

Social Computing and Big Data Analytics

社群運算與大數據分析

Data Science and Big Data Analytics:

Discovering, Analyzing, Visualizing and Presenting Data

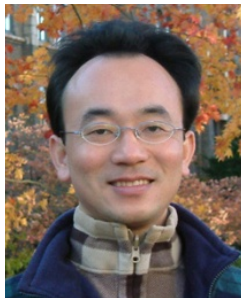
(資料科學與大數據分析：

探索、分析、視覺化與呈現資料)

1052SCBDA02

MIS MBA (M2226) (8606)

Wed, 8,9, (15:10-17:00) (B505)



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/>

2017-02-22



課程大綱 (Syllabus)

- | 週次 (Week) | 日期 (Date) | 內容 (Subject/Topics) |
|-----------|------------|--|
| 1 | 2017/02/15 | Course Orientation for Social Computing and Big Data Analytics
(社群運算與大數據分析課程介紹) |
| 2 | 2017/02/22 | Data Science and Big Data Analytics:
Discovering, Analyzing, Visualizing and Presenting Data
(資料科學與大數據分析：
探索、分析、視覺化與呈現資料) |
| 3 | 2017/03/01 | Fundamental Big Data: MapReduce Paradigm,
Hadoop and Spark Ecosystem
(大數據基礎：MapReduce典範、
Hadoop與Spark生態系統) |

課程大綱 (Syllabus)

週次 (Week)	日期 (Date)	內容 (Subject/Topics)
4	2017/03/08	Big Data Processing Platforms with SMACK: Spark, Mesos, Akka, Cassandra and Kafka (大數據處理平台SMACK： Spark, Mesos, Akka, Cassandra, Kafka)
5	2017/03/15	Big Data Analytics with Numpy in Python (Python Numpy 大數據分析)
6	2017/03/22	Finance Big Data Analytics with Pandas in Python (Python Pandas 財務大數據分析)
7	2017/03/29	Text Mining Techniques and Natural Language Processing (文字探勘分析技術與自然語言處理)
8	2017/04/05	Off-campus study (教學行政觀摩日)

課程大綱 (Syllabus)

週次 (Week)	日期 (Date)	內容 (Subject/Topics)
9	2017/04/12	Social Media Marketing Analytics (社群媒體行銷分析)
10	2017/04/19	期中報告 (Midterm Project Report)
11	2017/04/26	Deep Learning with Theano and Keras in Python (Python Theano 和 Keras 深度學習)
12	2017/05/03	Deep Learning with Google TensorFlow (Google TensorFlow 深度學習)
13	2017/05/10	Sentiment Analysis on Social Media with Deep Learning (深度學習社群媒體情感分析)

課程大綱 (Syllabus)

週次 (Week)	日期 (Date)	內容 (Subject/Topics)
14	2017/05/17	Social Network Analysis (社會網絡分析)
15	2017/05/24	Measurements of Social Network (社會網絡量測)
16	2017/05/31	Tools of Social Network Analysis (社會網絡分析工具)
17	2017/06/07	Final Project Presentation I (期末報告 I)
18	2017/06/14	Final Project Presentation II (期末報告 II)

2017/02/22

Data Science and

Big Data Analytics:

Discovering, Analyzing,

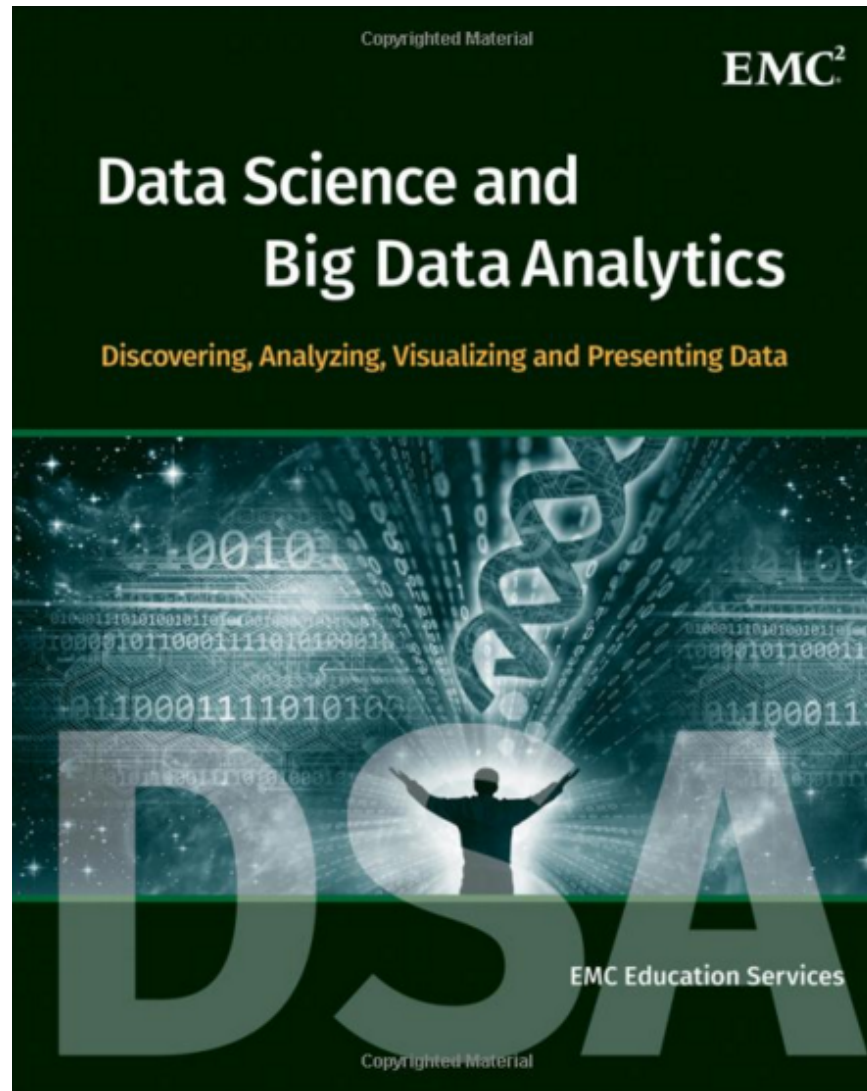
Visualizing and Presenting Data

(資料科學與大數據分析：

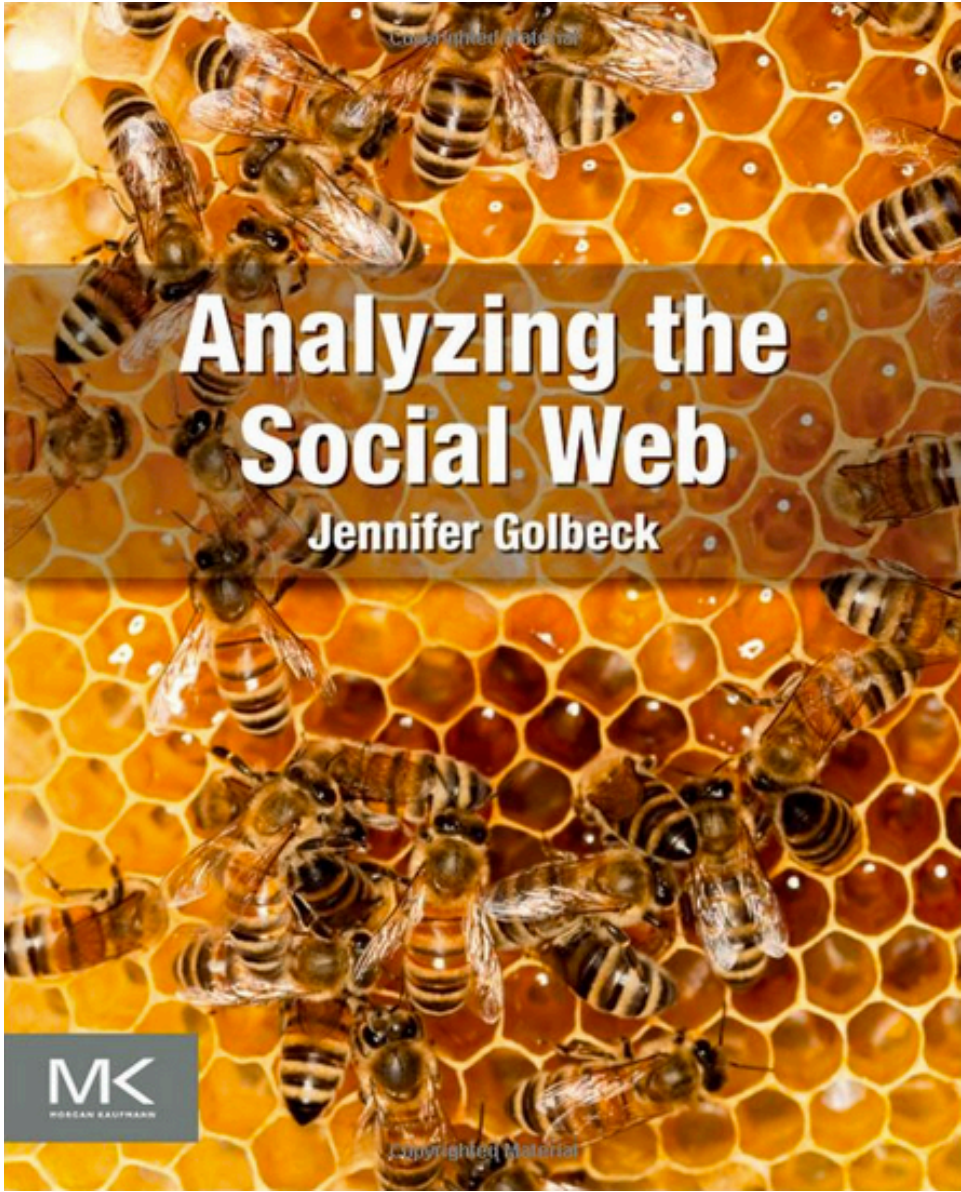
探索、分析、

視覺化與呈現資料)

