

# Big Data Mining

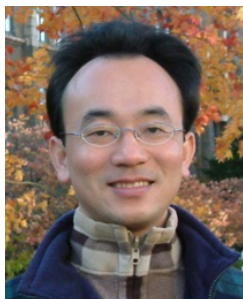
## 巨量資料探勘

# Course Orientation for Big Data Mining (巨量資料探勘課程介紹)

1072DM01

MI4 (M2244) (2849)

Wed 6, 7 (13:10-15:00) (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/>

2019-02-20



# 淡江大學107學年度第2學期

## 課程教學計畫表

### Spring 2019 (2019.02 - 2019.06)

- 課程名稱：巨量資料探勘 (Big Data Mining)
- 授課教師：戴敏育 (Min-Yuh Day)
- 開課系級：資管四P (TLMXB4P) (M2244) (2849)
- 開課資料：選修 單學期 2 學分 (2 Credits, Elective)
- 上課時間：週三 6,7 (Wed 13:10-15:00)
- 上課教室：B206

# 課程簡介

- 本課程介紹巨量資料探勘 (Big Data Mining) 的基礎概念及應用技術。
- 課程內容包括
  - 巨量資料探勘 (Big Data Mining)
  - AI人工智慧與大數據分析 (Artificial Intelligence and Big Data Analytics)
  - 關連分析 (Association Analysis)
  - 分類與預測 (Classification and Prediction)
  - 分群分析 (Cluster Analysis)
  - 機器學習與深度學習 (Machine Learning and Deep Learning)
  - SAS企業資料採礦實務 (SAS Enterprise Miner)
  - 巨量資料探勘個案分析與實作

# Course Introduction

- This course introduces the **fundamental concepts** and **applications technology** of **big data mining**.
- Topics include
  - Big Data Mining
  - Artificial Intelligence and Big Data Analytics
  - Association Analysis
  - Classification and Prediction
  - Cluster Analysis
  - Machine Learning and Deep Learning
  - Data Mining Using SAS Enterprise Miner (SAS EM)
  - Case Study and Implementation of Big Data Mining



# 課程目標 (Objective)

- 瞭解及應用 巨量資料探勘基本概念與技術。
- Understand and apply the fundamental concepts and technology of big data mining

# 課程大綱 (Syllabus)

週次 (Week)	日期 (Date)	內容 (Subject/Topics)
1	2019/02/20	巨量資料探勘課程介紹 (Course Orientation for Big Data Mining)
2	2019/02/27	AI人工智慧與大數據分析 (Artificial Intelligence and Big Data Analytics)
3	2019/03/06	分群分析 (Cluster Analysis)
4	2019/03/13	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	2019/03/20	關連分析 (Association Analysis)
6	2019/03/27	個案分析與實作二 (SAS EM 關連分析) : Case Study 2 (Association Analysis using SAS EM)
7	2019/04/03	教學行政觀摩日 (Off-campus study)
8	2019/04/10	分類與預測 (Classification and Prediction)

# 課程大綱 (Syllabus)

週次 (Week)	日期 (Date)	內容 (Subject/Topics)
9	2019/04/17	期中報告 (Midterm Project Presentation)
10	2019/04/24	期中考試週 (Midterm Exam)
11	2019/05/01	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)
12	2019/05/08	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)
13	2019/05/15	機器學習與深度學習 (Machine Learning and Deep Learning)
14	2019/05/22	期末報告 (Final Project Presentation)
15	2019/05/29	畢業考試週 (Final Exam)

# 教學方法與評量方法

- 教學方法
  - 講述、討論、賞析、模擬、實作、問題解決
- 評量方法
  - 紙筆測驗、實作、報告、上課表現

# 教材課本

- 教材課本

- 講義 (Slides)

- 資料採礦運用：以SAS Enterprise Miner為工具，李淑娟，2015，SAS賽仕電腦軟體

- 參考書籍

- Big Data, Data Mining, and Machine Learning: Value Creation for Business Leaders and Practitioners, Jared Dean, Wiley, 2014
  - Data Science for Business: What you need to know about data mining and data-analytic thinking, Foster Provost and Tom Fawcett, O'Reilly, 2013
  - Applied Analytics Using SAS Enterprise Mining, Jim Georges, Jeff Thompson and Chip Wells, SAS, 2010
  - Data Mining: Concepts and Techniques, Third Edition, Jiawei Han, Micheline Kamber and Jian Pei, Morgan Kaufmann, 2011
  - Learning Data Mining with Python - Second Edition, Robert Layton, Packt Publishing, 2017

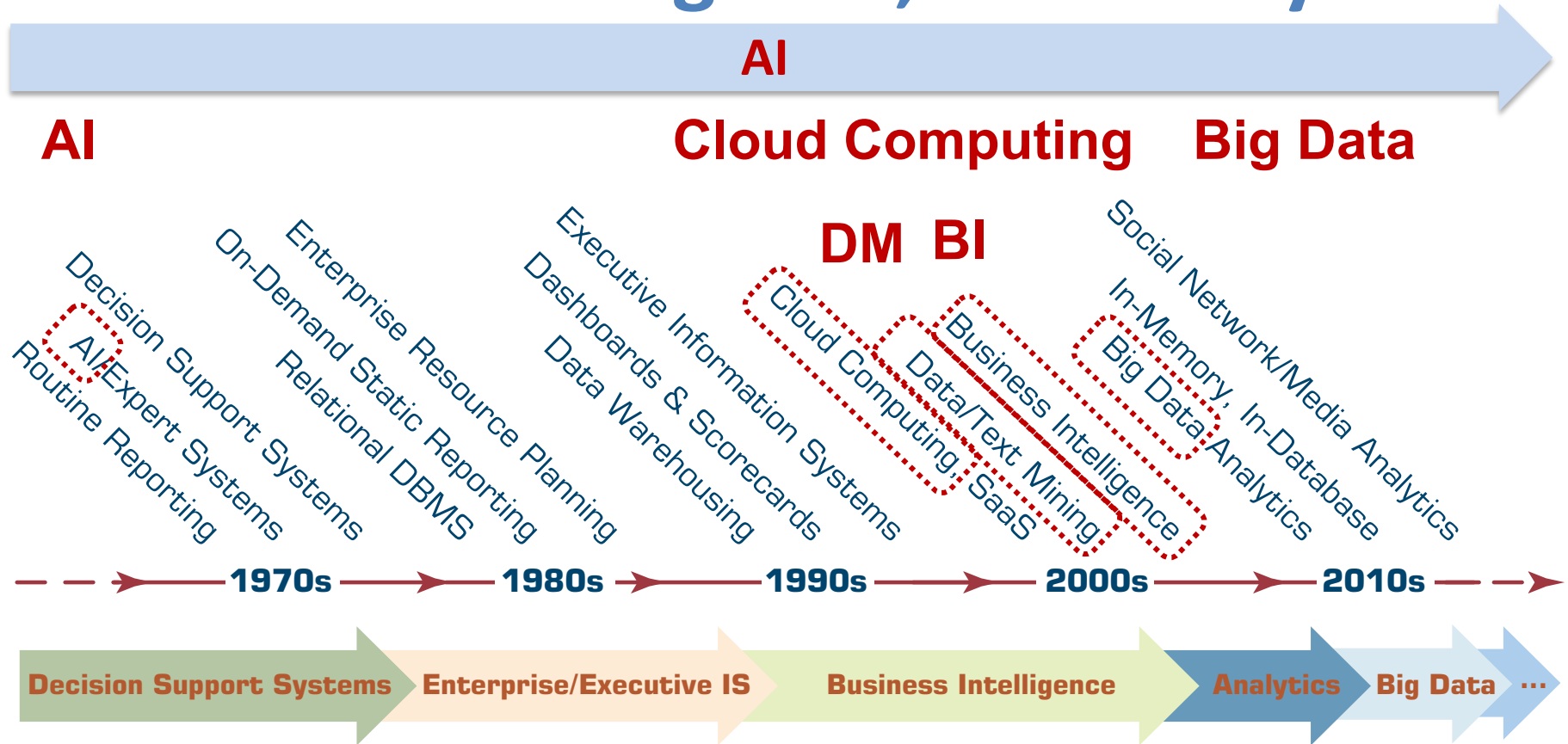
# 作業與學期成績計算方式

- 作業篇數
  - 3篇
- 學期成績計算方式
  - 期中評量：30 %
  - 期末評量：30 %
  - 其他（課堂參與及報告討論表現）：40 %

# Team Term Project

- Term Project Topics
  - Big Data mining
  - Big Data Analytics
  - Business Intelligence
  - FinTech
- 3-4 人為一組
  - 分組名單於 2019/02/27 (三) 課程下課時繳交
  - 由班代統一收集協調分組名單

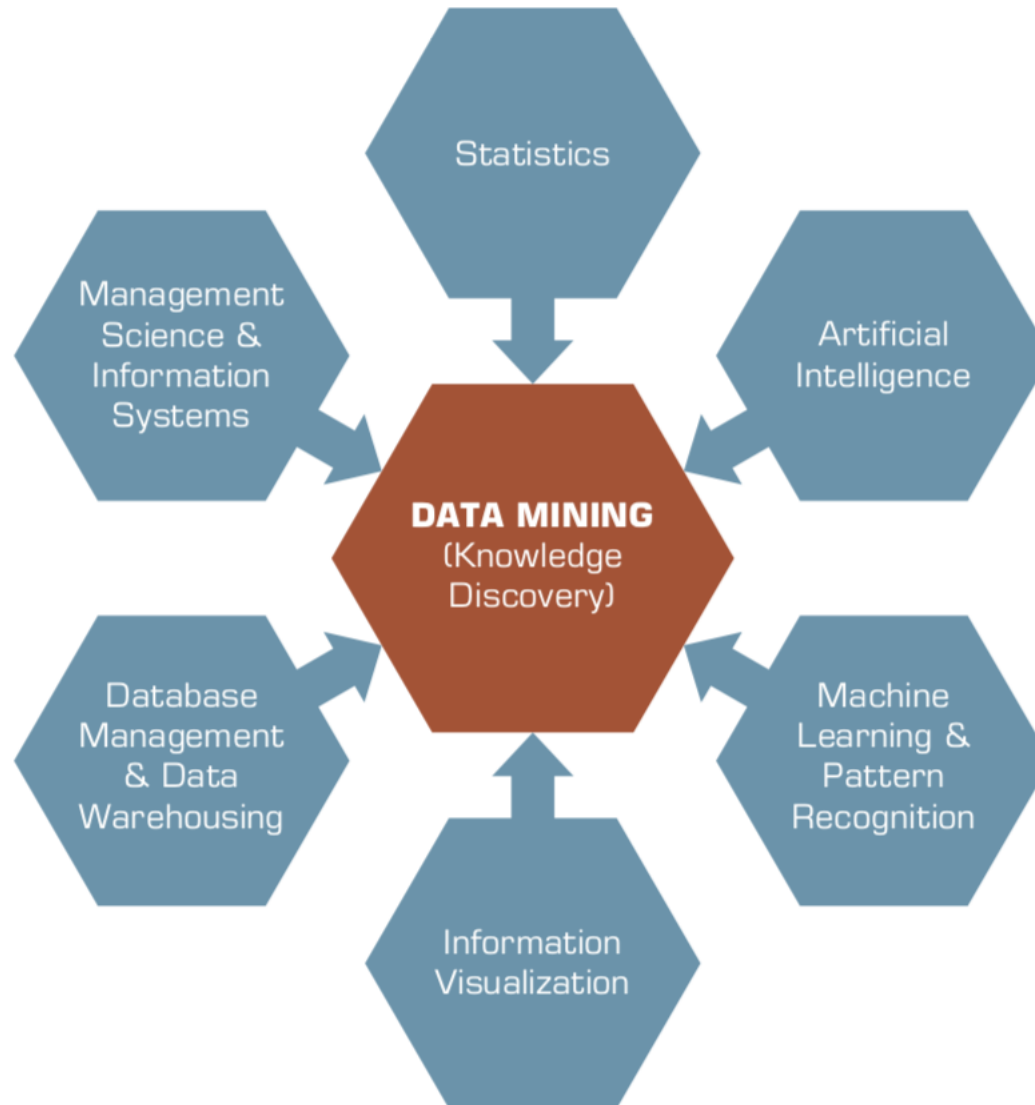
# AI, Big Data, Cloud Computing Evolution of Decision Support, Business Intelligence, and Analytics





# Data Mining

## Is a Blend of Multiple Disciplines

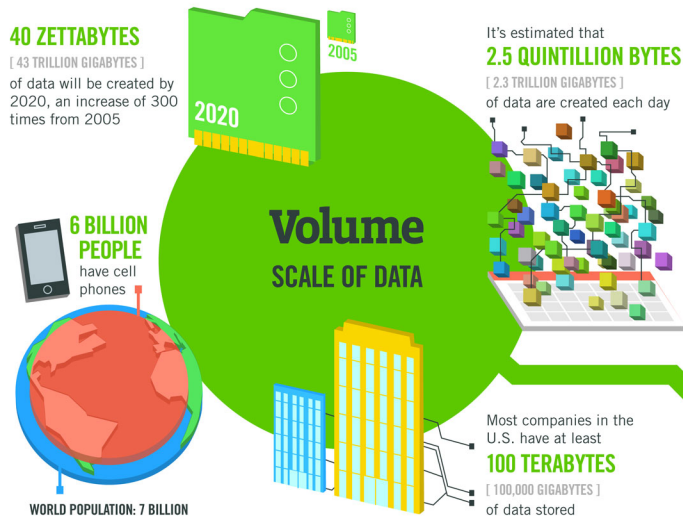


# Data Mining Tasks & Methods

Data Mining Tasks & Methods	Data Mining Algorithms	Learning Type
<b>Prediction</b>		
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
<b>Association</b>		
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
<b>Segmentation</b>		
Clustering	k-means, Expectation Maximization (EM)	Unsupervised
Outlier analysis	k-means, Expectation Maximization (EM)	Unsupervised

