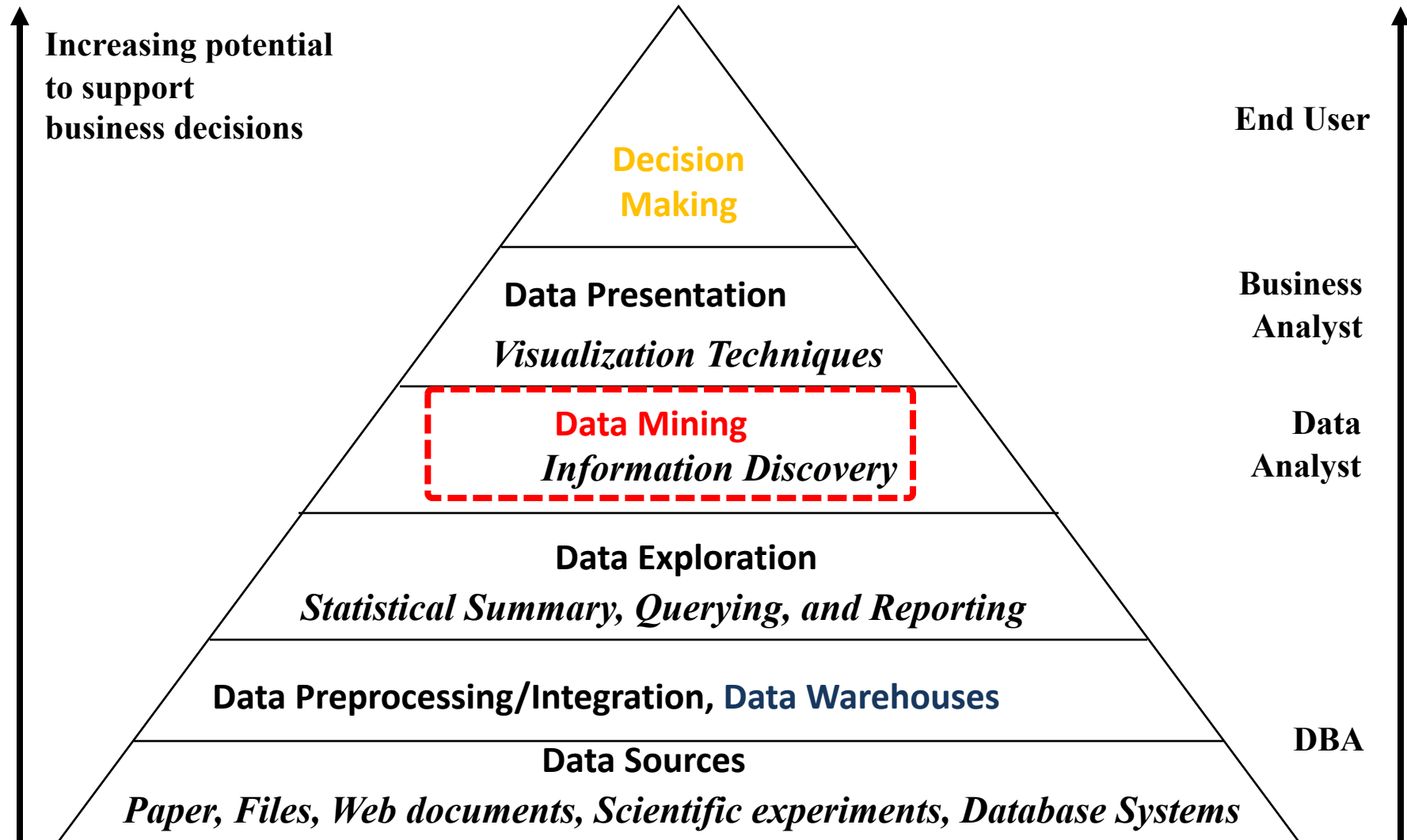- **Large-scale visualization**

- **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

# Business Intelligence and Data Mining



**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

# Data Mining at the Intersection of Many Disciplines

# 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?