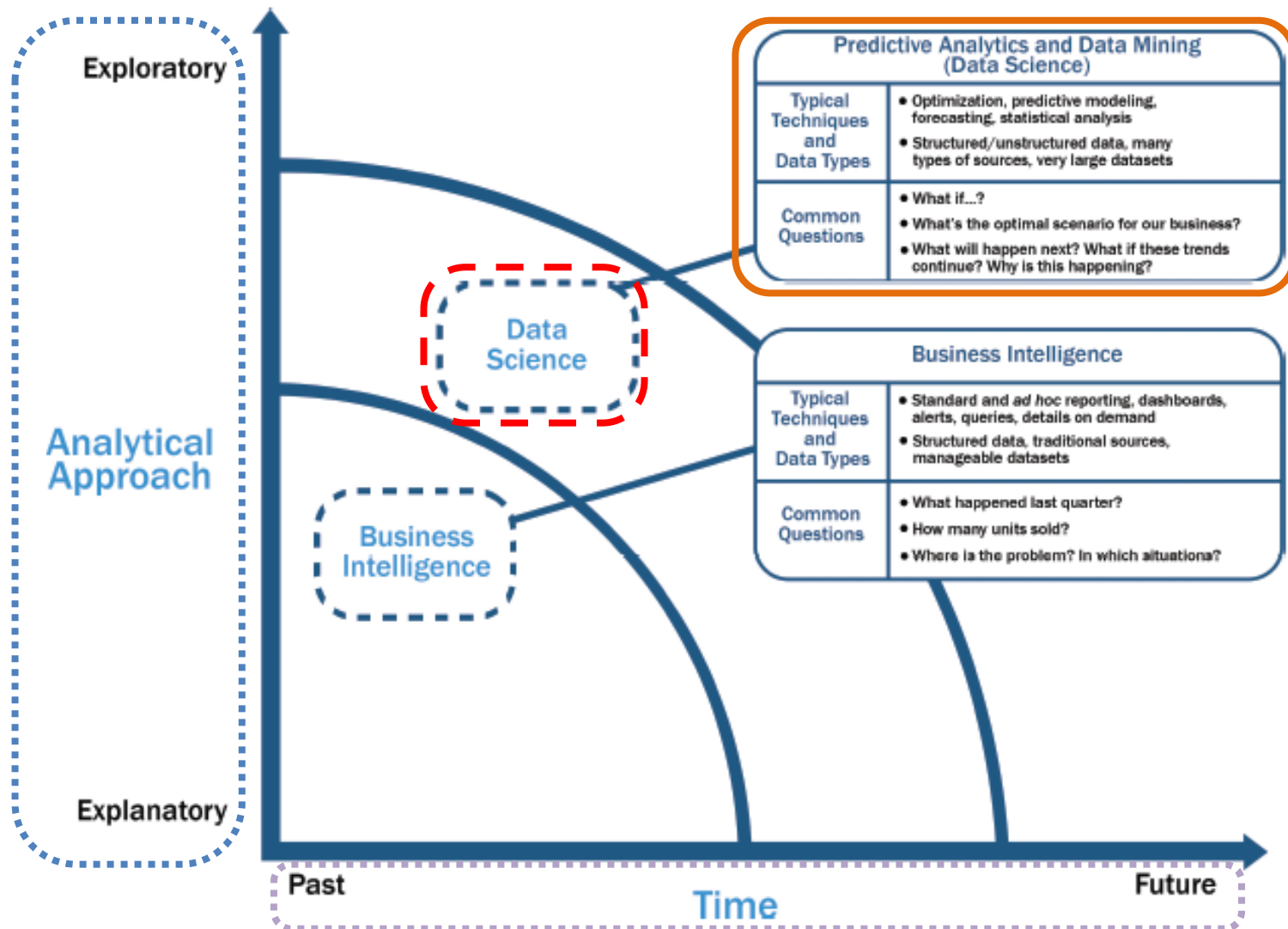
**EMC Education Services,
Data Science and Big Data Analytics:
Discovering, Analyzing, Visualizing and Presenting Data,
Wiley, 2015**



Jennifer Golbeck (2013), *Analyzing the Social Web*, Morgan Kaufmann



Data Science and Business Intelligence



Data Science and Business Intelligence



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?

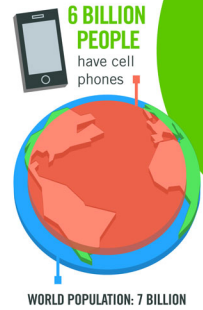
What will happen next?

What if these trends continue?

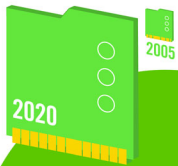
Why is this happening?