**Big Data**  
**Analytics**  
and  
**Data Mining**

# Big Data 4 V



## 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



As of 2011, the global size of data in healthcare was estimated to be

**150 EXABYTES**  
[161 BILLION GIGABYTES]

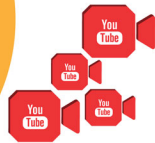


**30 BILLION PIECES OF CONTENT** are shared on Facebook every month



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



**Variety**  
DIFFERENT FORMS OF 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

**Velocity**  
ANALYSIS OF STREAMING DATA

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



**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**

**Veracity**  
UNCERTAINTY OF DATA

in one survey were unsure of how much of their data was inaccurate

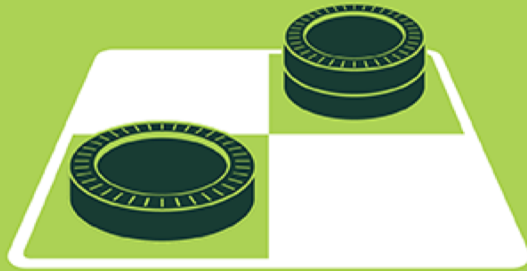
**value**

# 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

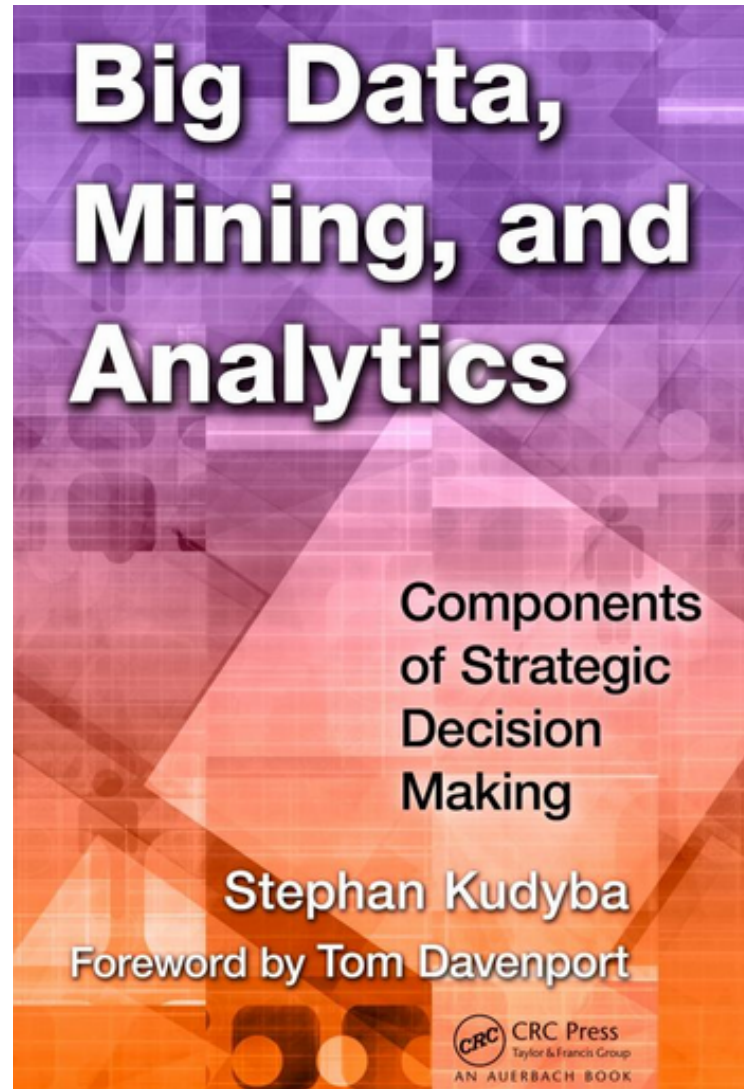
1990's

2000's

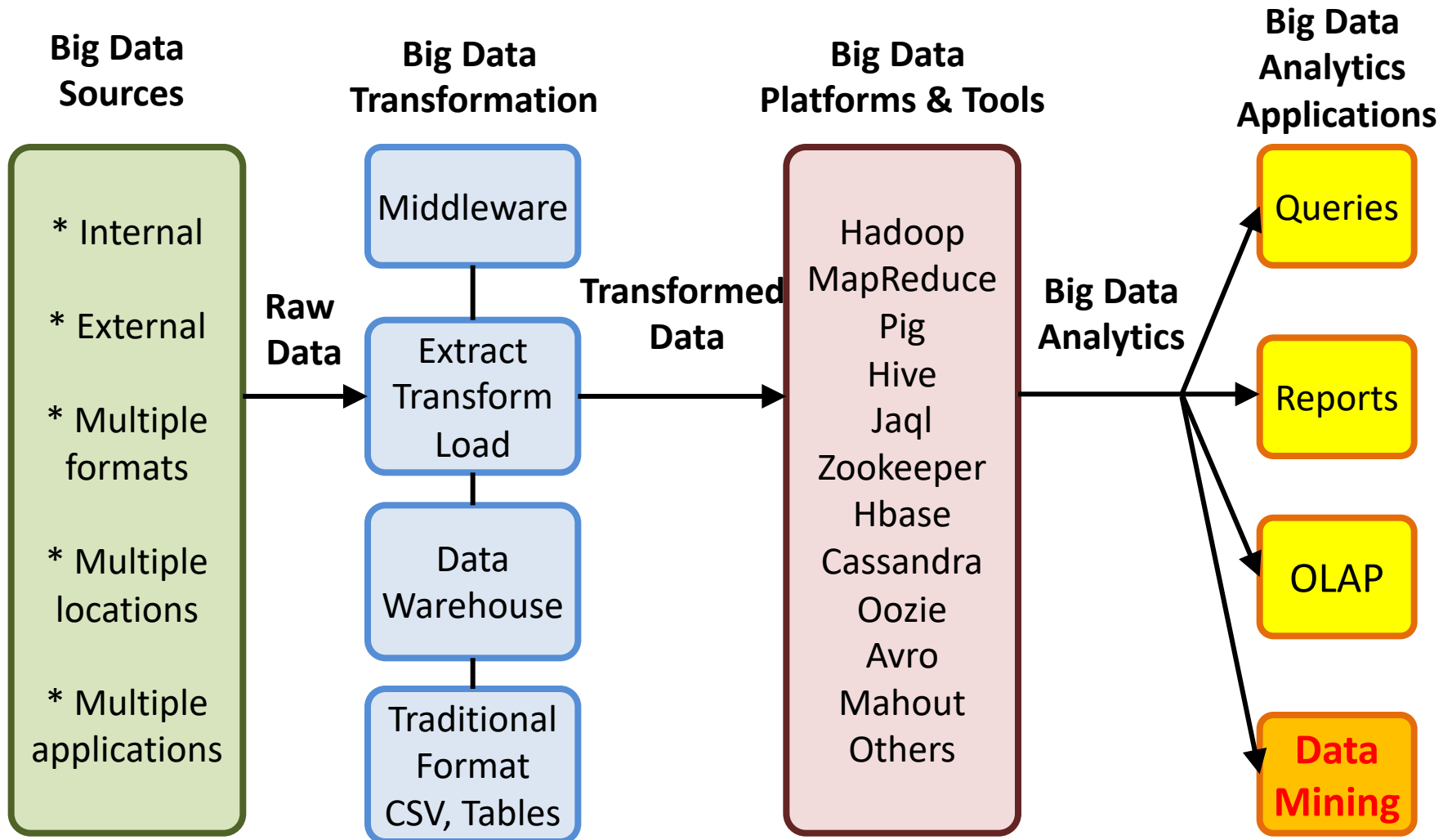
2010's

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.

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



# Architecture of Big Data Analytics



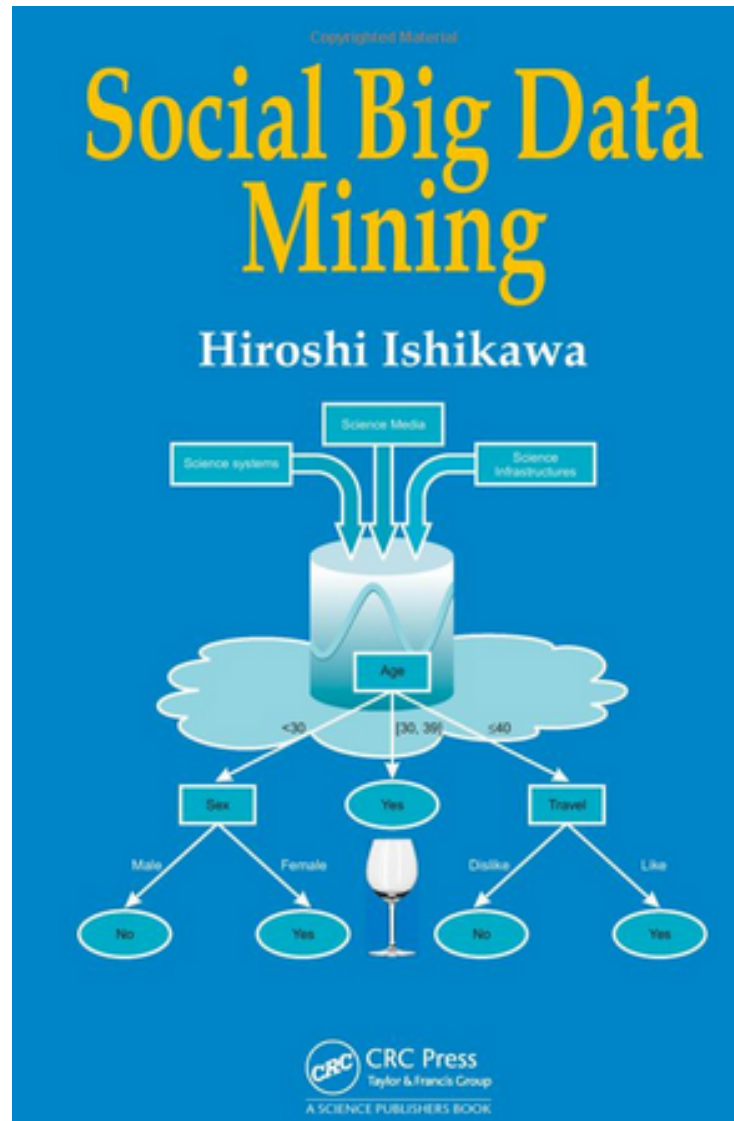


# Architecture of Big Data Analytics



# Social Big Data Mining

(Hiroshi Ishikawa, 2015)

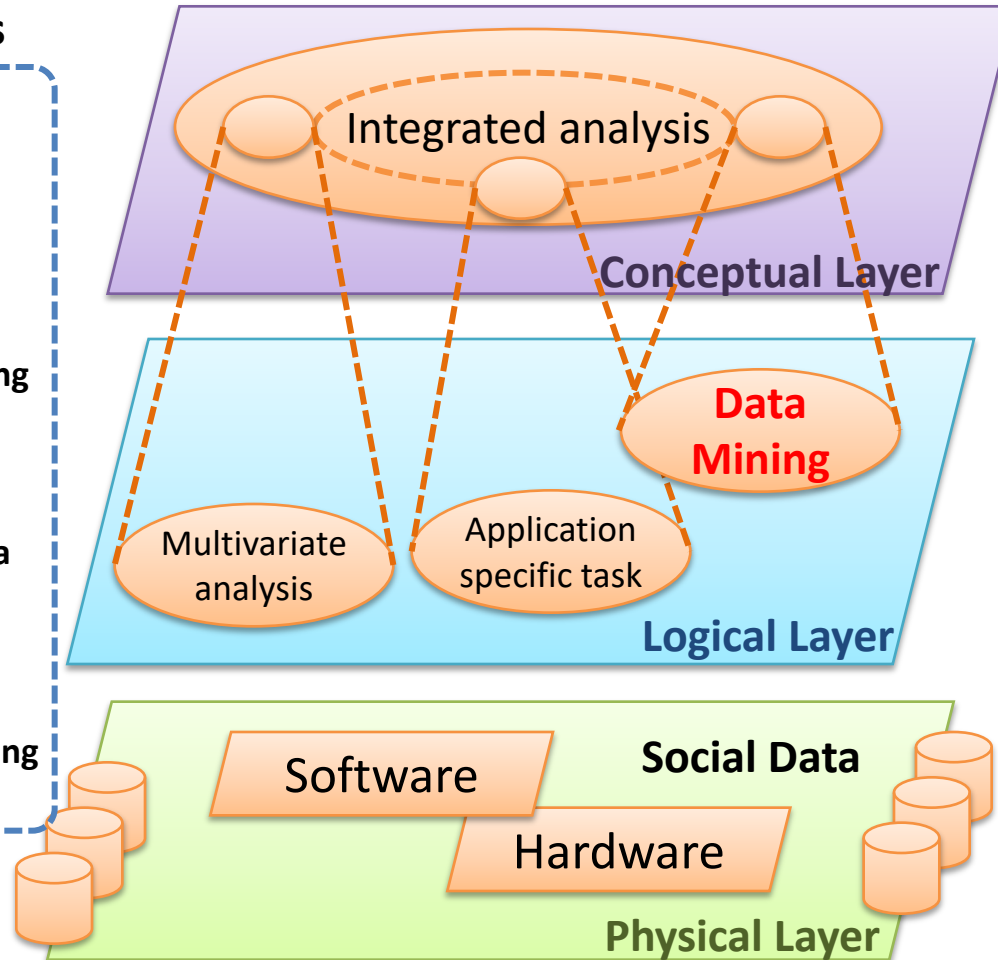


# Architecture for Social Big Data Mining

(Hiroshi Ishikawa, 2015)

## Enabling Technologies

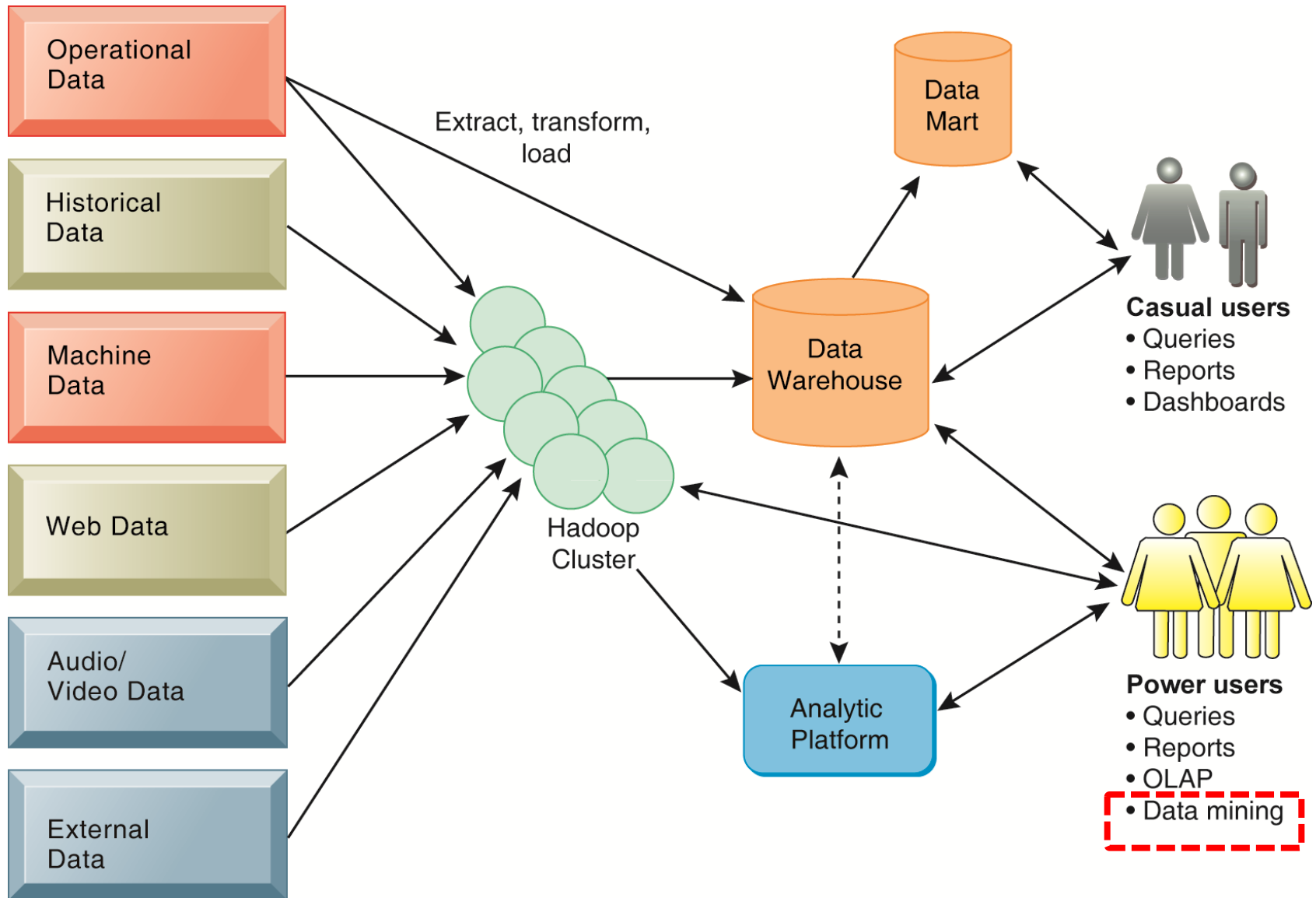
- Integrated analysis model
- Natural Language Processing
- Information Extraction
- Anomaly Detection
- Discovery of relationships among heterogeneous data
- Large-scale visualization
- Parallel distributed processing