# 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

| | |
|---|---|
| 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% |

Legend: 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.**
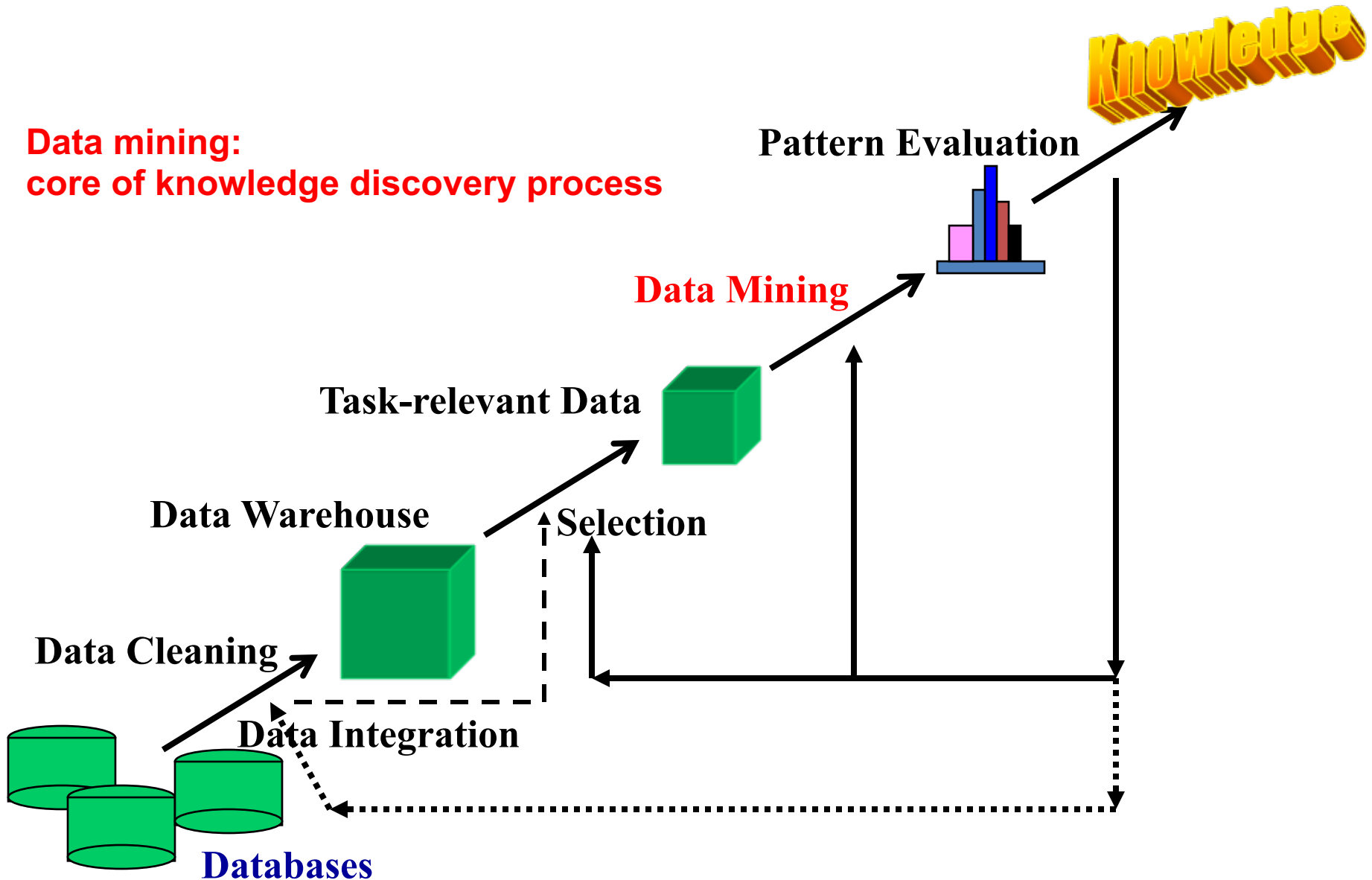Communications of the ACM, 39(11), 27-34.