Big Data 4 V

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



Volume SCALE OF DATA



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



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

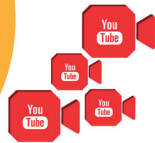


30 BILLION PIECES OF CONTENT
are shared on Facebook every month



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



Velocity ANALYSIS OF STREAMING DATA



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



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

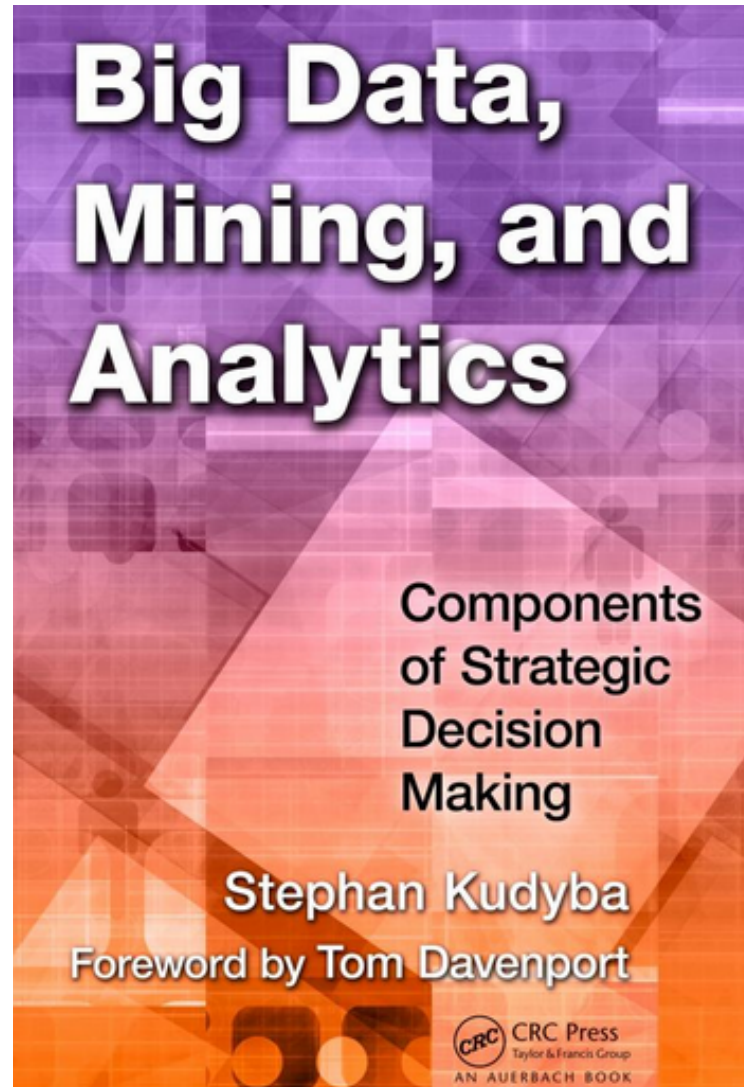
Veracity UNCERTAINTY OF DATA

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

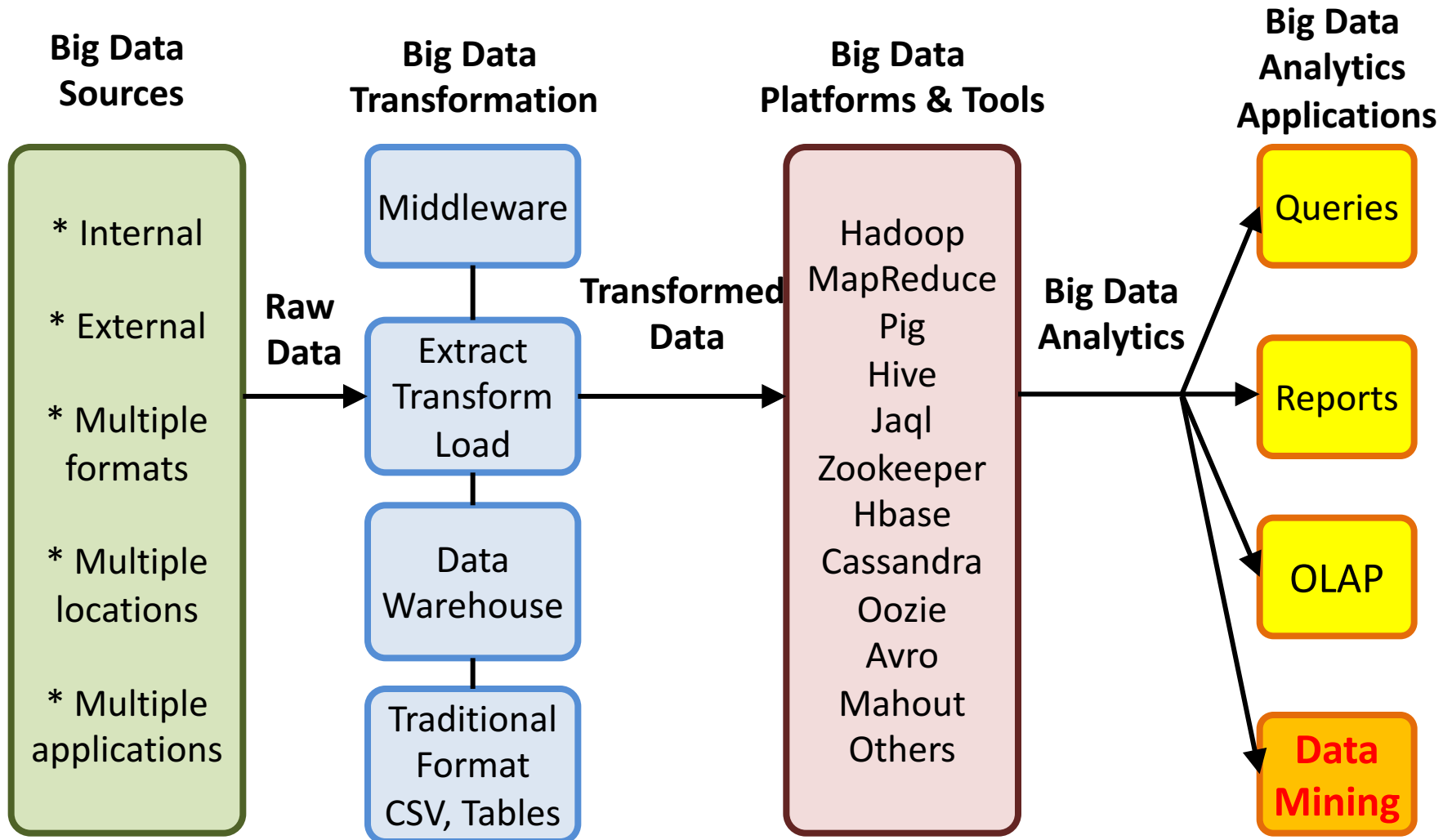
value

Big Data
Analytics
and
Data Mining

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



Architecture of Big Data Analytics

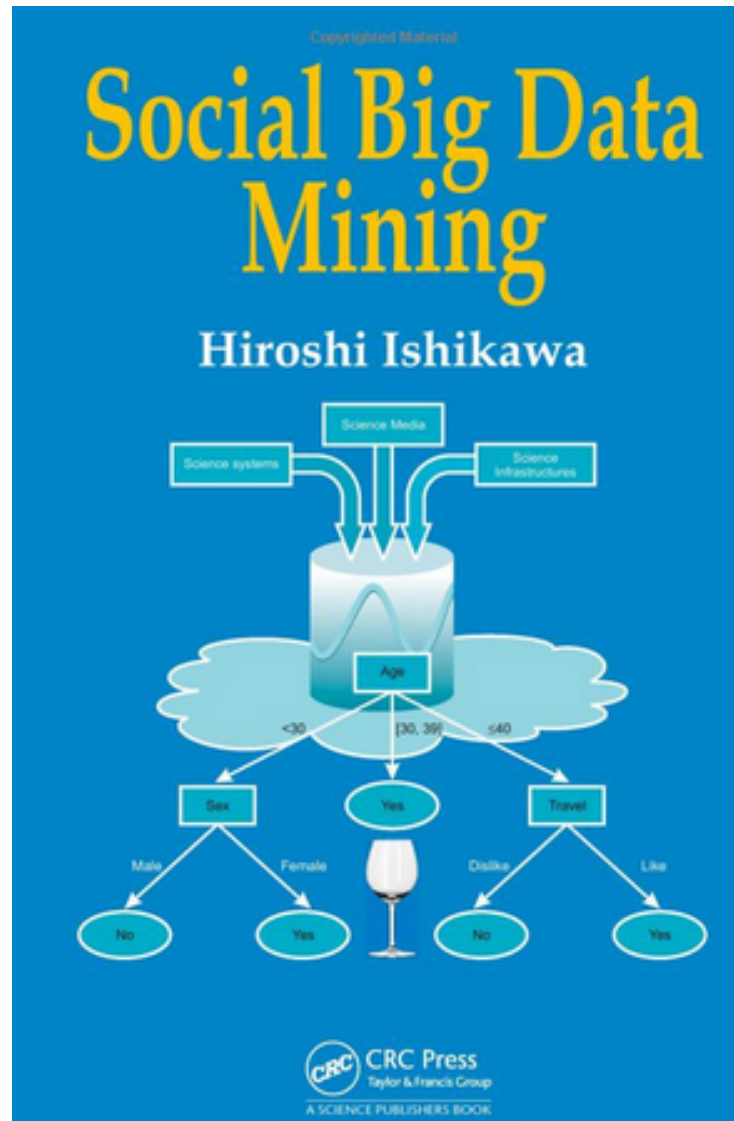


Architecture of Big Data Analytics



Social Big Data Mining

(Hiroshi Ishikawa, 2015)

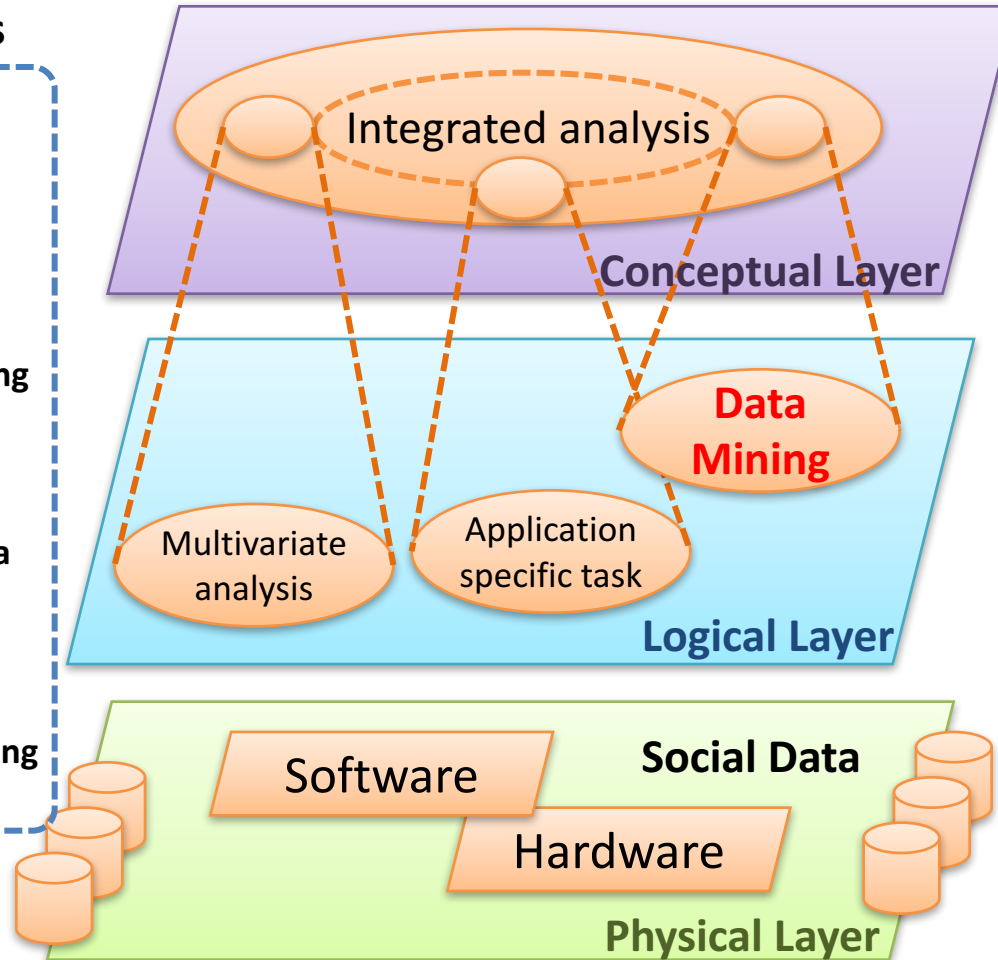


Architecture for Social Big Data Mining

(Hiroshi Ishikawa, 2015)