## Analysts

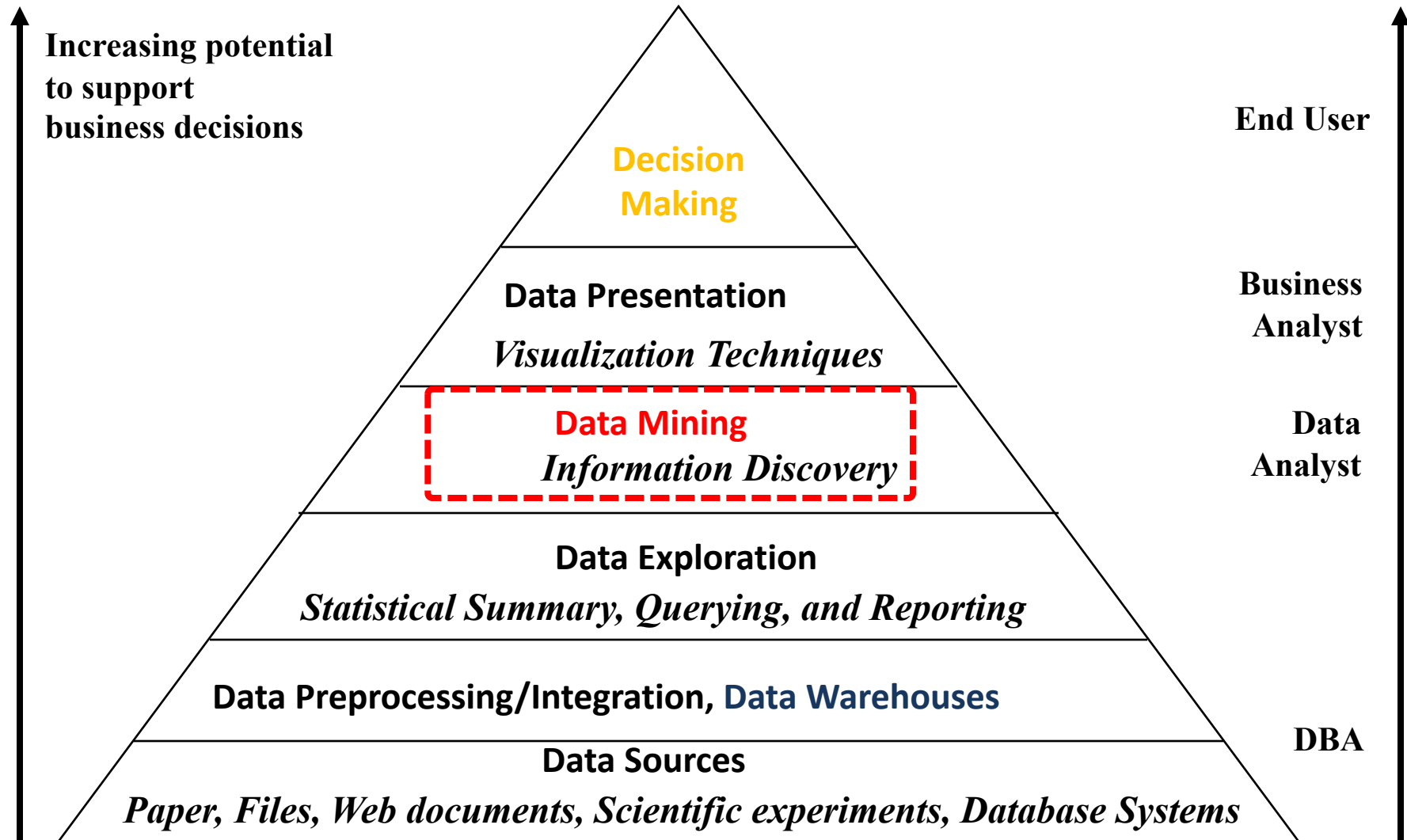
- Model Construction
- Explanation by Model
- Construction and confirmation of individual hypothesis
- Description and execution of application-specific task

# Business Intelligence (BI) Infrastructure



# Data Warehouse

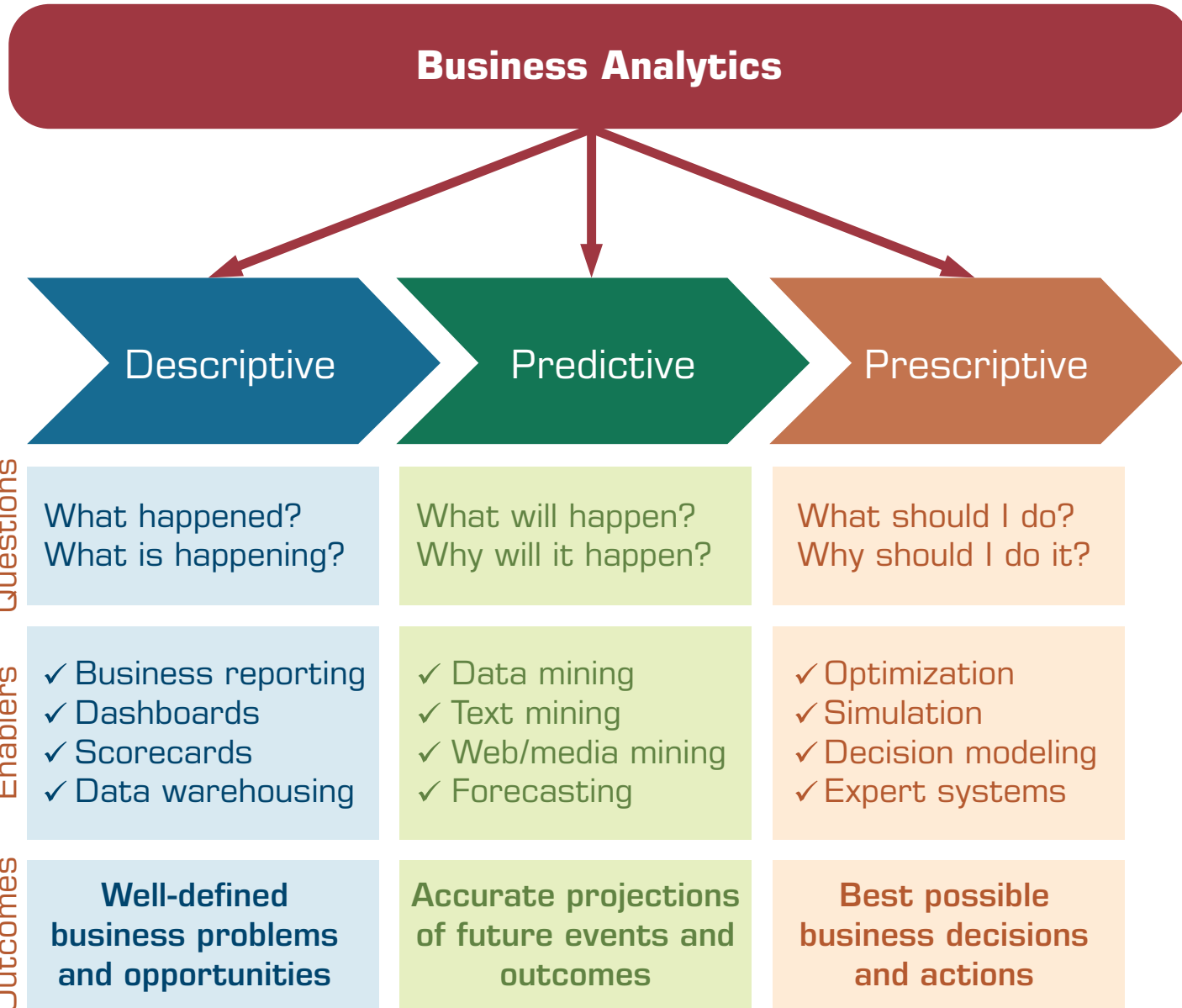
## Data Mining and Business Intelligence



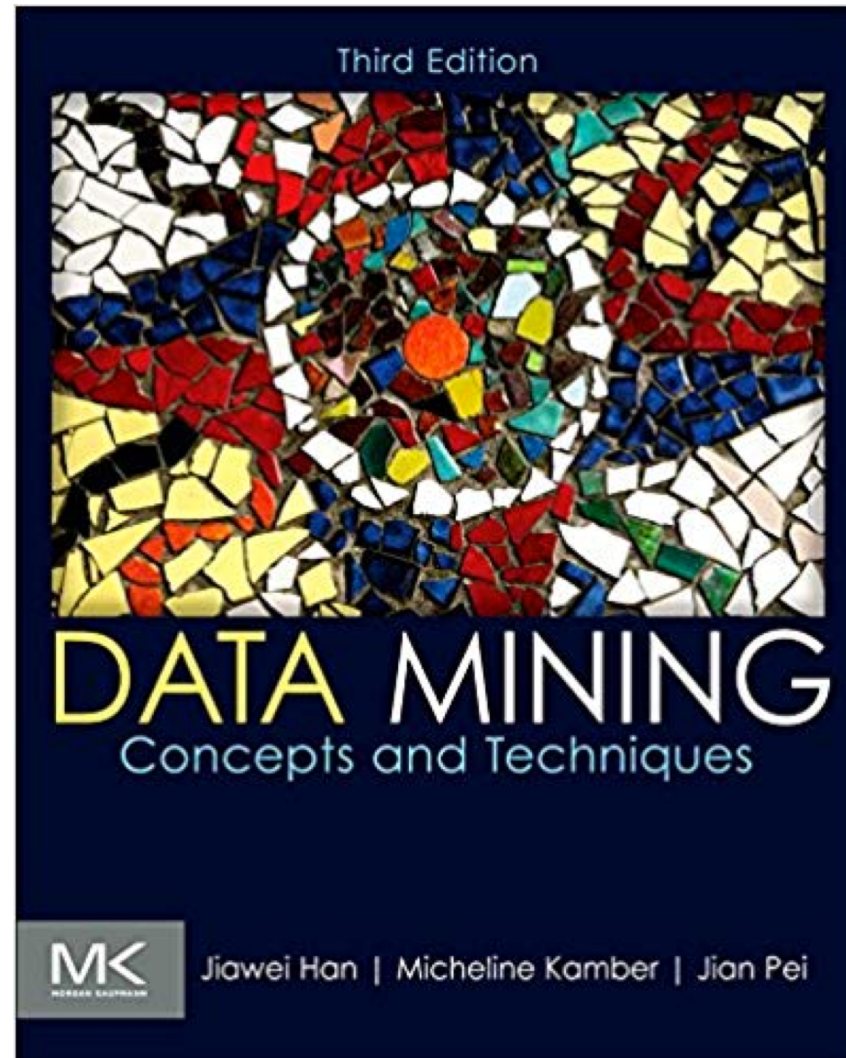
# The Evolution of BI Capabilities



# Three Types of Analytics



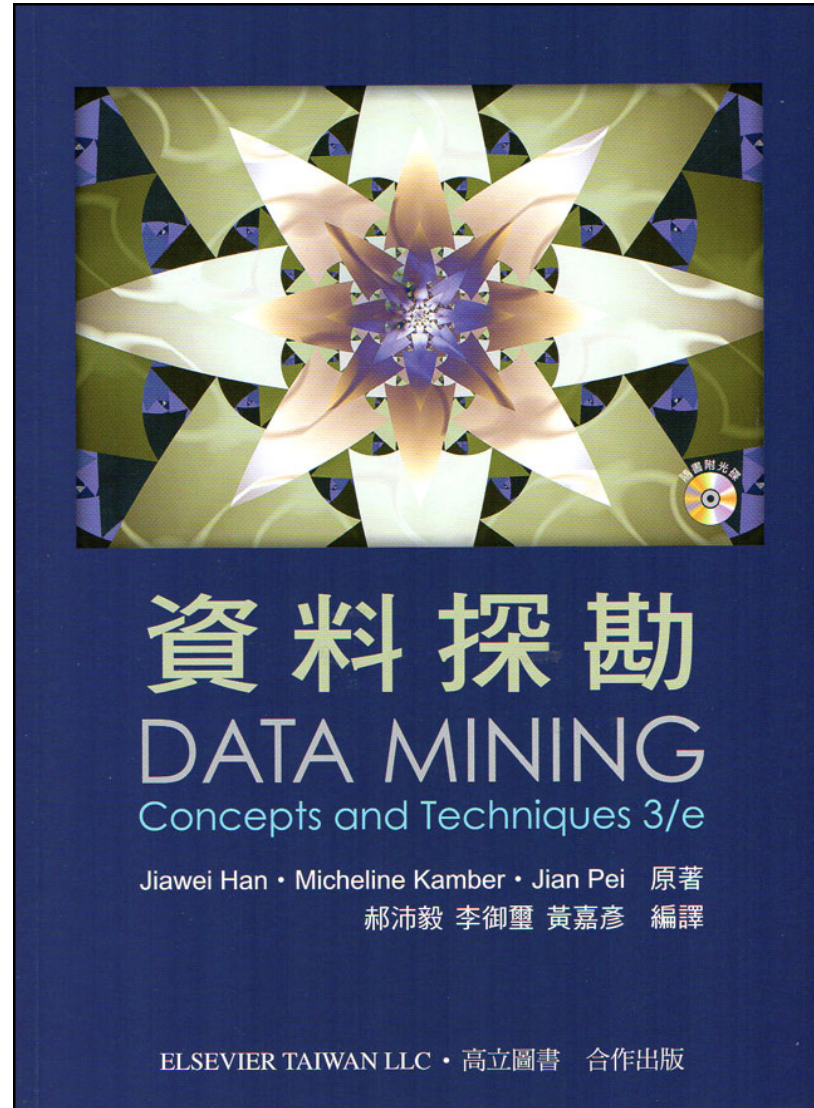
**Data Mining: Concepts and Techniques, Third Edition,  
Jiawei Han, Micheline Kamber and Jian Pei,  
Morgan Kaufmann, 2011**