# 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

**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
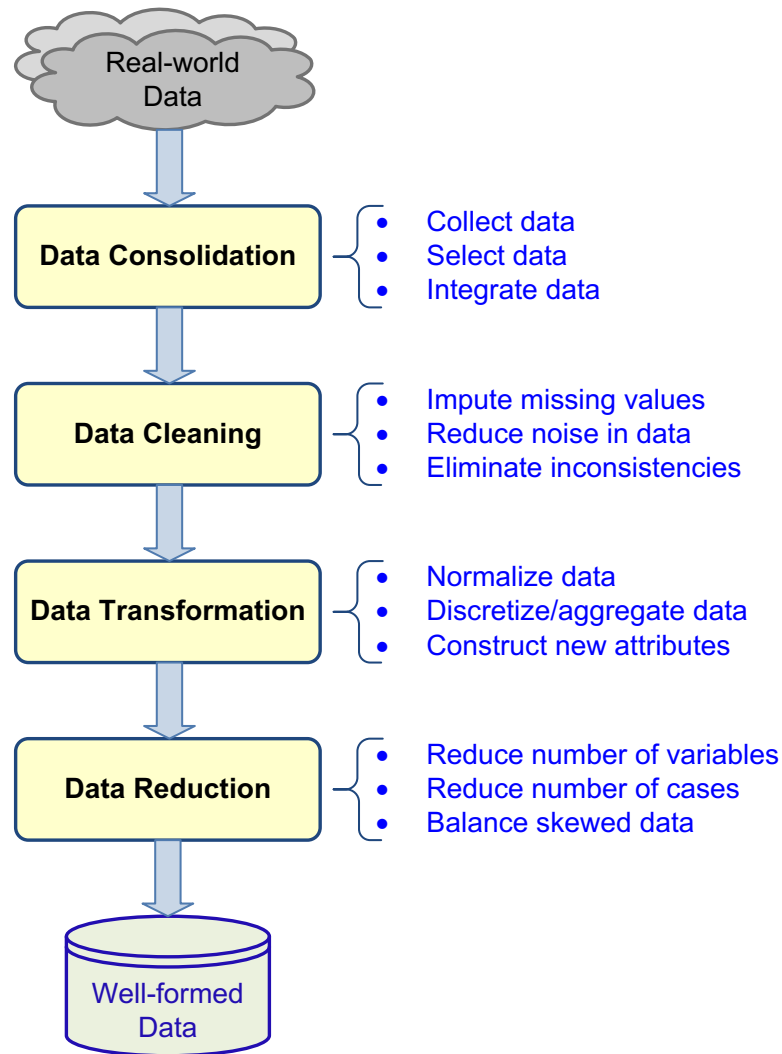
Step 5: Testing and Evaluation