Enabling Technologies

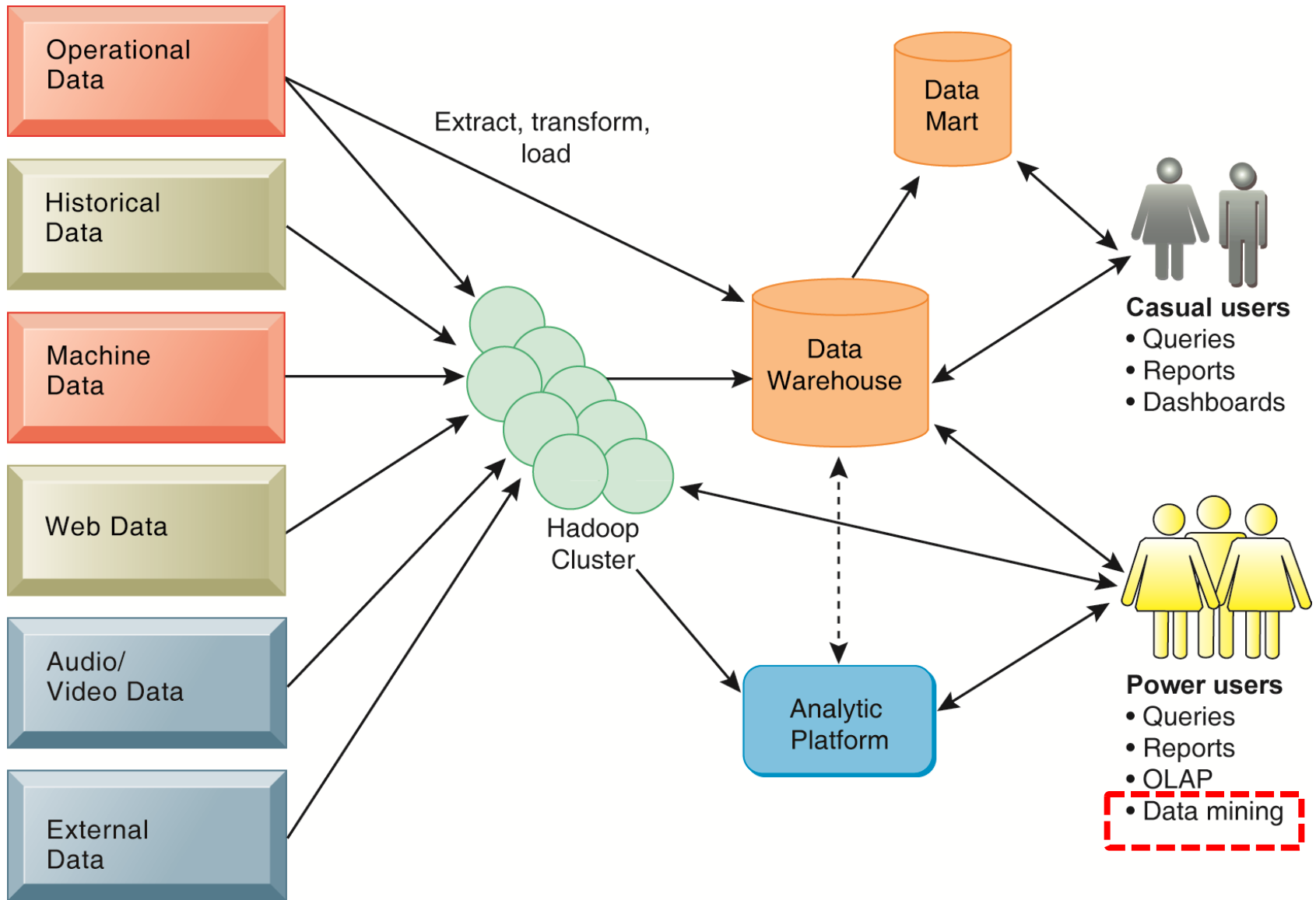
- Integrated analysis model
- Natural Language Processing
- Information Extraction
- Anomaly Detection
- Discovery of relationships among heterogeneous data
- Large-scale visualization
- Parallel distributed processing



Analysts

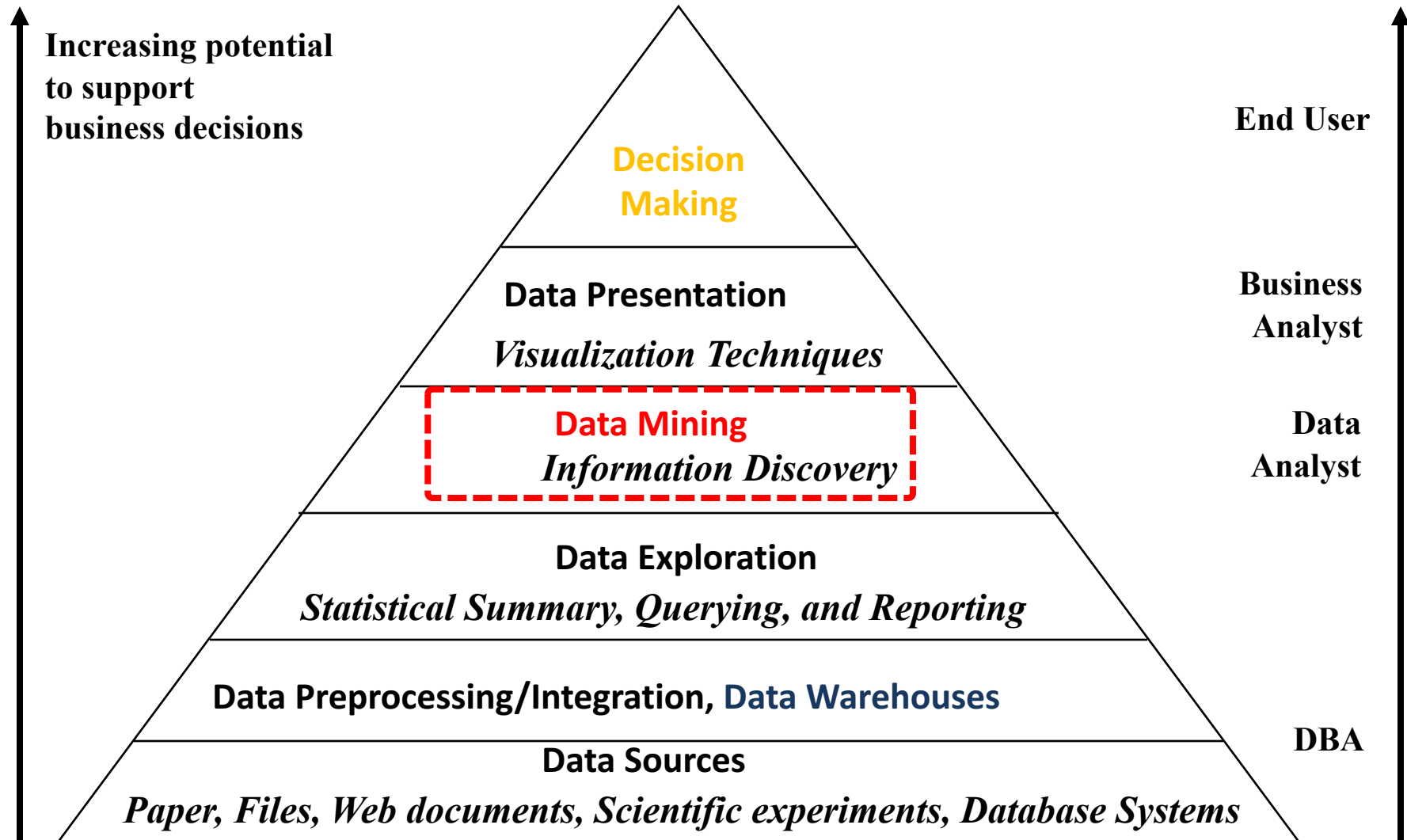
- Model Construction
- Explanation by Model
- Construction and confirmation of individual hypothesis
- Description and execution of application-specific task