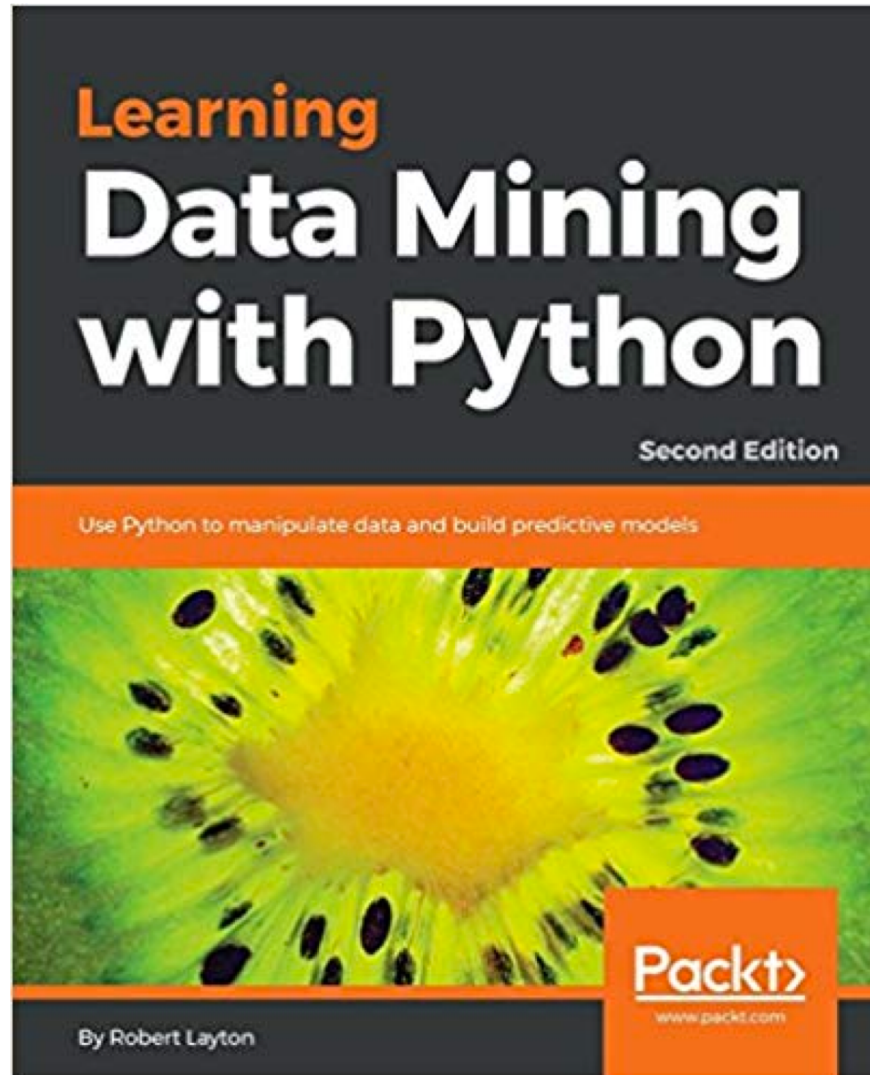


# 郝沛毅, 李御璽, 黃嘉彥 編譯, 資料探勘

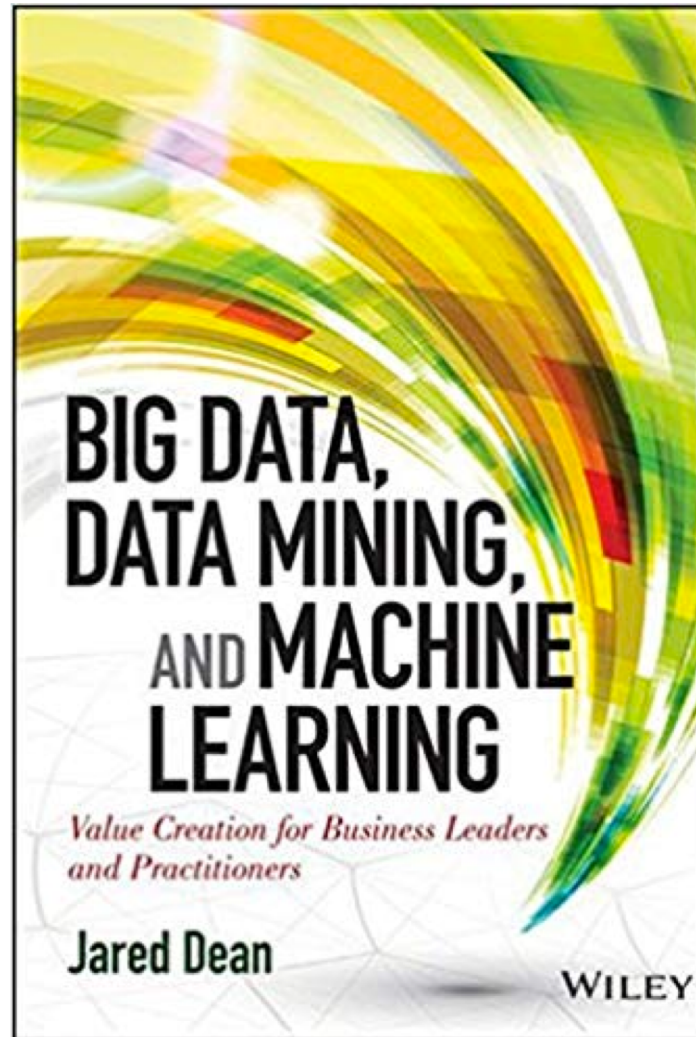
(Jiawei Han, Micheline Kamber, Jian Pei, Data Mining - Concepts and Techniques 3/e),  
高立圖書, 2014



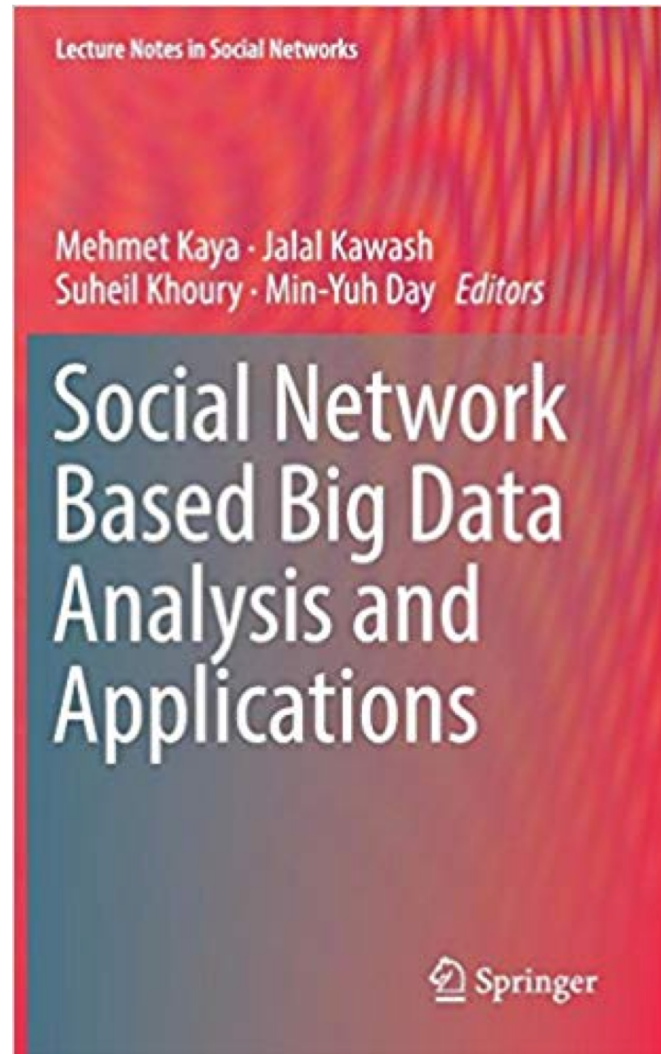
Learning Data Mining with Python - Second Edition,  
Robert Layton,  
Packt Publishing, 2017



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

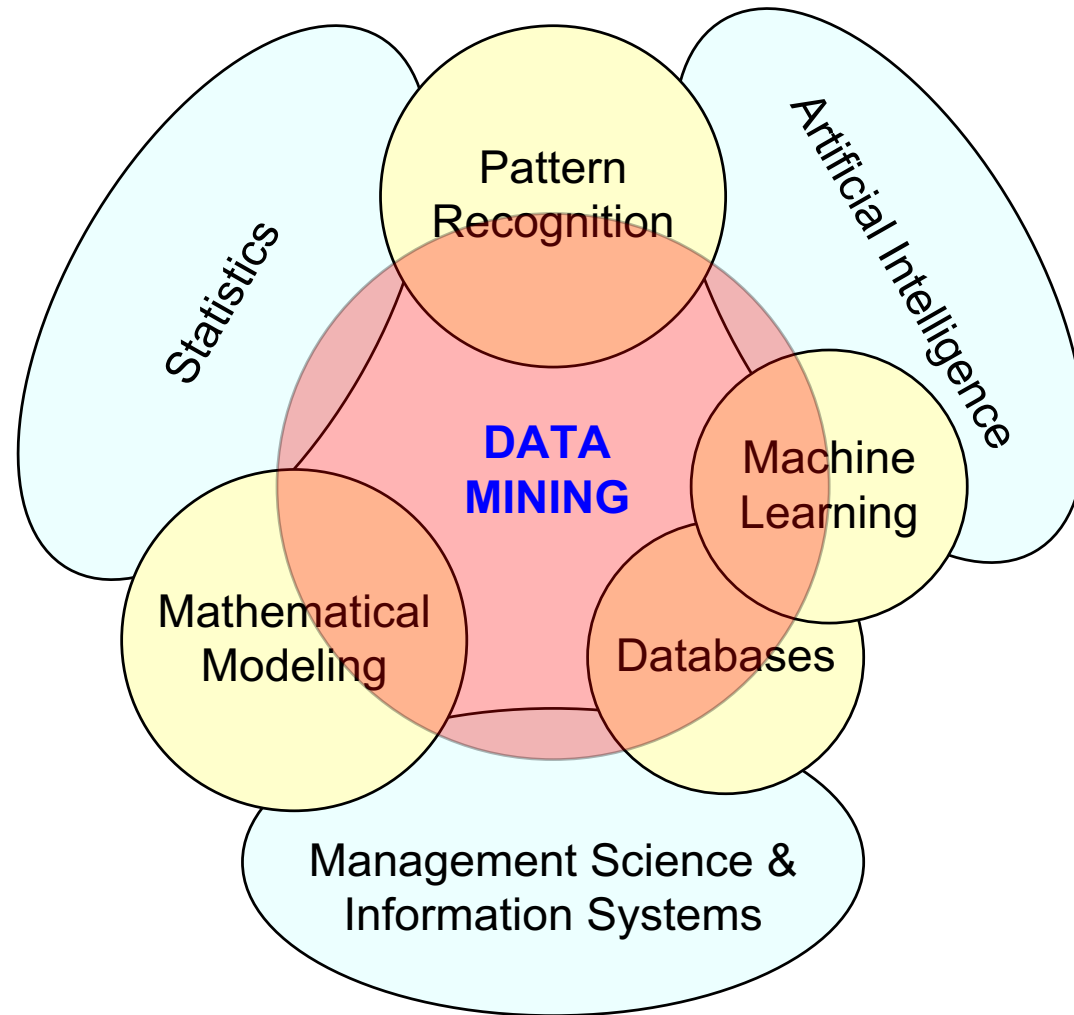


**Social Network Based Big Data Analysis and Applications,  
Lecture Notes in Social Networks,  
Mehmet Kaya, Jalal Kawash, Suheil Khoury, Min-Yuh Day,  
Springer International Publishing, 2018.**





# Data Mining at the Intersection of Many Disciplines





# Data Mining:

Core **Analytics** Process

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

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

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

Knowledge Discovery in Databases creates the context for developing the tools needed to control the flood of data facing organizations that depend on ever-growing databases of business, manufacturing, scientific, and personal information.

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

AS WE MARCH INTO THE AGE of digital information, the problem of data overload looms ominously ahead. Our ability to analyze and understand massive datasets lags far behind our ability to gather and store the data. A new generation of computational techniques and tools is required to support the extraction of useful knowledge from the rapidly growing volumes of data. These techniques and tools are the subject of the emerging field of knowledge discovery in databases (KDD) and data mining.

Large databases of digital information are ubiquitous. Data from the neighborhood store's checkout register, your bank's credit card authorization device, records in your doctor's office, patterns in your telephone calls,

Usama Fayyad,  
Gregory Piatetsky-Shapiro,  
and Padhraic Smyth

and many more applications generate streams of digital records archived in huge databases, sometimes in so-called data warehouses.

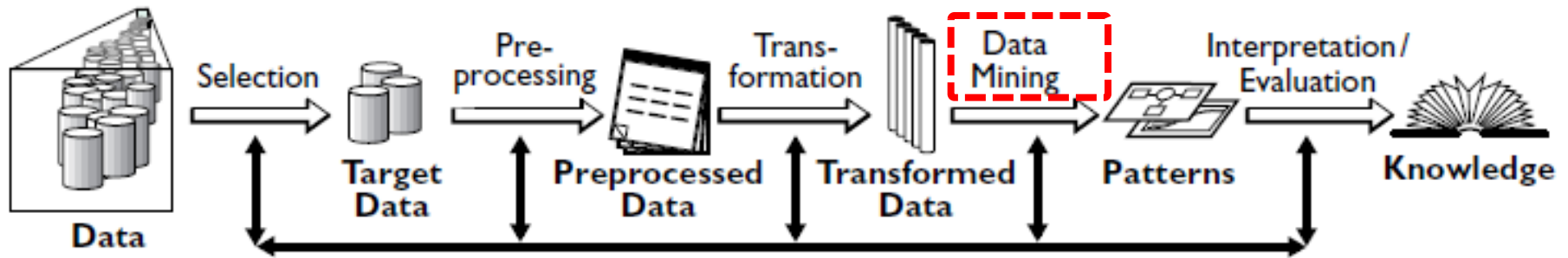
Current hardware and database technology allow efficient and inexpensive reliable data storage and access. However, whether the context is business, medicine, science, or government, the datasets themselves (in raw form) are of little direct value. What is of value is the knowledge that can be inferred from the data and put to use. For example, the marketing database of a consumer



# Data Mining

## Knowledge Discovery in Databases (KDD) Process

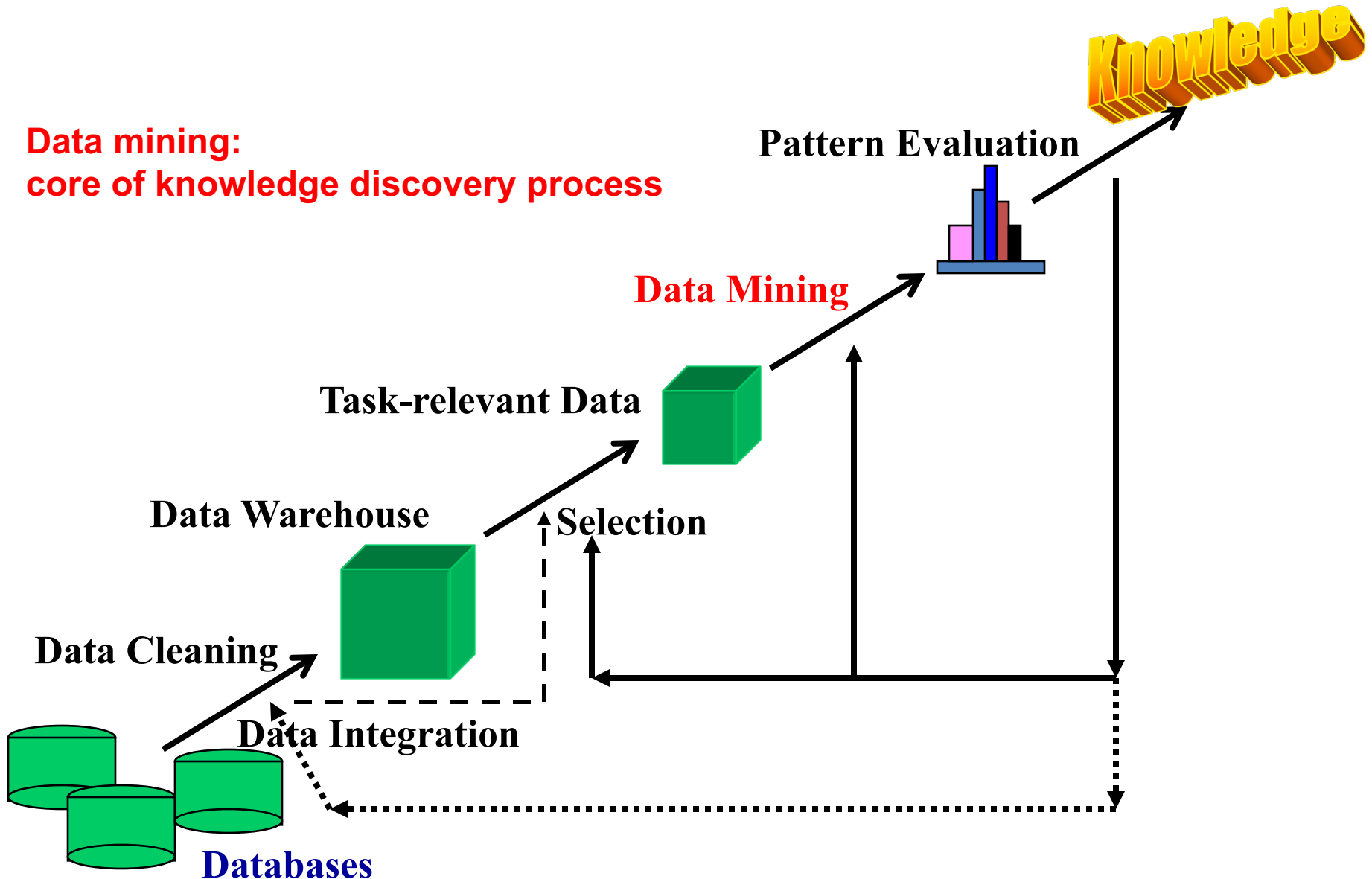
(Fayyad et al., 1996)