Step 6: Deployment

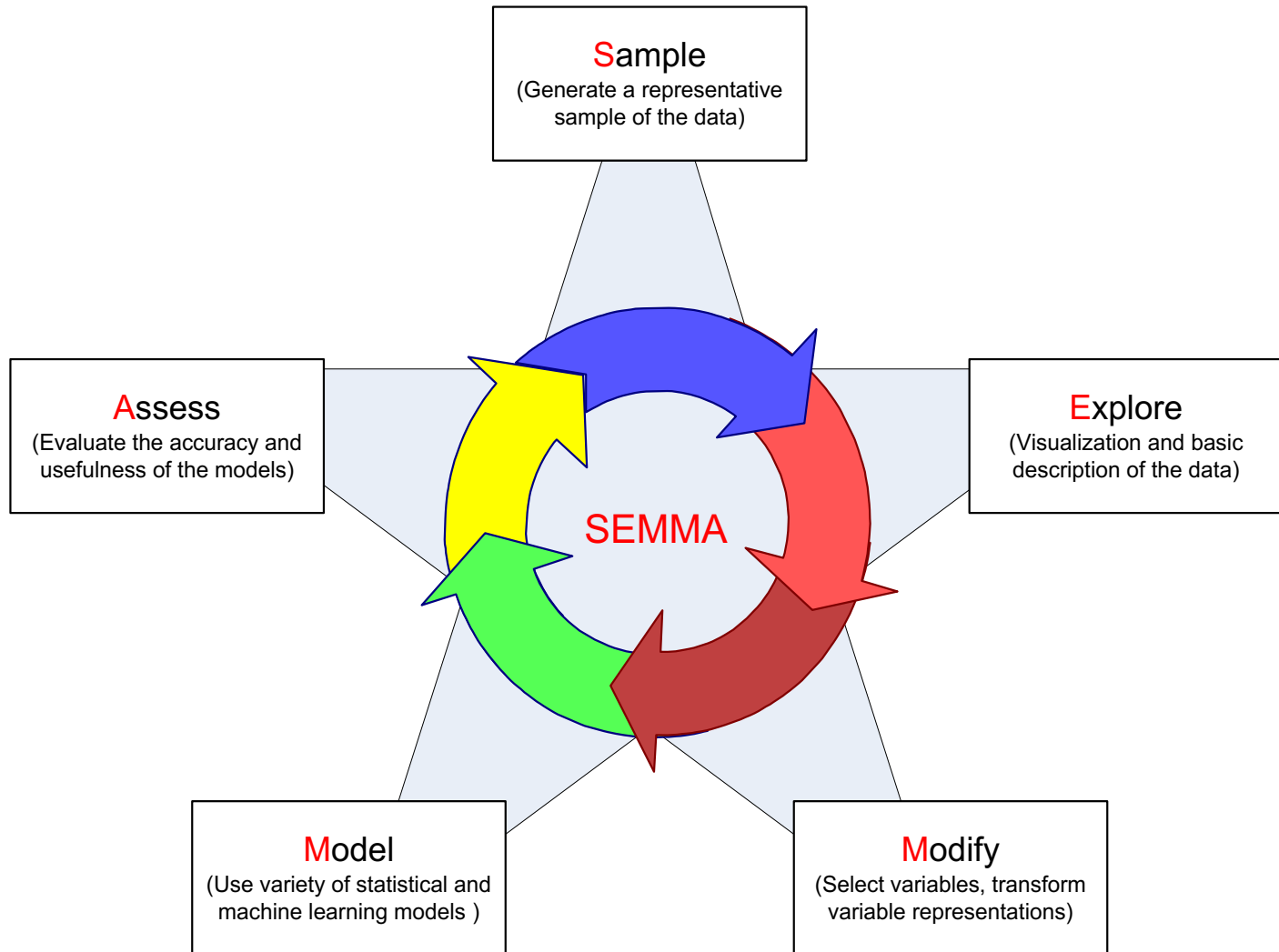- The process is highly repetitive and experimental (DM: art versus science?)