Business Intelligence (BI) Infrastructure

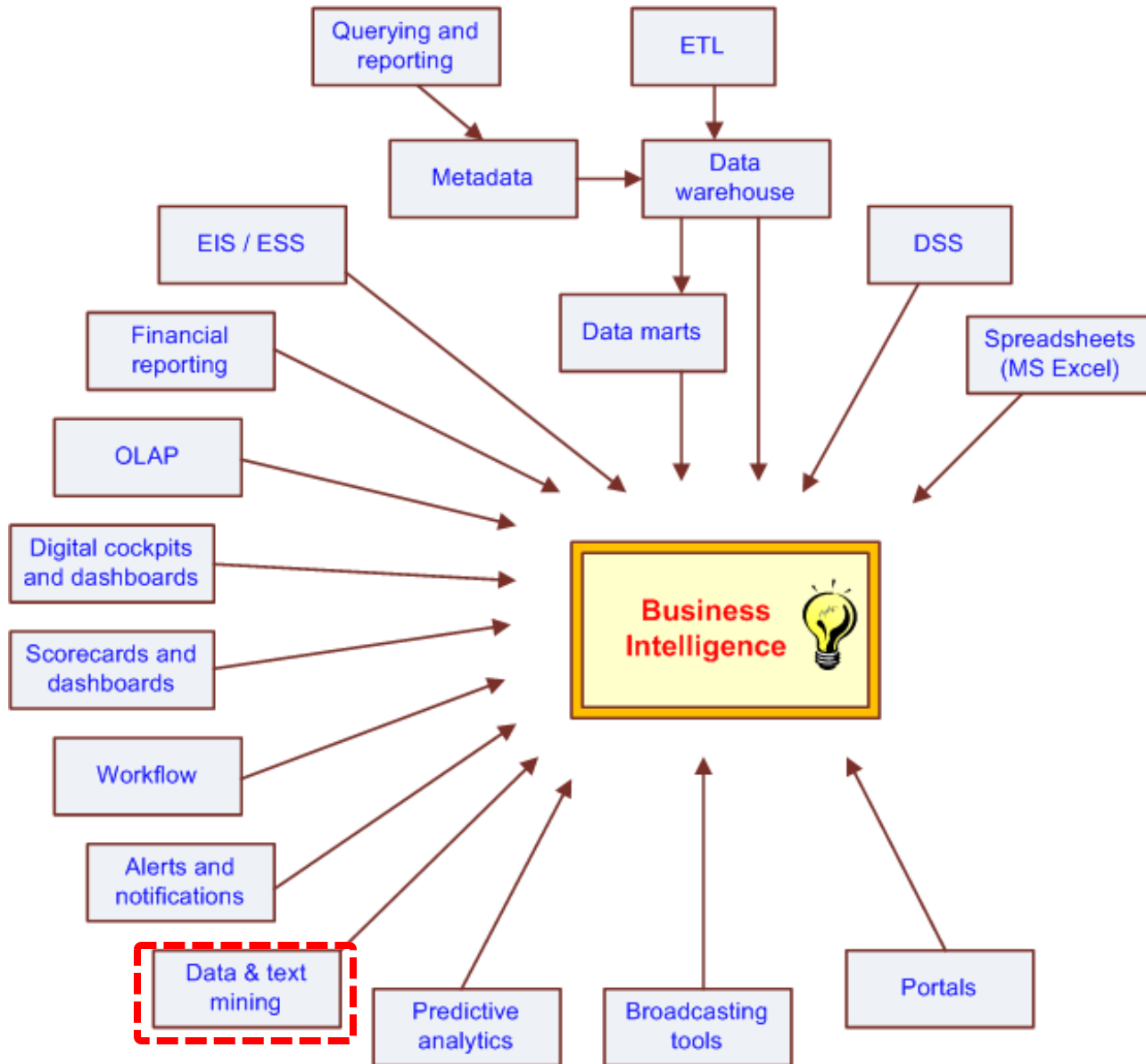


Data Warehouse

Data Mining and Business Intelligence



The Evolution of BI Capabilities



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)