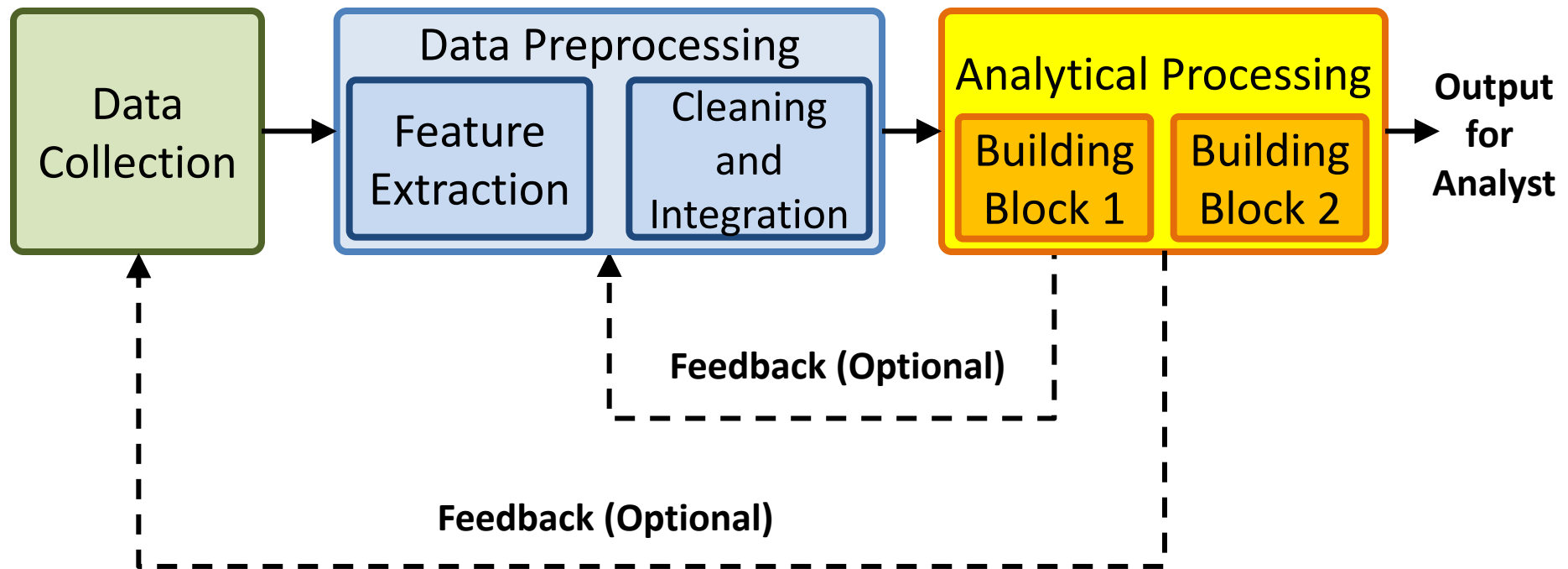
# Knowledge Discovery (KDD) Process

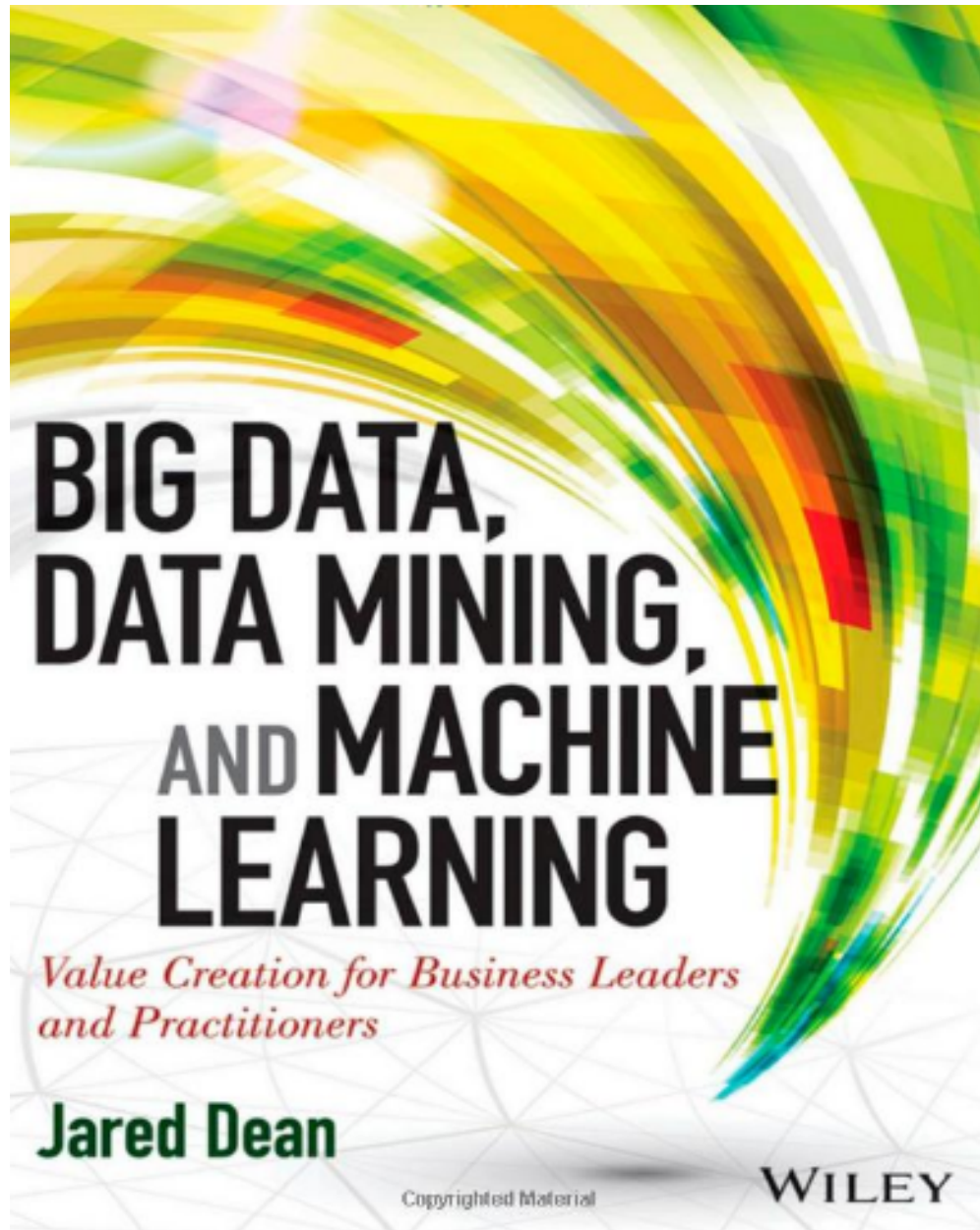
**Data mining:**  
core of knowledge discovery process



# Data Mining Processing Pipeline

(Charu Aggarwal, 2015)

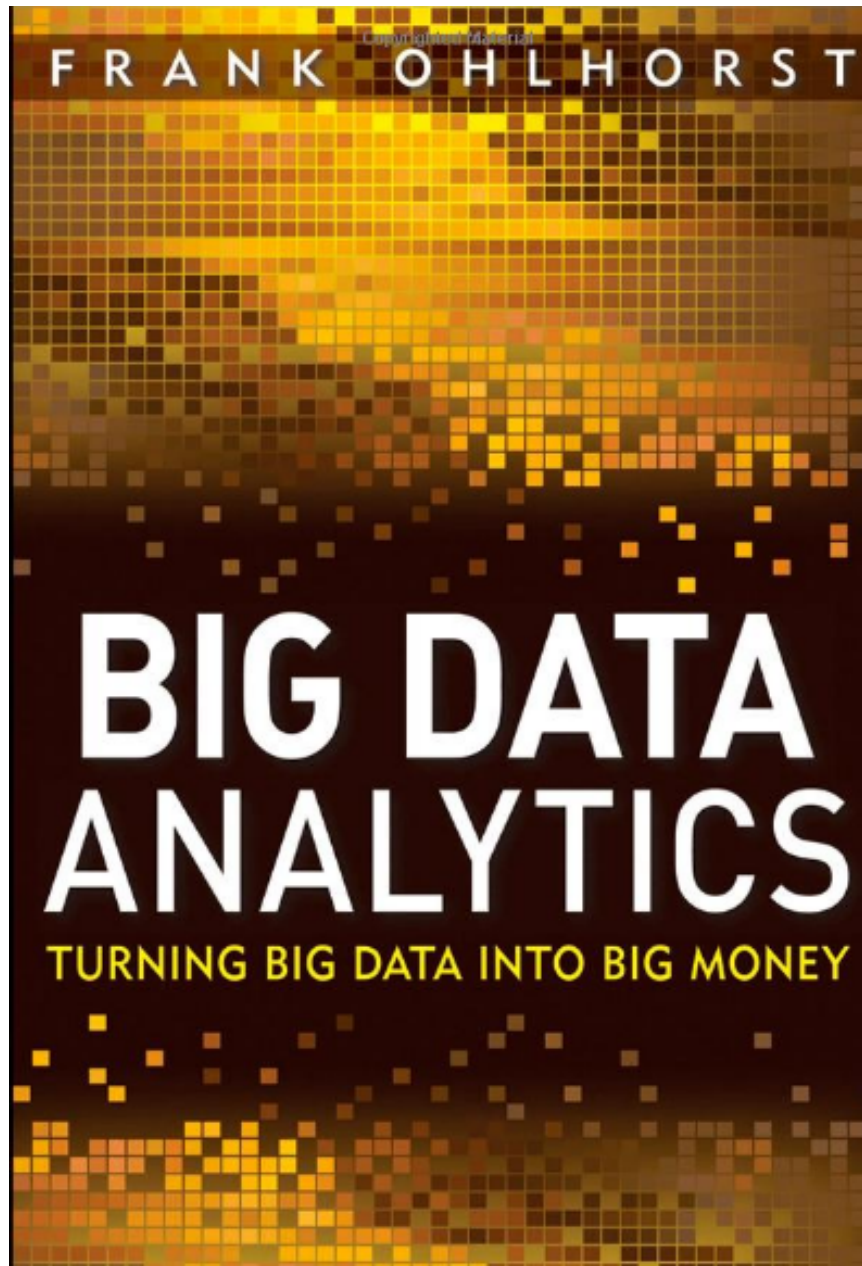




# Deep Learning

## Intelligence from Big Data









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<sup>2</sup>I<sup>2</sup>

DETECT

Anomalies  
Communities  
Typologies

PREDICT

Trending  
Real-time  
Prediction

DECIDE

Simulation  
Optimization



## BIG DATA – Batch



## BIG DATA – Real Time



Complex by nature



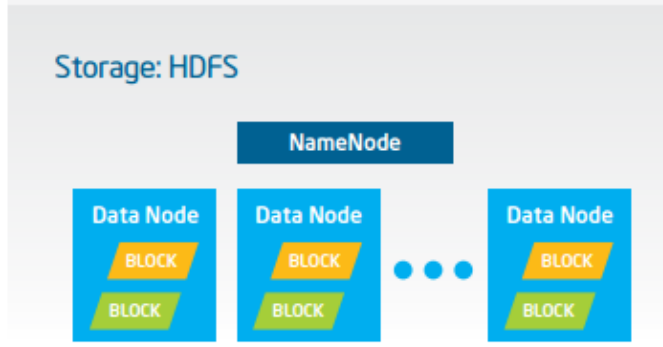
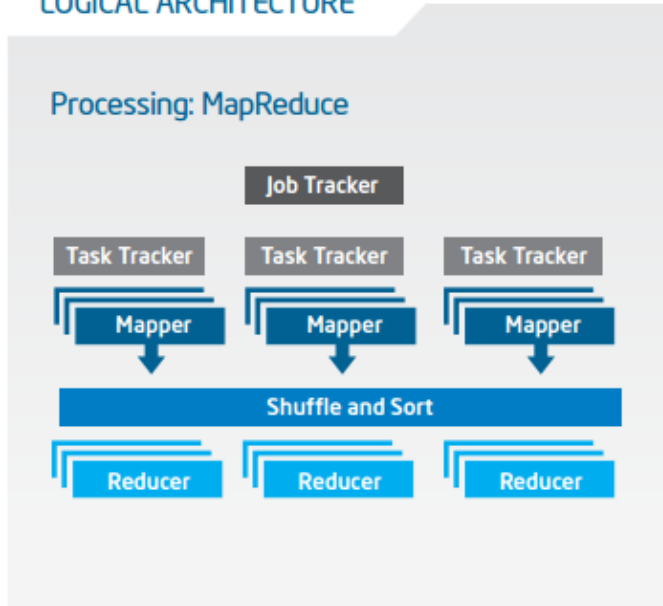
## DATA