# Data Preparation – A Critical DM Task

# Data Mining Process: SEMMA



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

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

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

SEMMA

**Model** (Use variety of statistical and machine learning models )
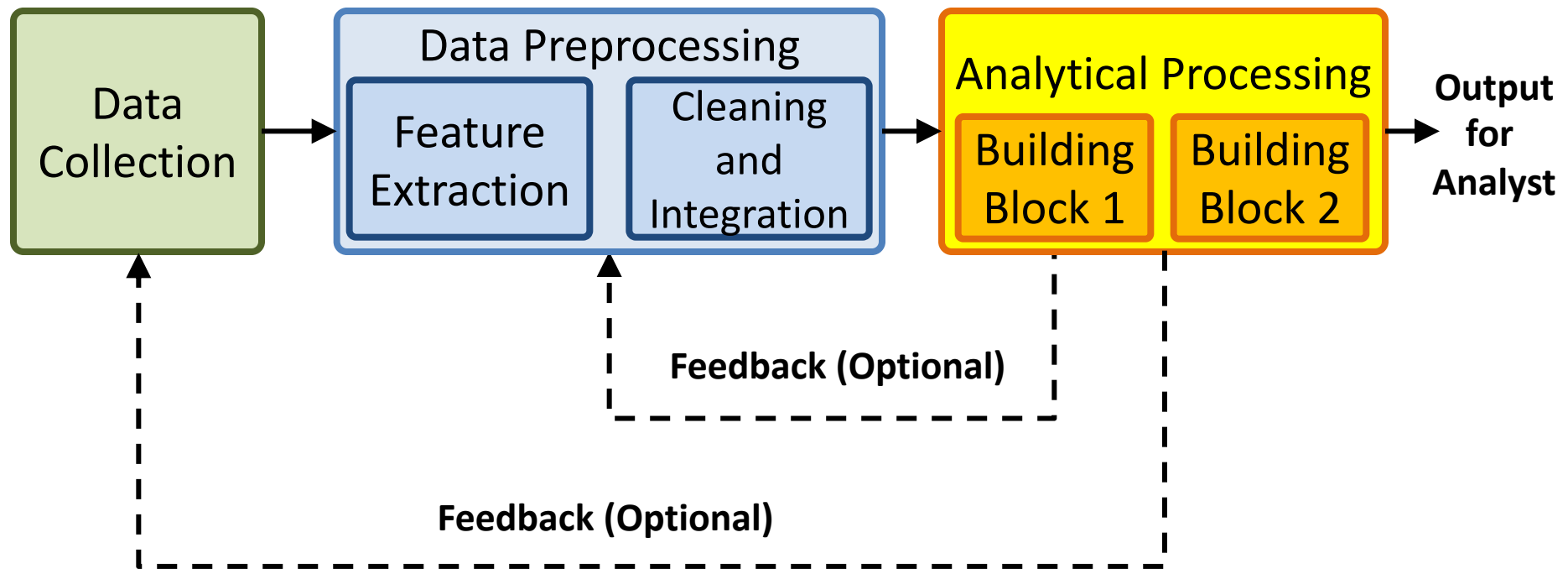
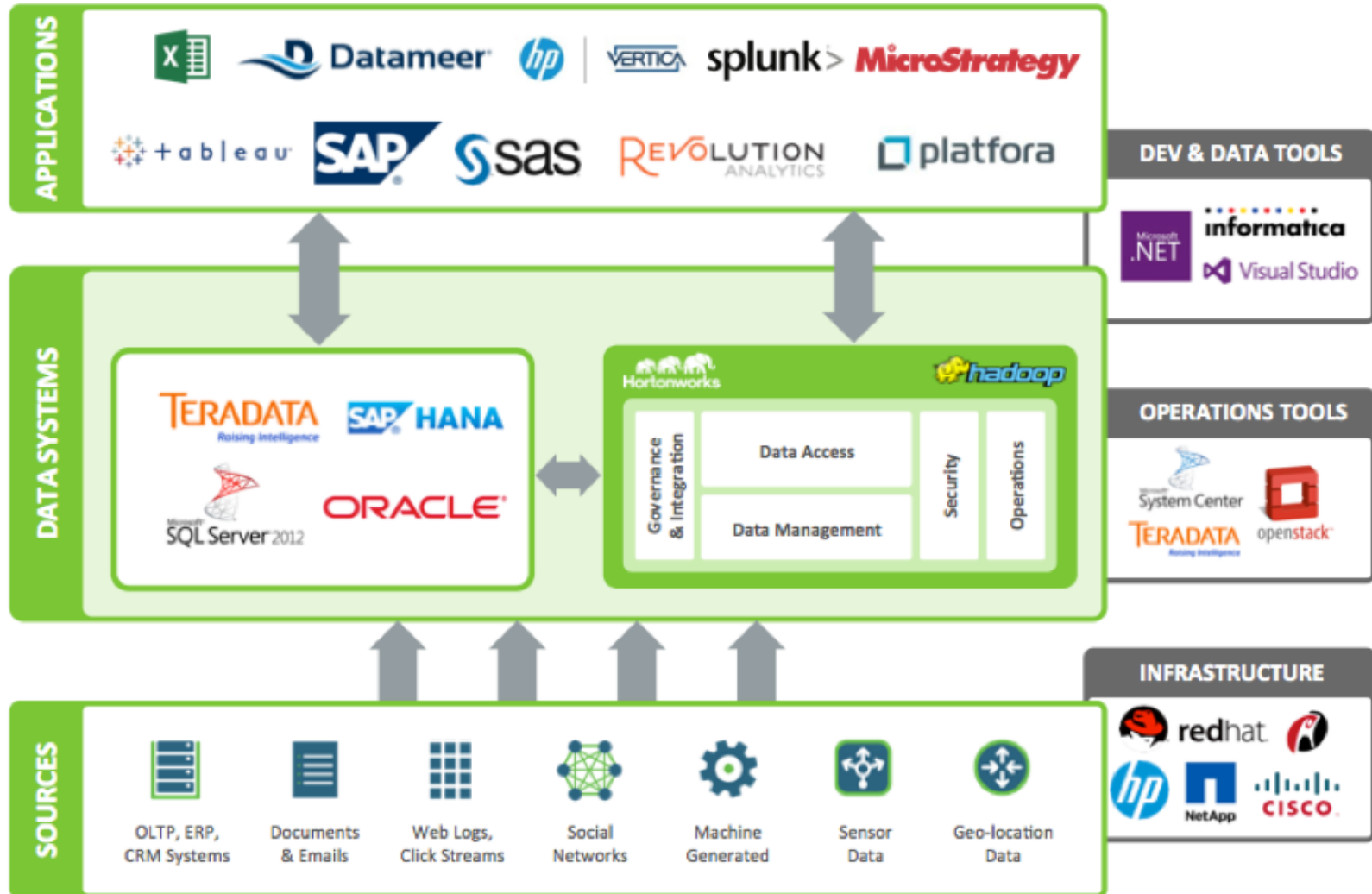**Modify** (Select variables, transform variable representations)