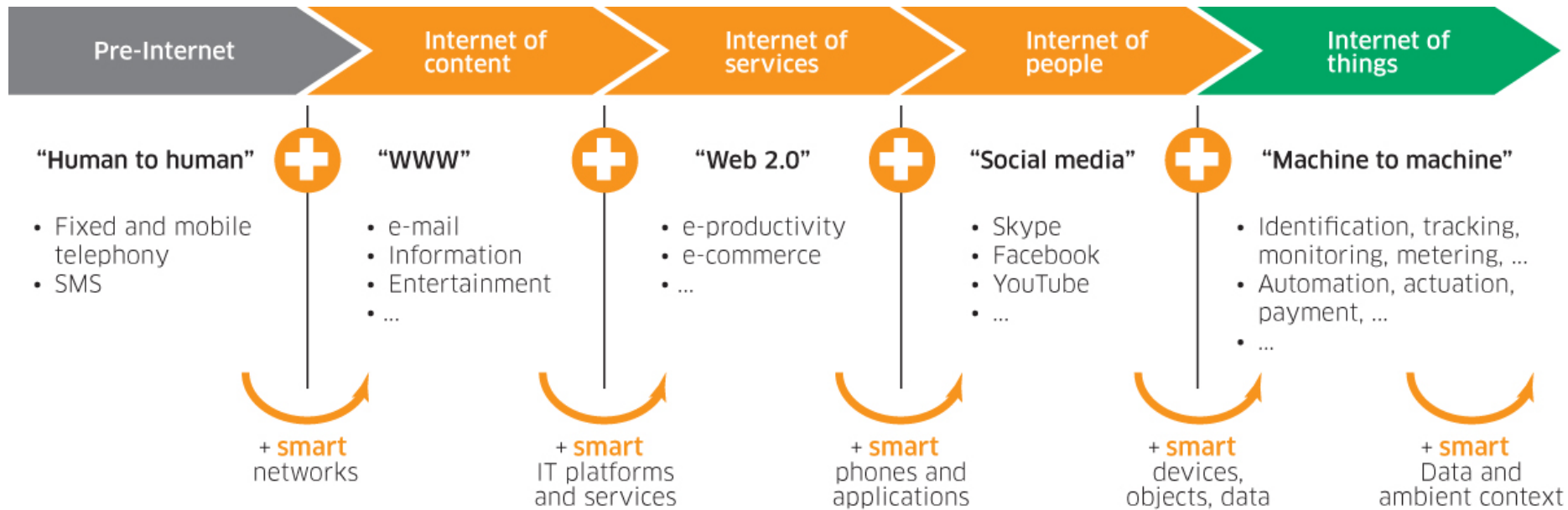
Big Data Analysis

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

Internet Evolution

Internet of People (IoP): Social Media

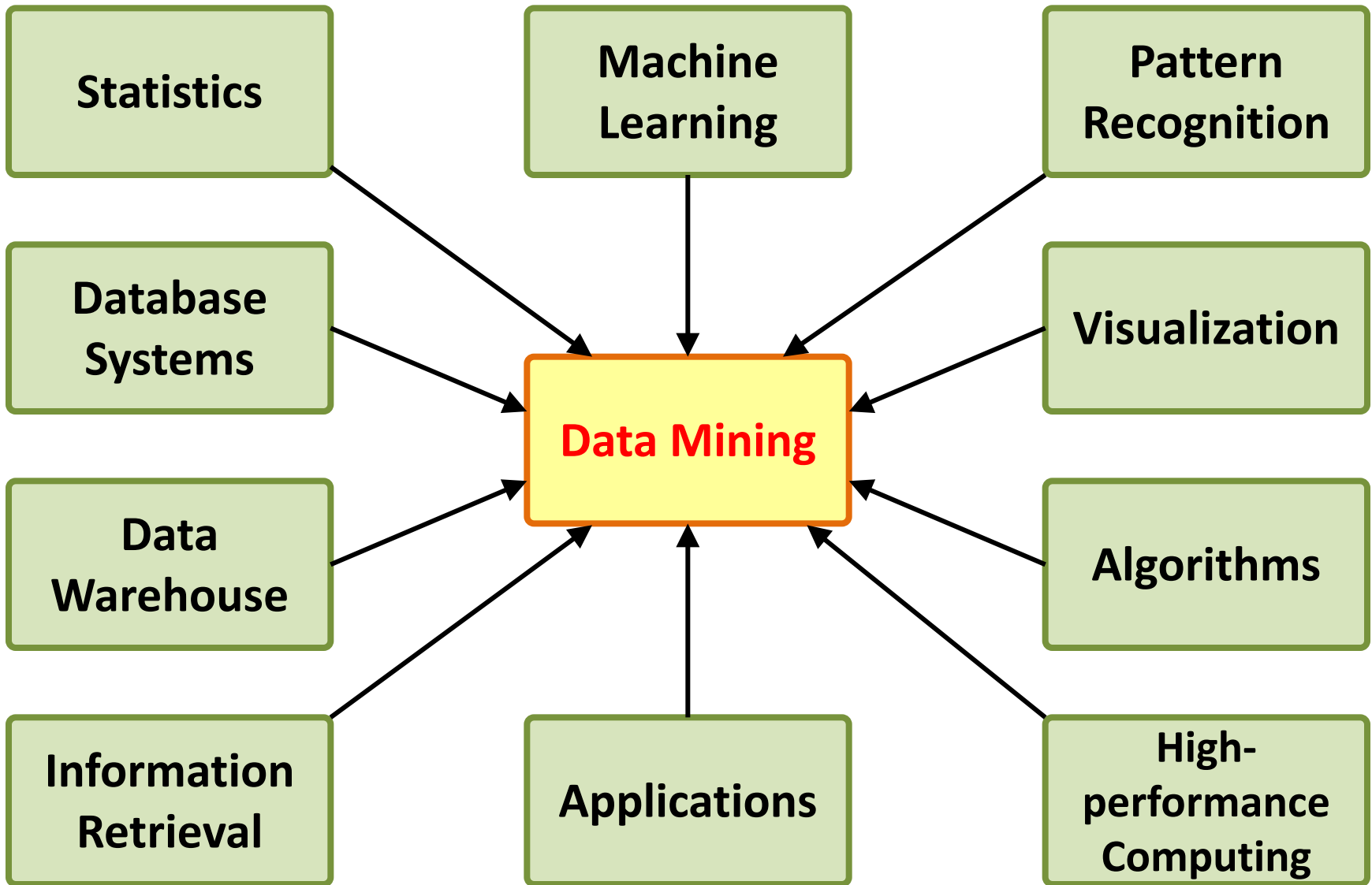
Internet of Things (IoT): Machine to Machine

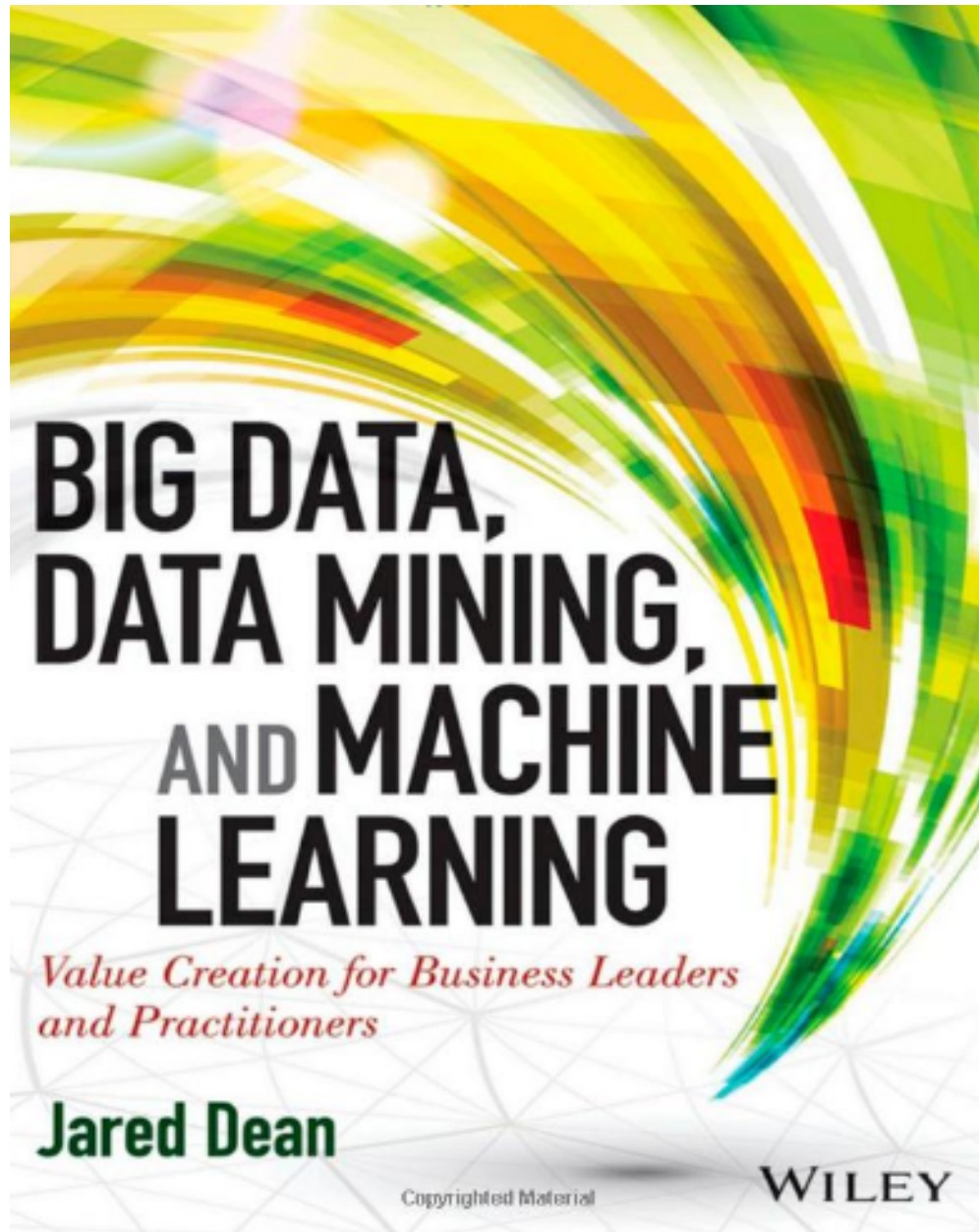


Source: Marc Jadoul (2015), The IoT: The next step in internet evolution, March 11, 2015

<http://www2.alcatel-lucent.com/techzine/iot-internet-of-things-next-step-evolution/>

Data Mining Technologies

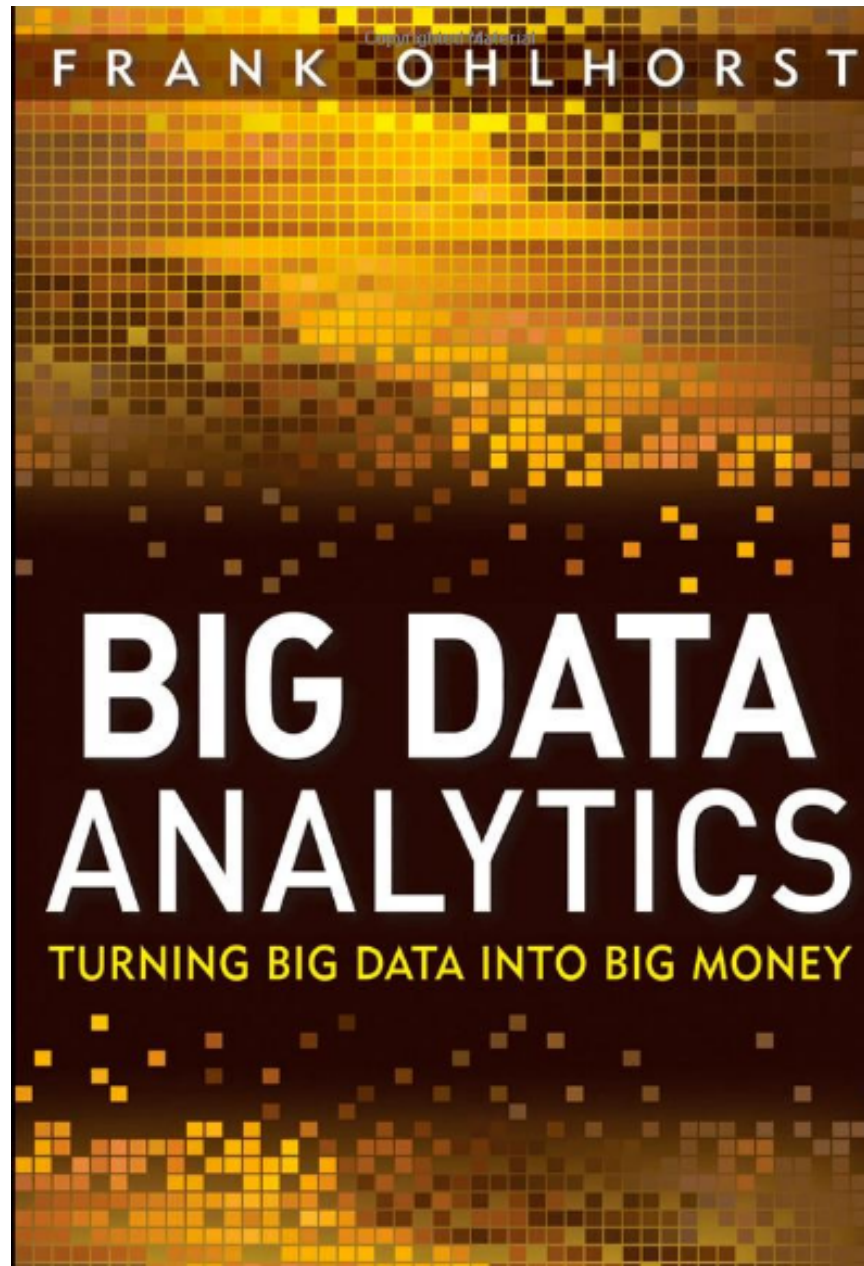




Deep Learning

Intelligence from Big Data







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

Big Data, Big Analytics:

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

Big Data, Prediction vs. Explanation

Big Data:

The Management Revolution

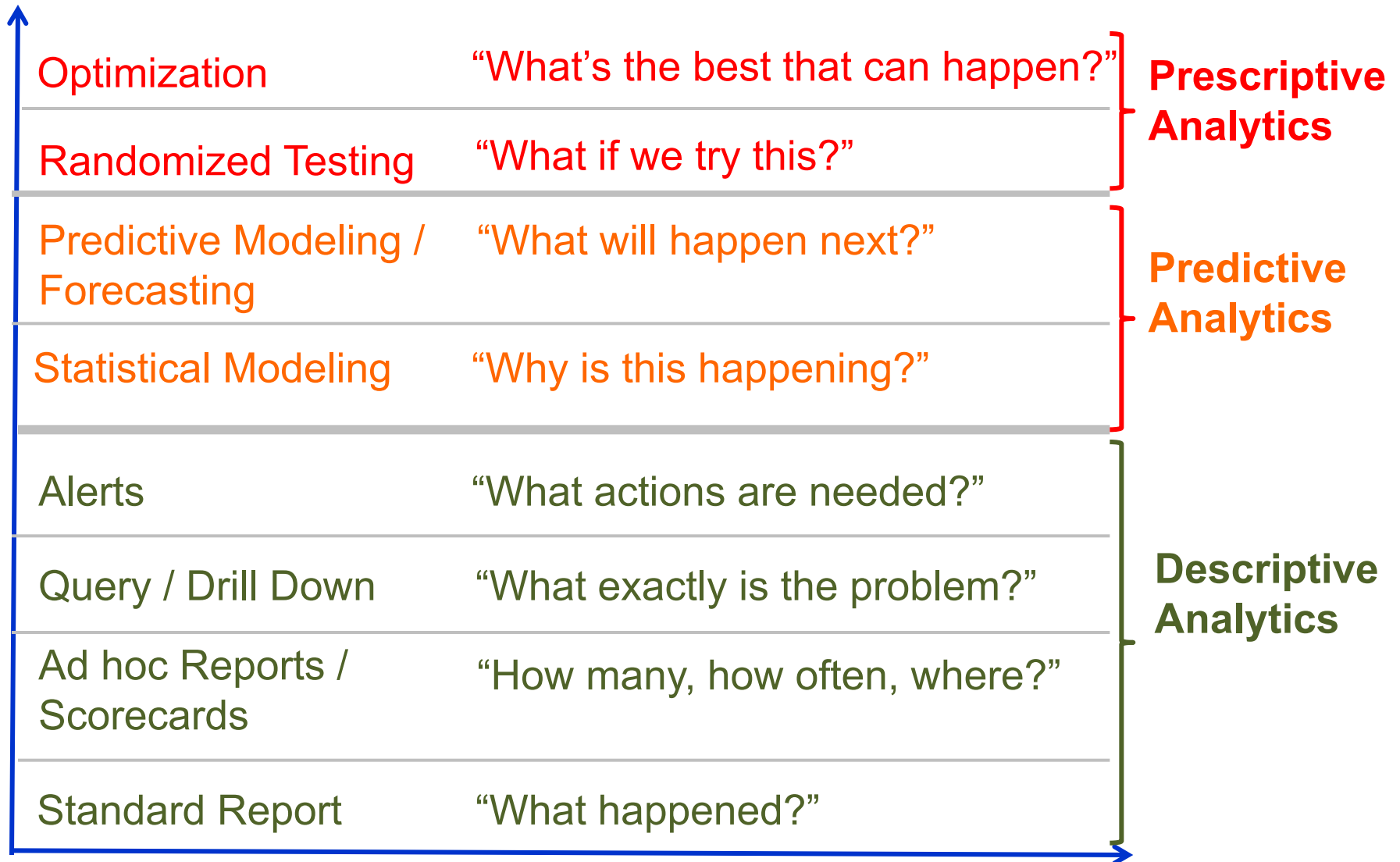
Business Intelligence and Enterprise Analytics

- Predictive analytics
- Data mining
- Business analytics
- Web analytics
- **Big-data** analytics

Three Types of Business Analytics

- Prescriptive Analytics
- Predictive Analytics
- Descriptive Analytics

Three Types of Business Analytics



Data Scientist:

The Sexiest Job of the 21st Century

**Meet the people who
can coax treasure out of
messy, unstructured data.**

*by Thomas H. Davenport
and D.J. Patil*

When Jonathan Goldman arrived for work in June 2006 at LinkedIn, the business networking site, the place still felt like a start-up. The company had just under 8 million accounts, and the number was growing quickly as existing members invited their friends and colleagues to join. But users weren't seeking out connections with the people who were already on the site at the rate executives had expected. Something was apparently missing in the social experience. As one LinkedIn manager put it, "It was like arriving at a conference reception and realizing you don't know anyone. So you just stand in the corner sipping your drink—and you probably leave early."