Complex by structure

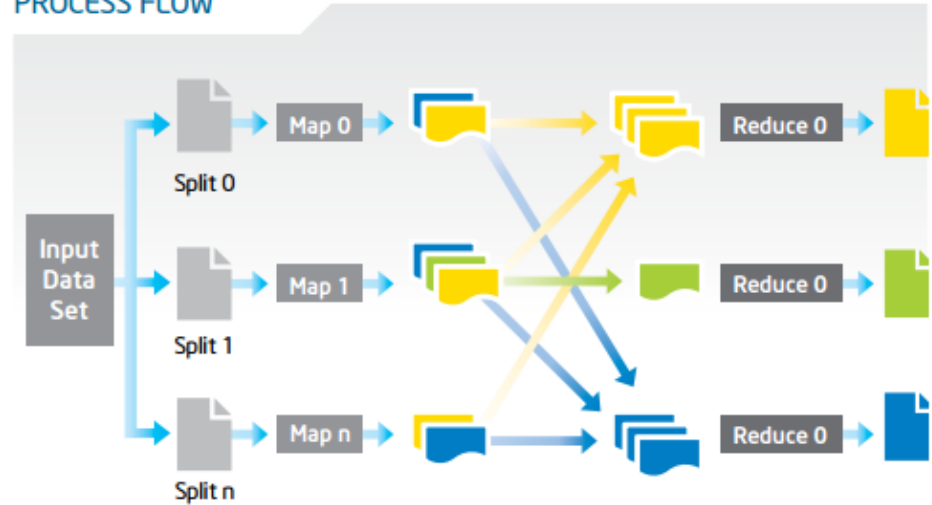


# Big Data with Hadoop Architecture

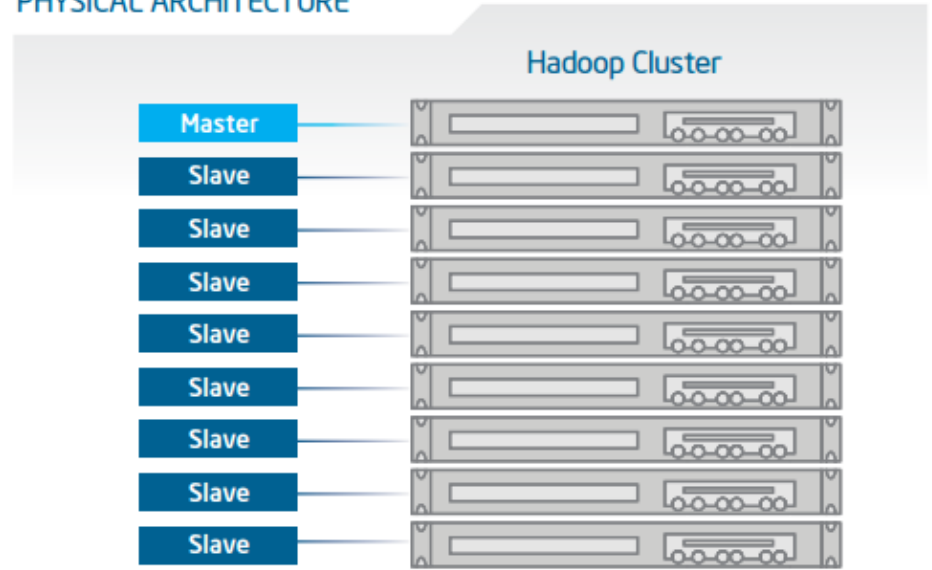
## LOGICAL ARCHITECTURE



## PROCESS FLOW



## PHYSICAL ARCHITECTURE

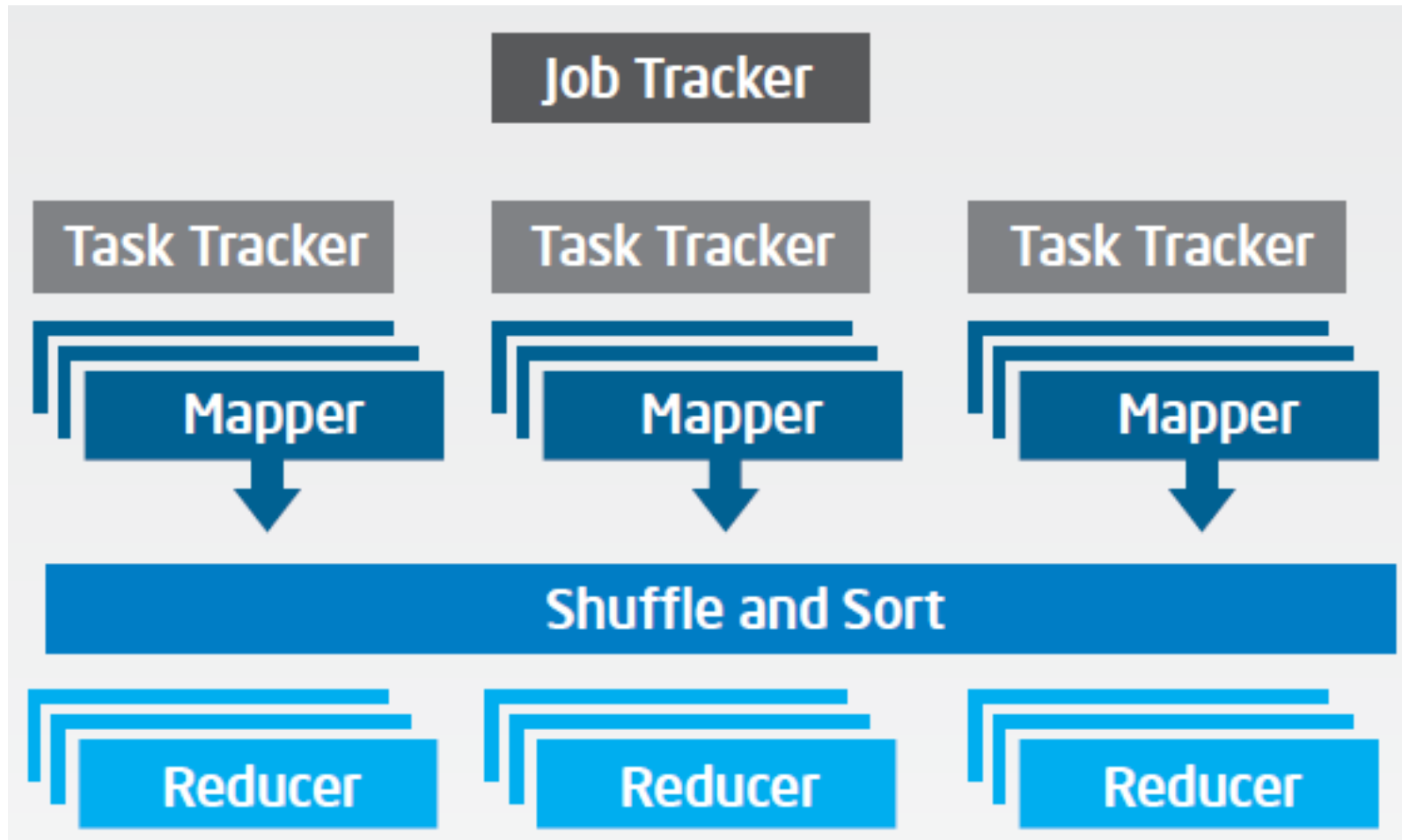




# Big Data with Hadoop Architecture

## Logical Architecture

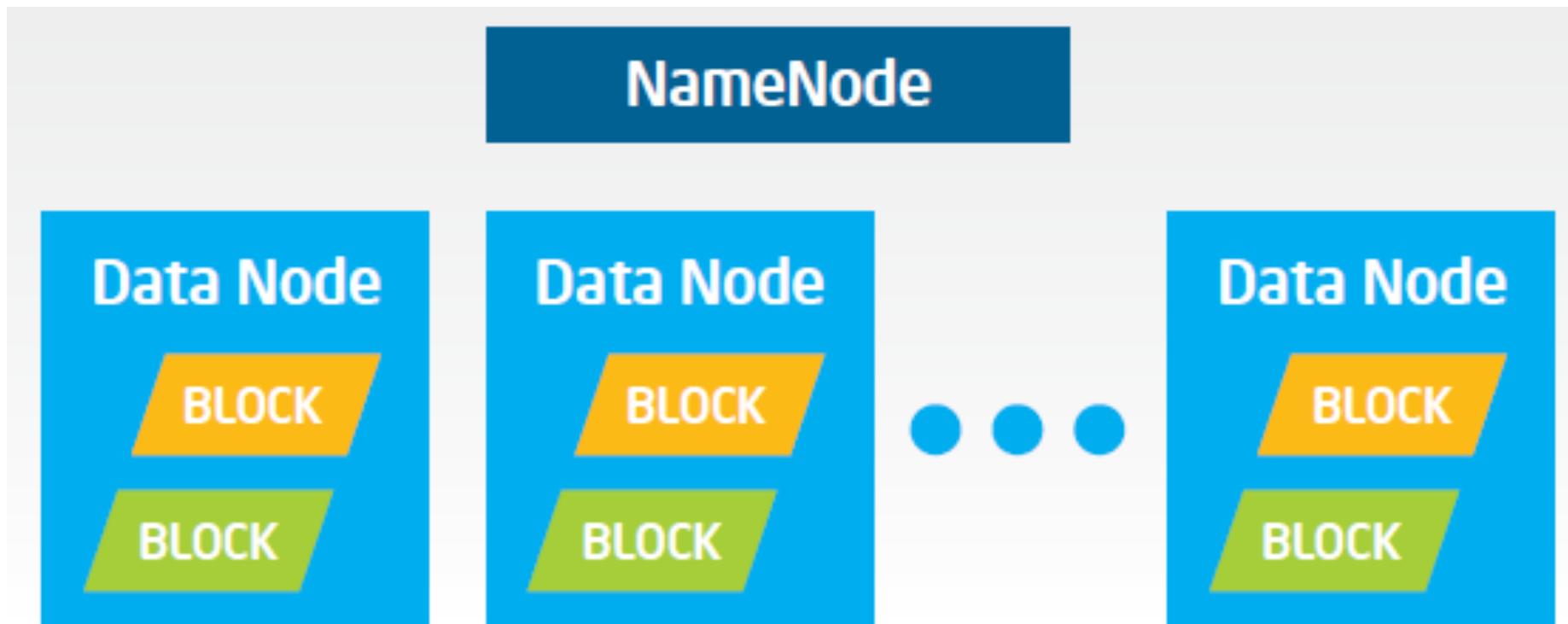
### Processing: MapReduce



# Big Data with Hadoop Architecture

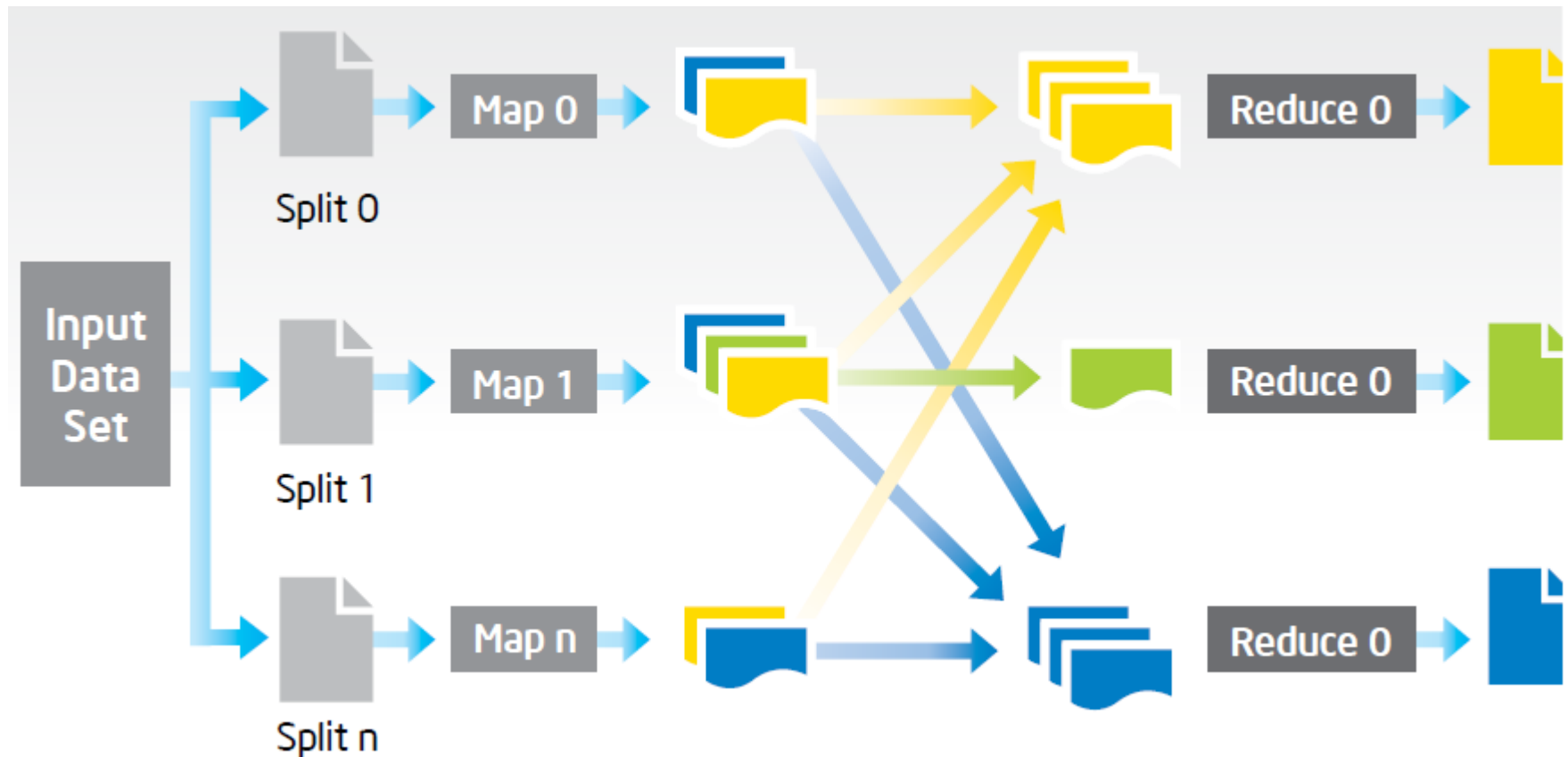
## Logical Architecture

Storage: HDFS



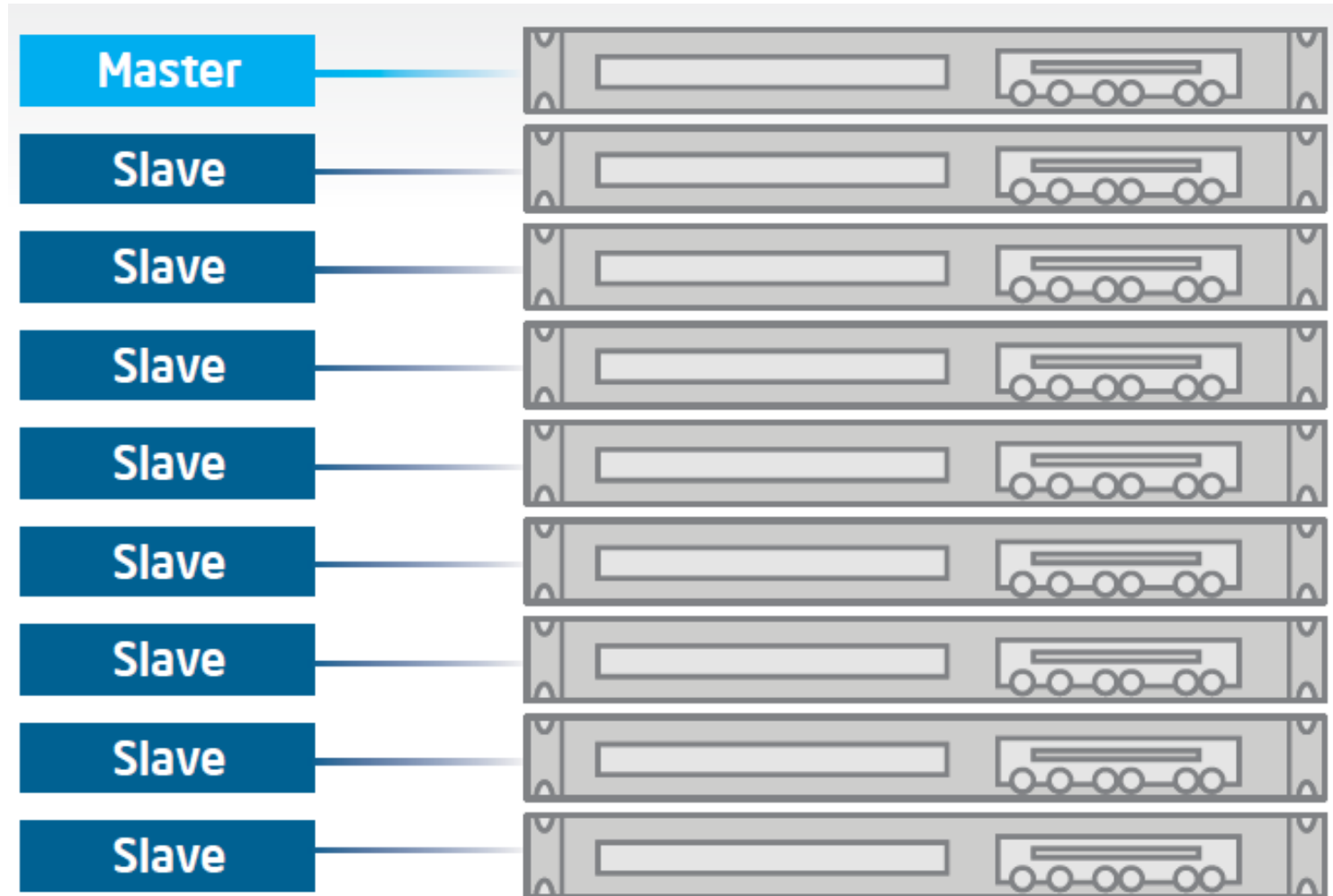
# Big Data with Hadoop Architecture

## Process Flow

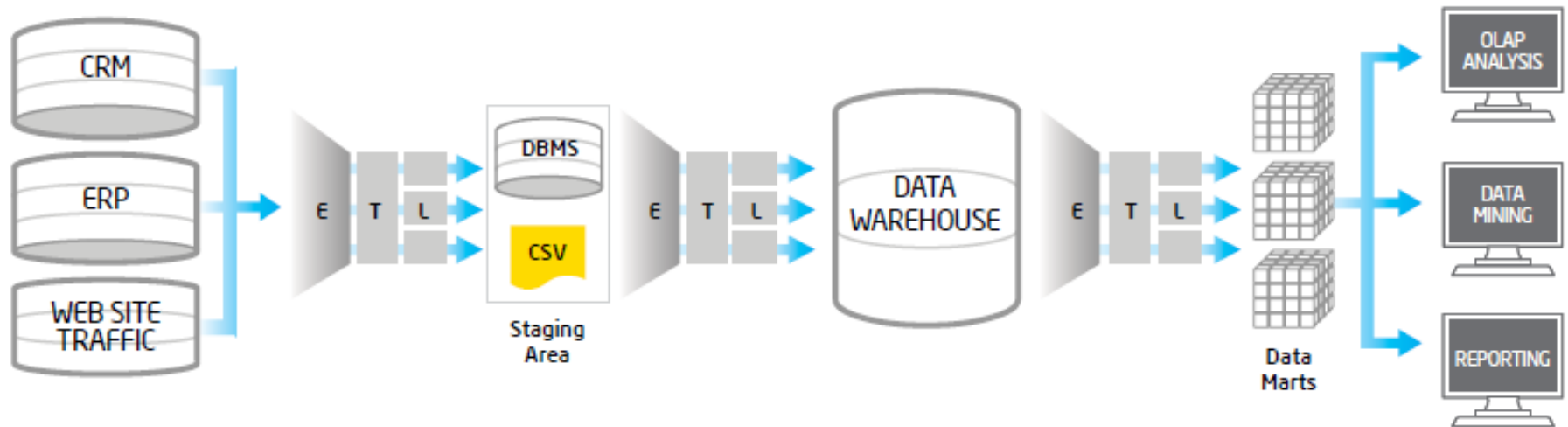


# Big Data with Hadoop Architecture

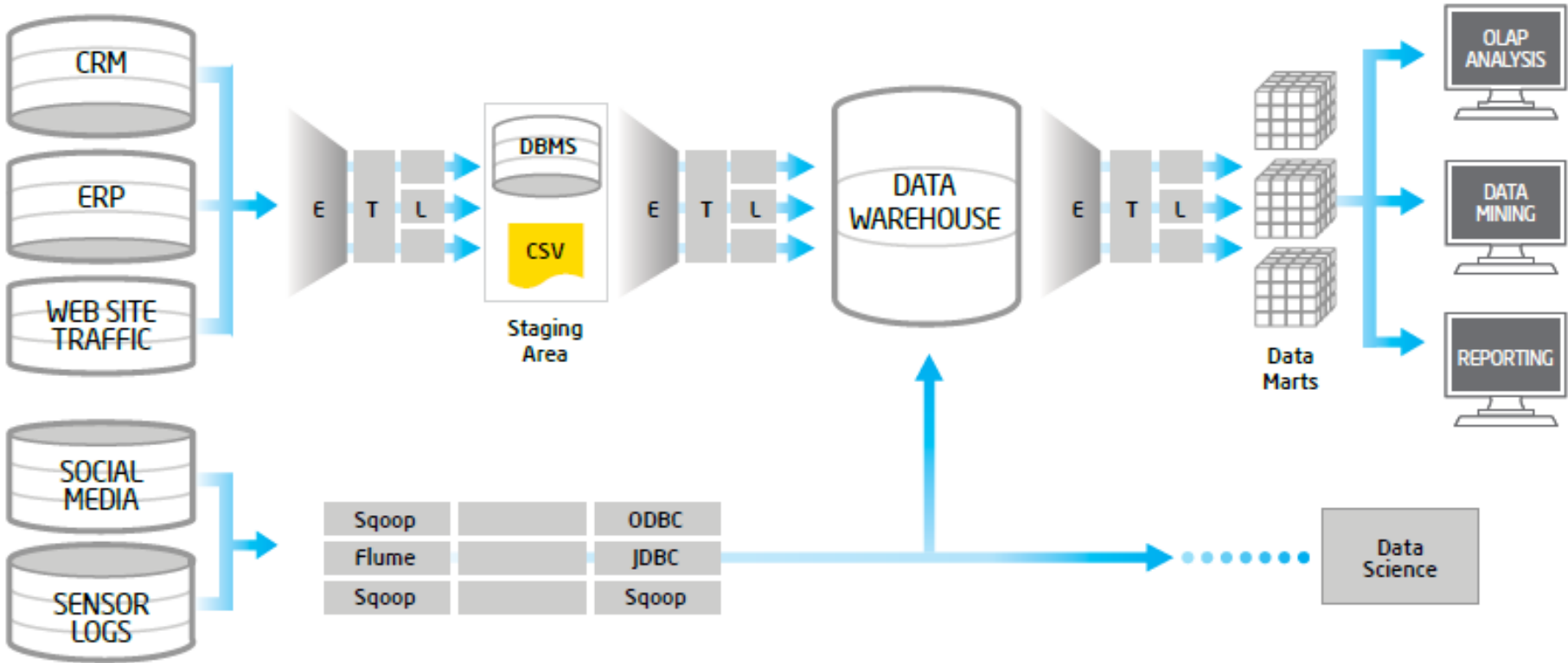
## Hadoop Cluster



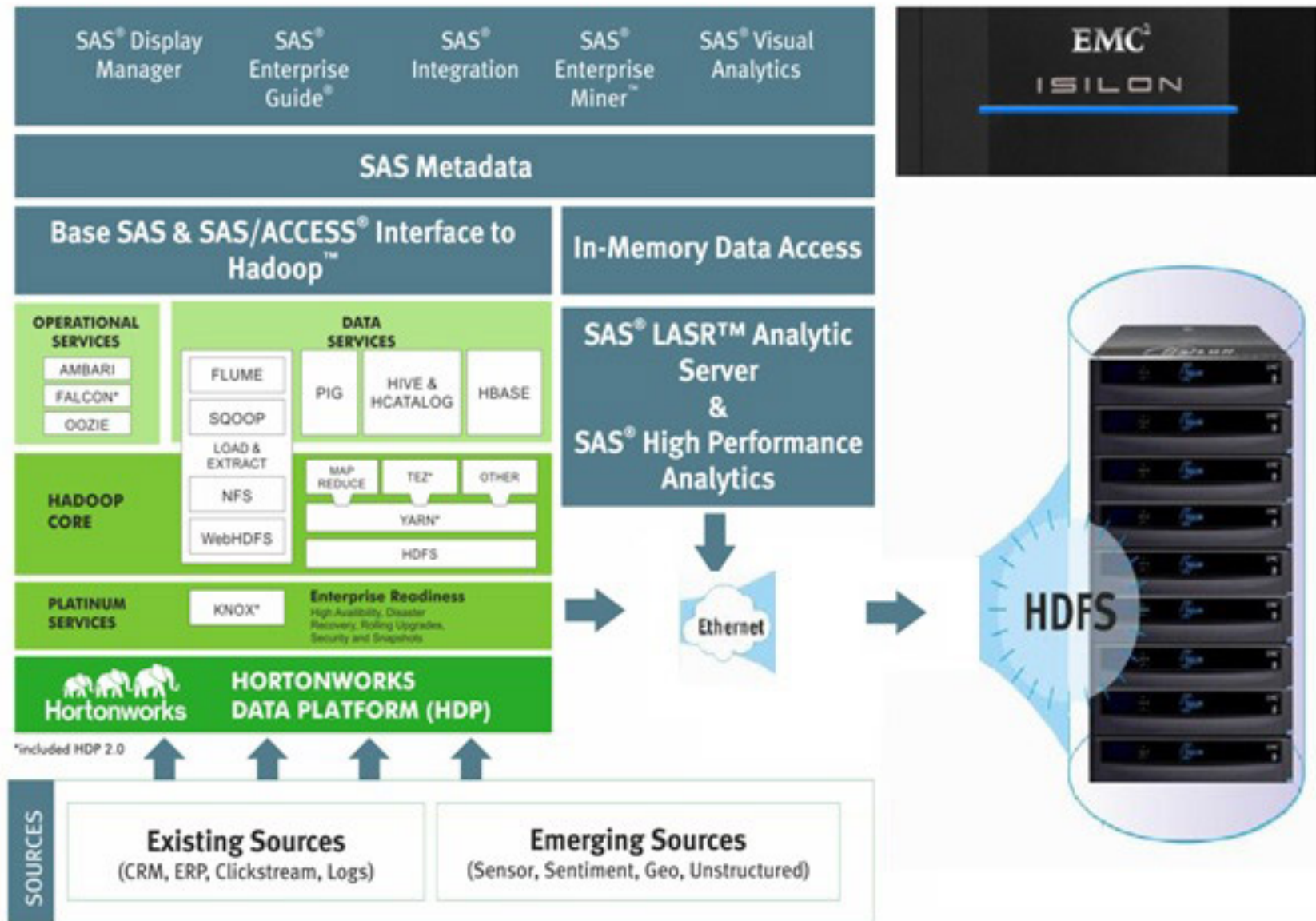
# Traditional ETL Architecture



# Offload ETL with Hadoop (Big Data Architecture)

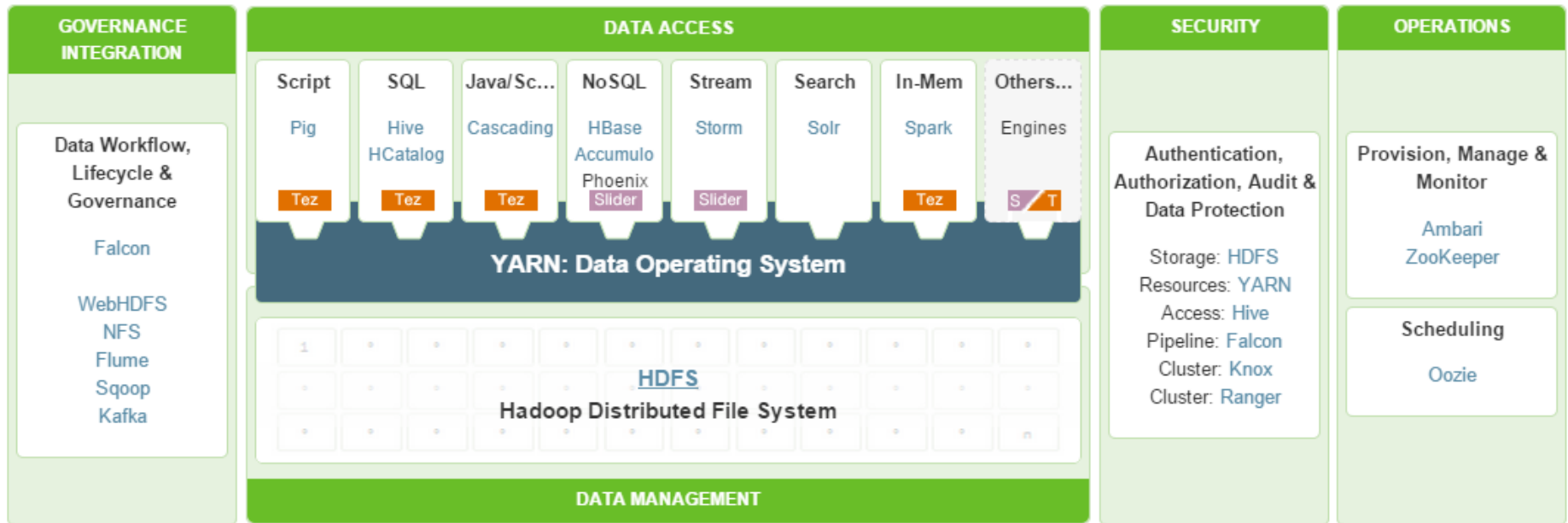


# Big Data Solution



# HDP

## A Complete Enterprise Hadoop Data Platform

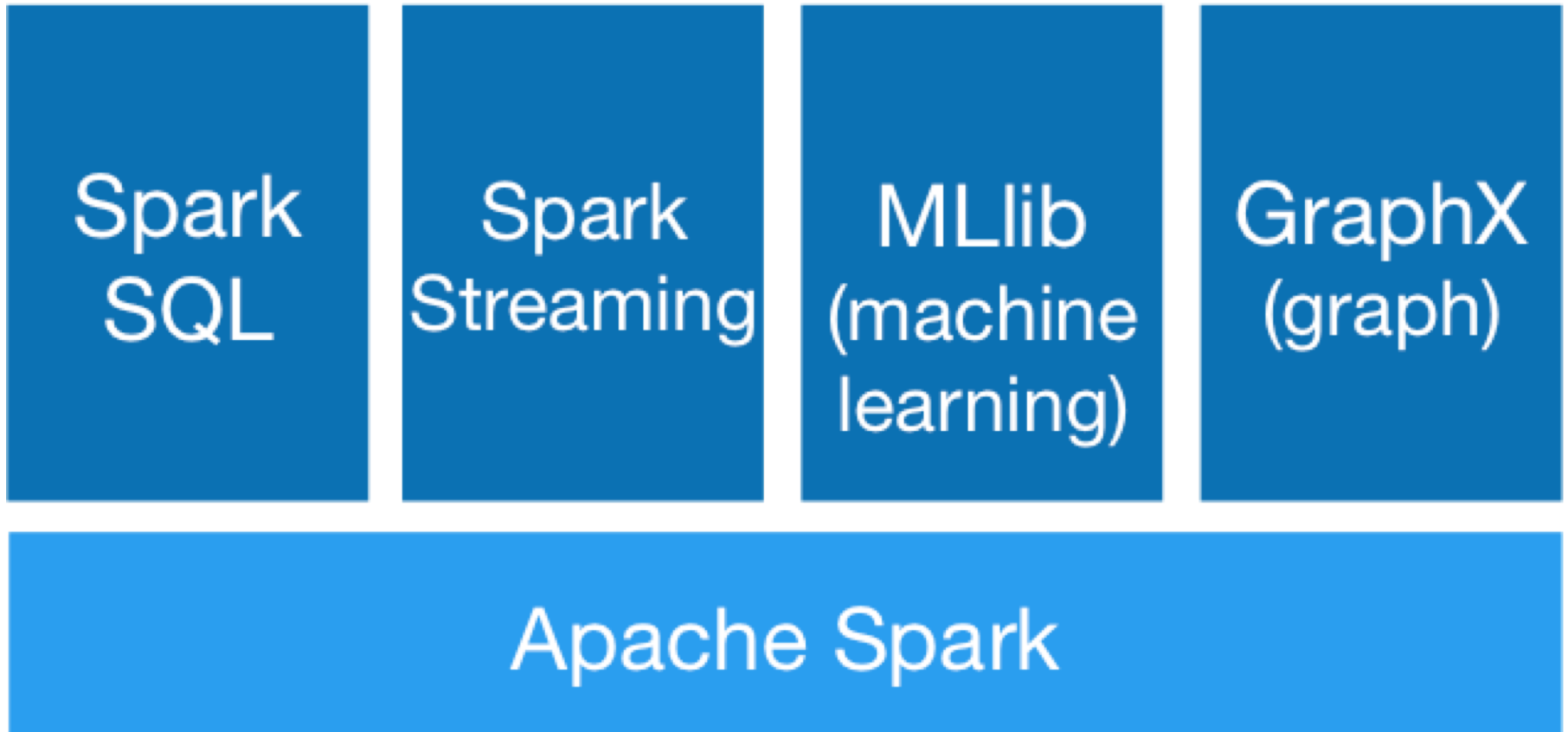