# Data Mining Processing Pipeline
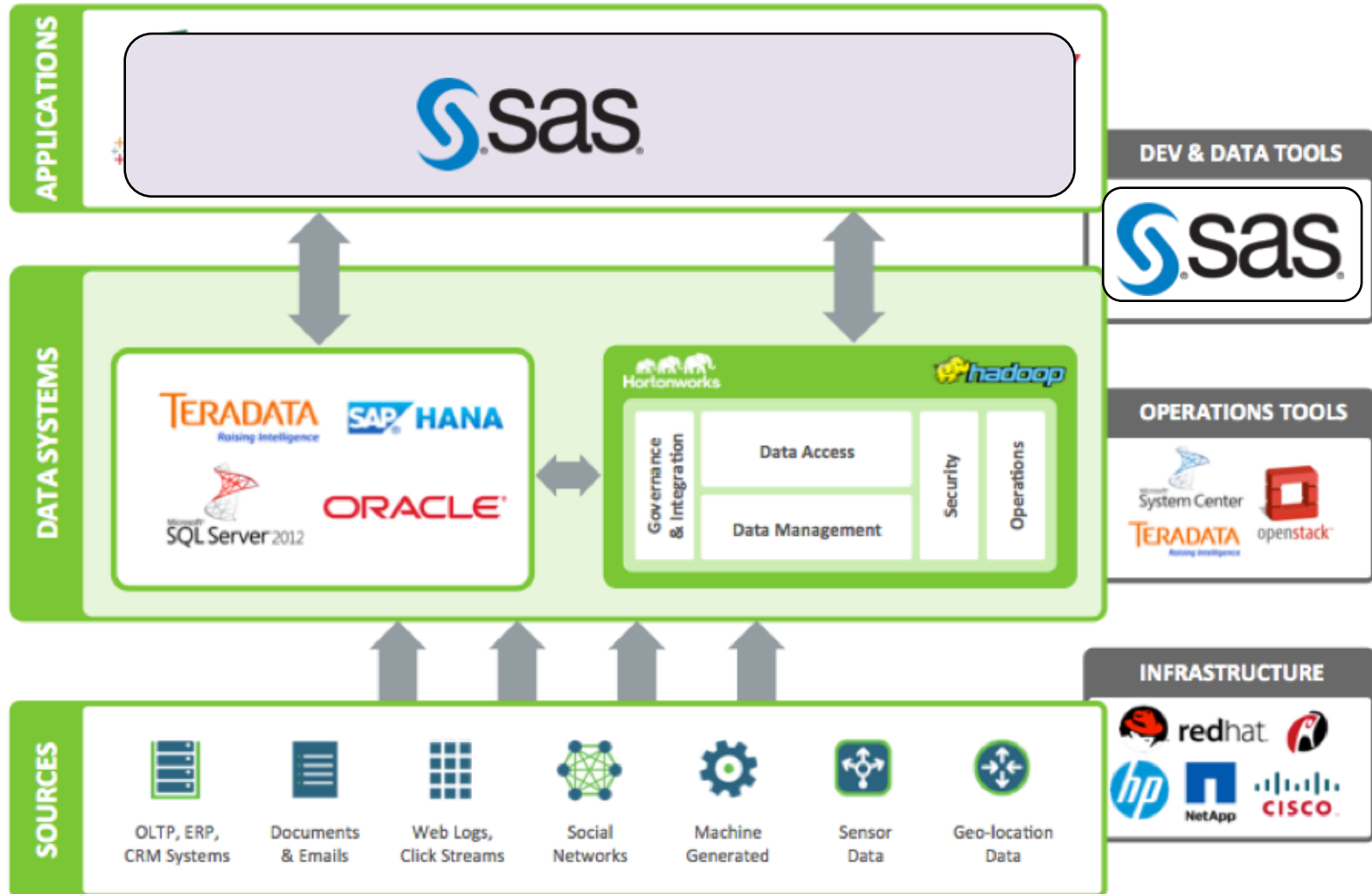
(Charu Aggarwal, 2015)

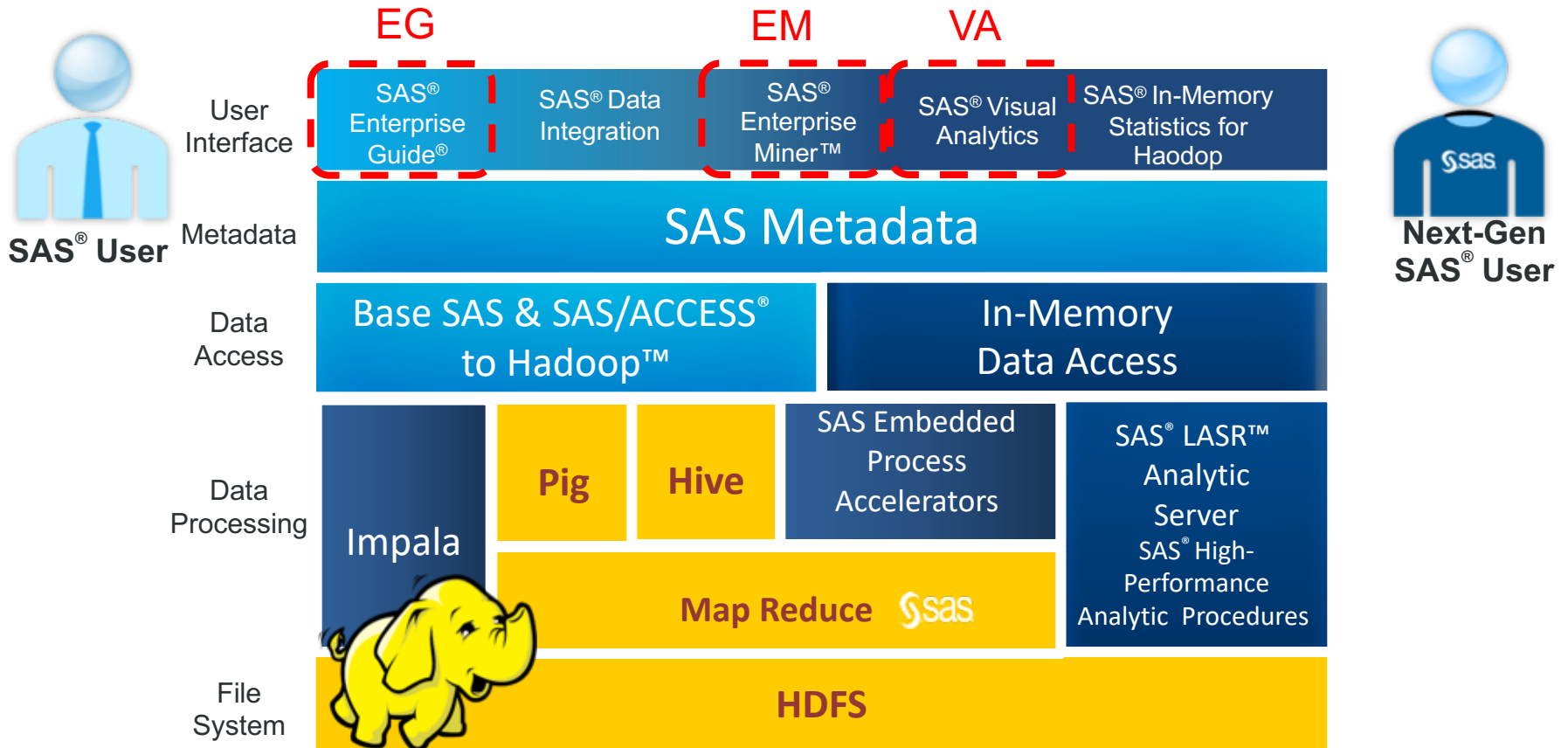# SAS Big data Strategy – SAS areas

# SAS Big data Strategy – SAS areas

# SAS® Within the HADOOP ECOSYSTEM

# Summary

- AI

- Big Data Analytics

# References

- Ramesh Sharda, Dursun Delen, and Efraim Turban (2017), Business Intelligence, Analytics, and Data Science: A Managerial Perspective, 4th Edition, Pearson.

- Jared Dean (2014), Big Data, Data Mining, and Machine Learning: Value Creation for Business Leaders and Practitioners, Wiley.

- Mehmet Kaya, Jalal Kawash, Suheil Khoury, and Min-Yuh Day (2018), Social Network Based Big Data Analysis and Applications, Lecture Notes in Social Networks, Springer International Publishing.

- Varun Grover, Roger HL Chiang, Ting-Peng Liang, and Dongsong Zhang (2018), "Creating Strategic Business Value from Big Data Analytics: A Research Framework", Journal of Management Information Systems, 35, no. 2, pp. 388-423.

- Ting-Peng Liang and Yu-Hsi Liu (2018), "Research Landscape of Business Intelligence and Big Data analytics: A bibliometrics study", Expert Systems with Applications, 111, no. 30, pp. 2-10.

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

- Javier Mata, Ignacio de Miguel, Ramón J. Durán, Noemí Merayo, Sandeep Kumar Singh, Admela Jukan, and Mohit Chamania  (2018), "Artificial intelligence (AI) methods in optical networks: A comprehensive survey", Optical Switching and Networking, 28, pp. 43-57

- Stephan Kudyba (2014), Big Data, Mining, and Analytics: Components of Strategic Decision Making, Auerbach Publications.