# Spark and Hadoop

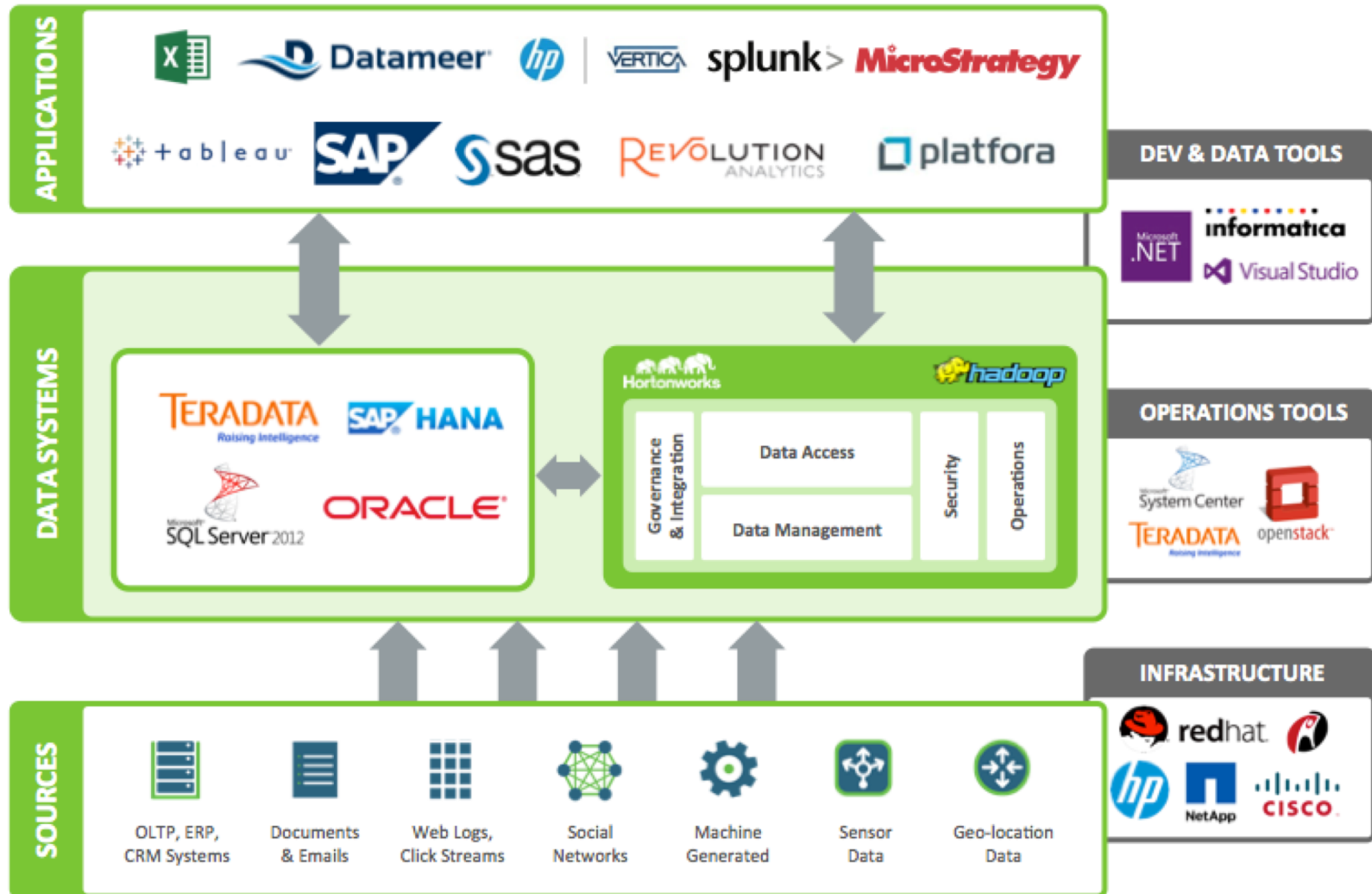


# Spark Ecosystem



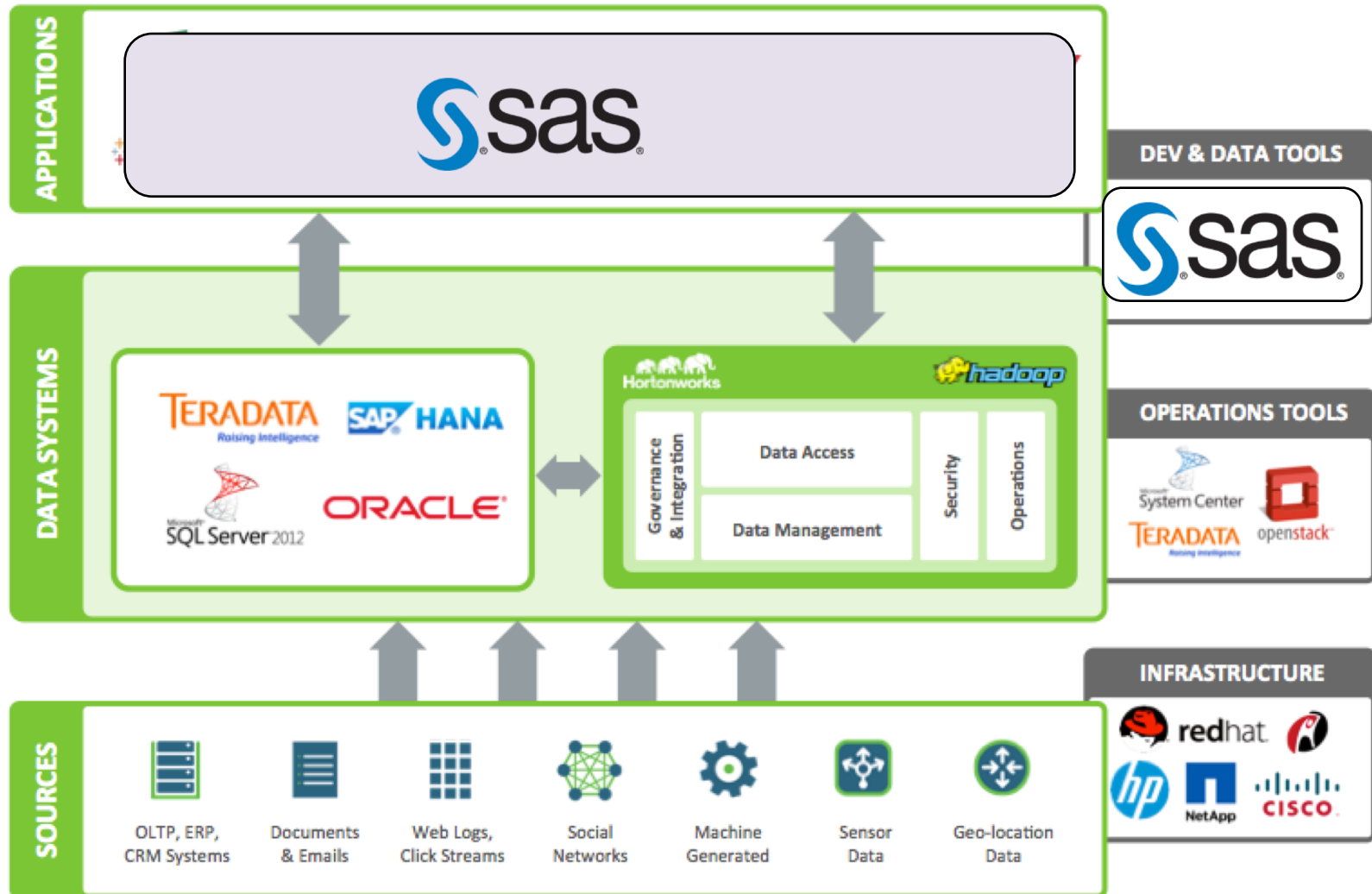
# SAS Big data Strategy

## – SAS areas

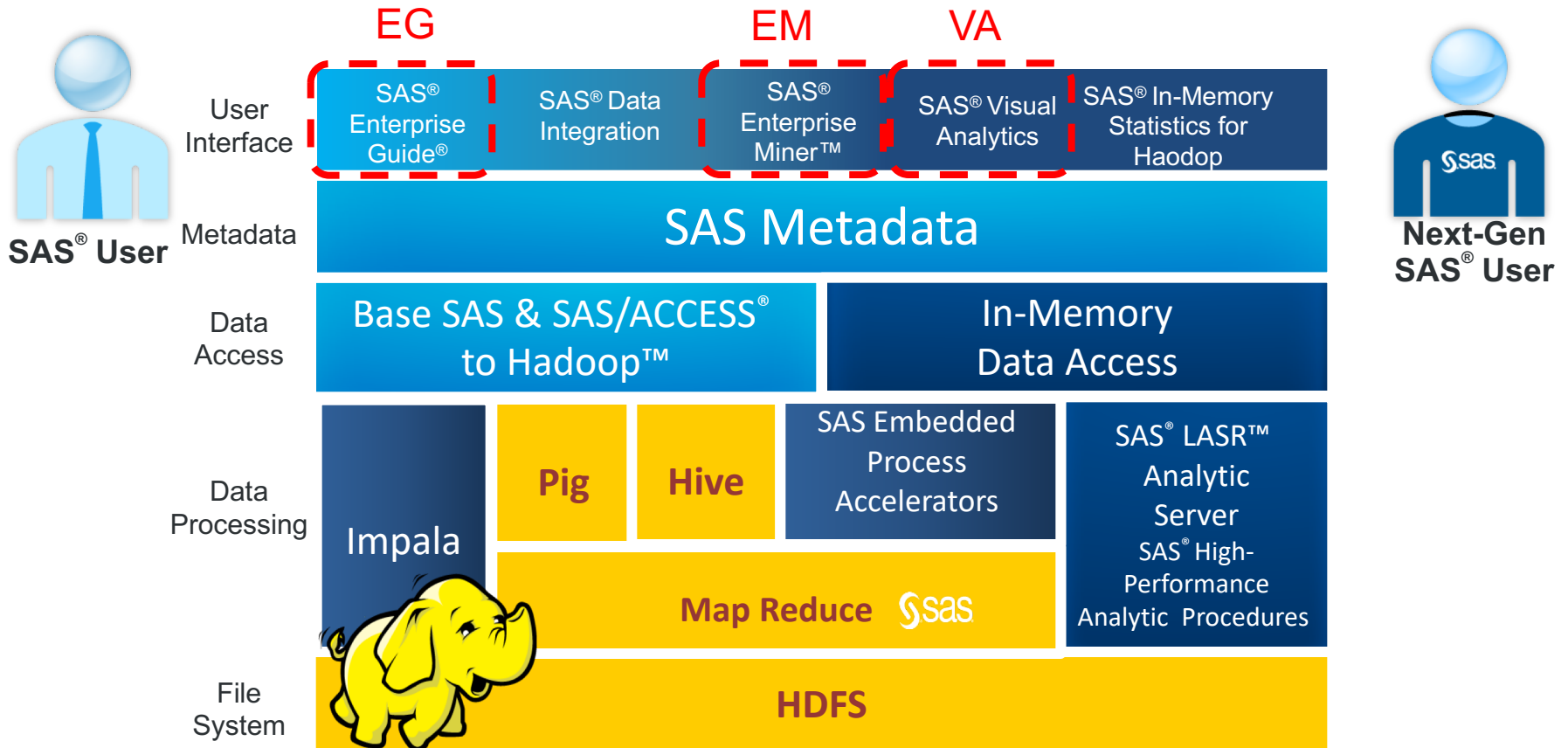


# SAS Big data Strategy

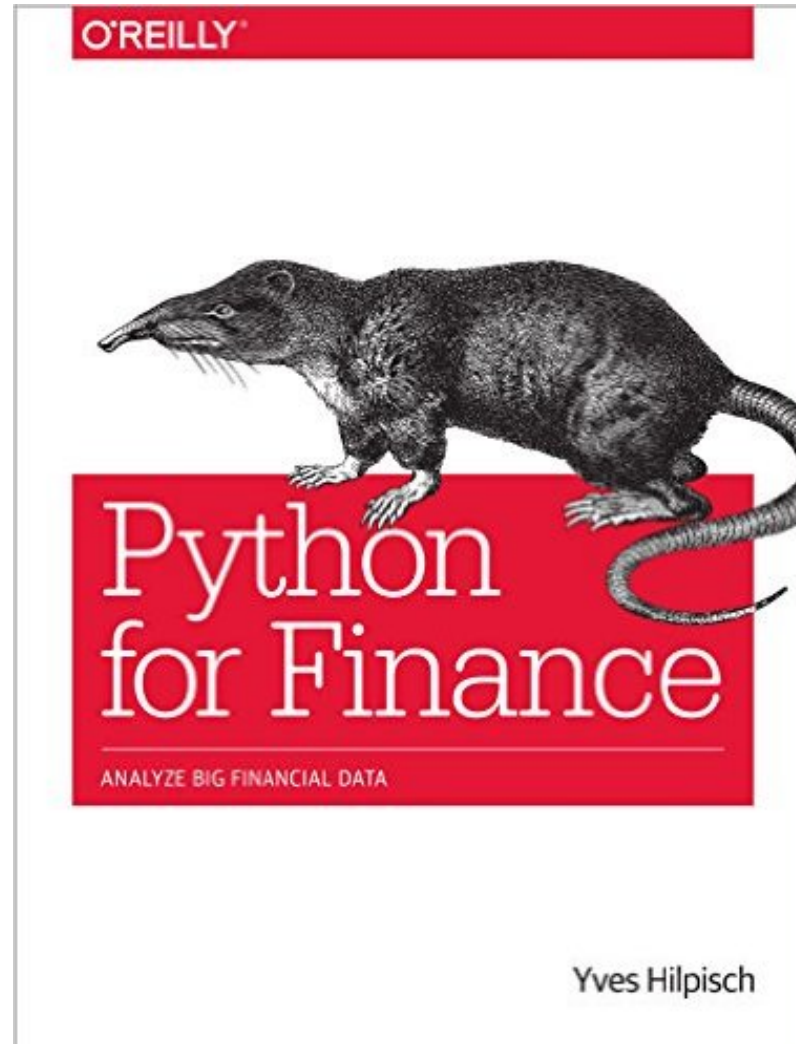
## – SAS areas



# SAS® Within the HADOOP ECOSYSTEM



# Yves Hilpisch, Python for Finance: Analyze Big Financial Data, O'Reilly, 2014



# Business Intelligence Trends

1. **Agile** Information Management (IM)
2. **Cloud** Business Intelligence (BI)
3. **Mobile** Business Intelligence (BI)
4. **Analytics**
5. **Big Data**

# Business Intelligence Trends: Computing and Service

- Cloud Computing and Service
- Mobile Computing and Service
- Social Computing and Service



# Business Intelligence and Analytics

- Business Intelligence 2.0 (BI 2.0)
  - Web Intelligence
  - Web Analytics
  - Web 2.0
  - Social Networking and Microblogging sites
- Data Trends
  - Big Data
- Platform Technology Trends
  - Cloud computing platform

# Business Intelligence and Analytics: Research Directions

## 1. Big Data Analytics

- Data analytics using Hadoop / MapReduce framework

## 2. Text Analytics

- From Information Extraction to Question Answering
- From Sentiment Analysis to Opinion Mining

## 3. Network Analysis

- Link mining
- Community Detection
- Social Recommendation



# Summary

- This course introduces the **fundamental concepts** and **applications technology** of **big data mining**.
- Topics include
  - Big Data Mining
  - Artificial Intelligence and Big Data Analytics
  - Association Analysis
  - Classification and Prediction
  - Cluster Analysis
  - Machine Learning and Deep Learning
  - Data Mining Using SAS Enterprise Miner (SAS EM)
  - Case Study and Implementation of Big Data Mining

# Contact Information

戴敏育 博士 (Min-Yuh Day, Ph.D.)

專任助理教授

淡江大學 資訊管理學系

電話：02-26215656 #2846

傳真：02-26209737

研究室：B929

地址：25137 新北市淡水區英專路151號

Email：myday@mail.tku.edu.tw

網址：<http://mail.tku.edu.tw/myday/>