Big Data



**Mobile
Sensors**



**Social
Media**



**Video
Surveillance**



**Video
Rendering**



**Smart
Grids**



**Geophysical
Exploration**

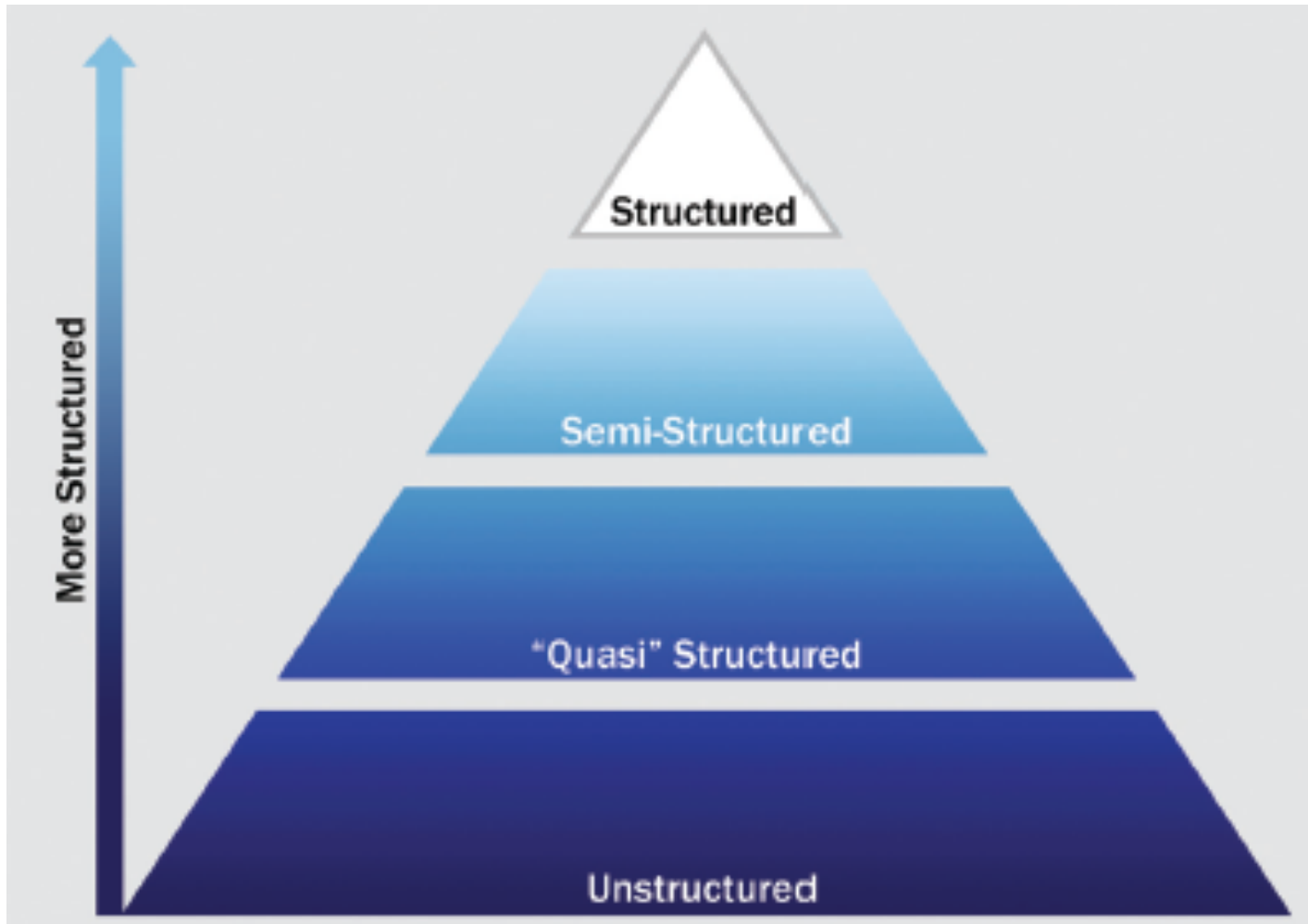


**Medical
Imaging**

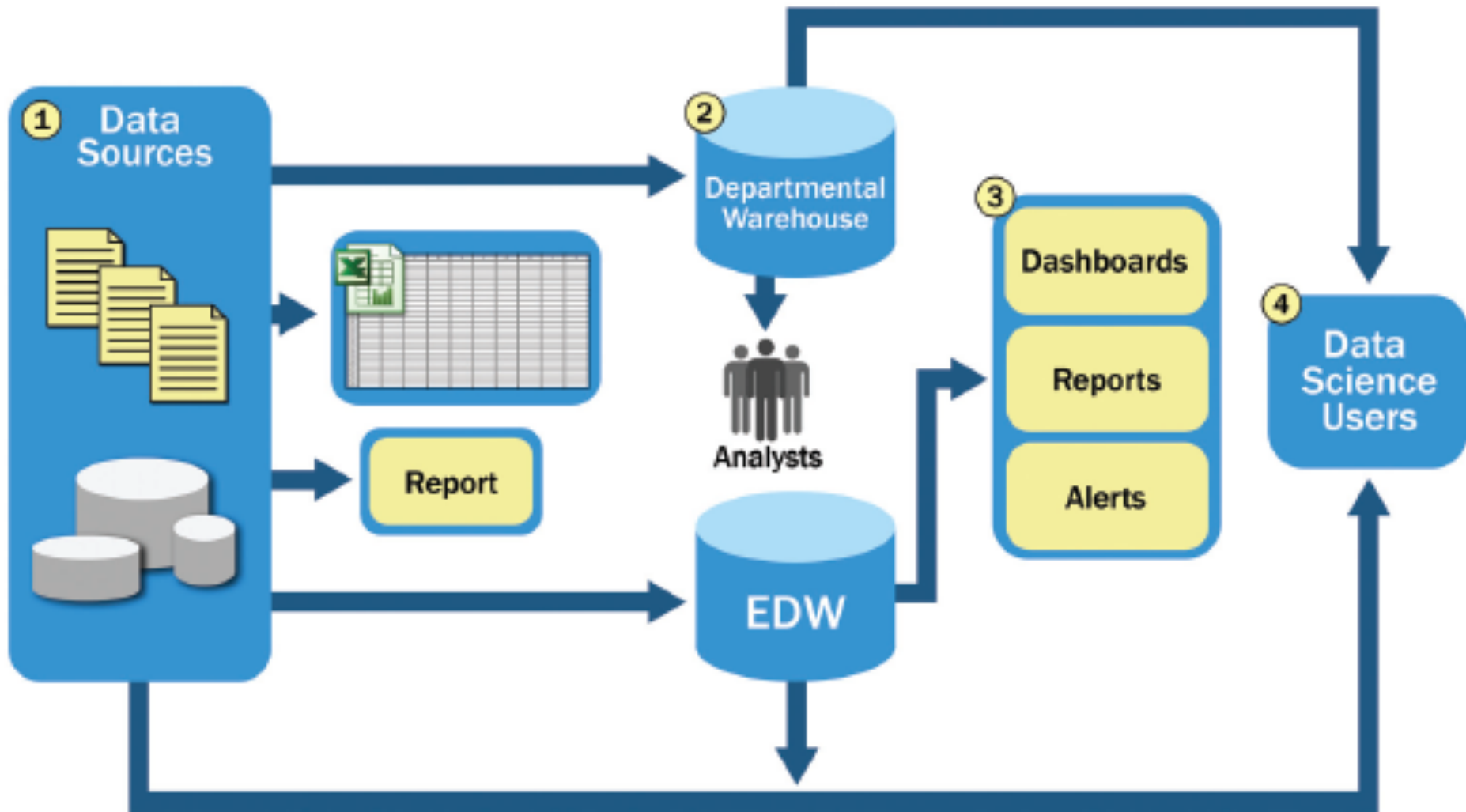


**Gene
Sequencing**

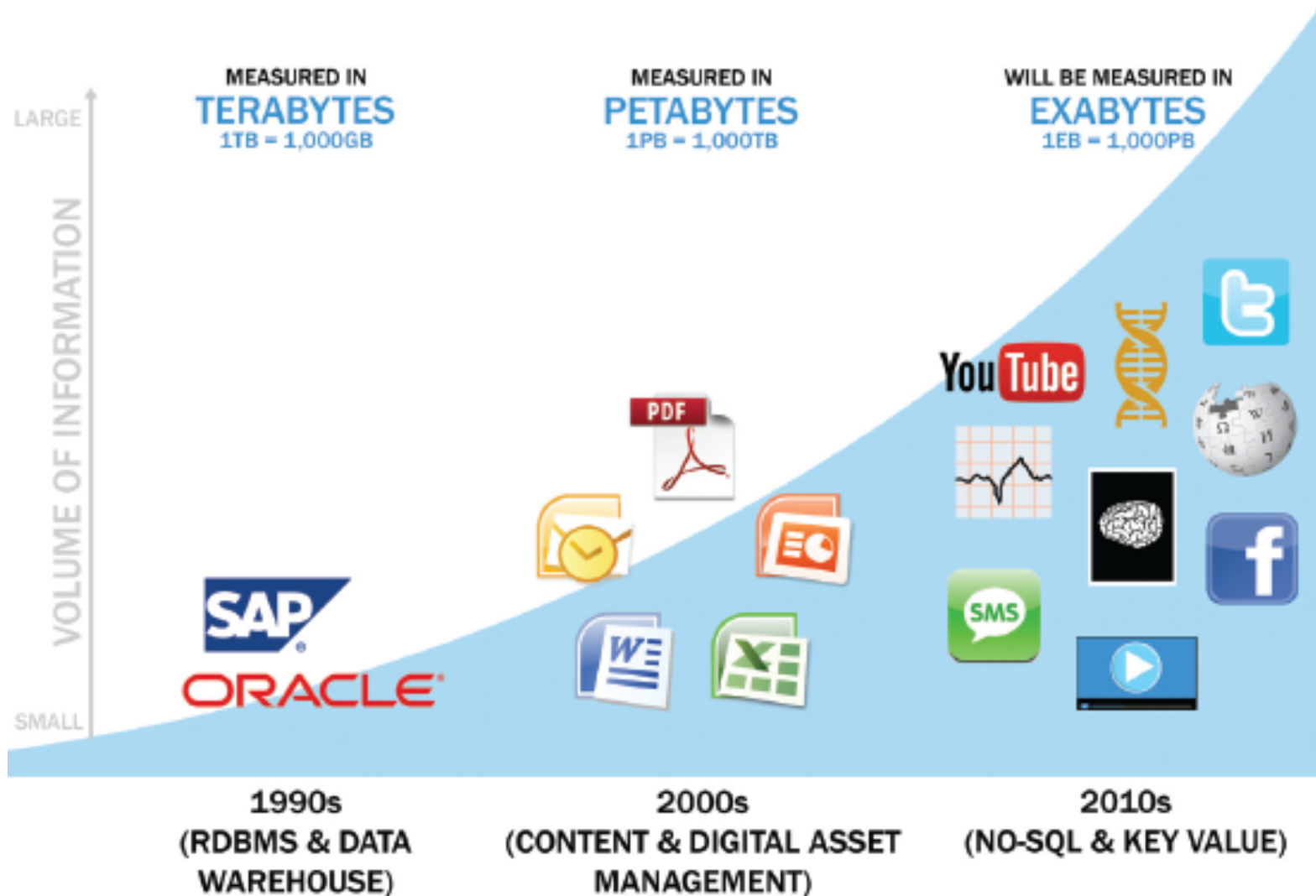
Big Data Growth is increasingly **unstructured**



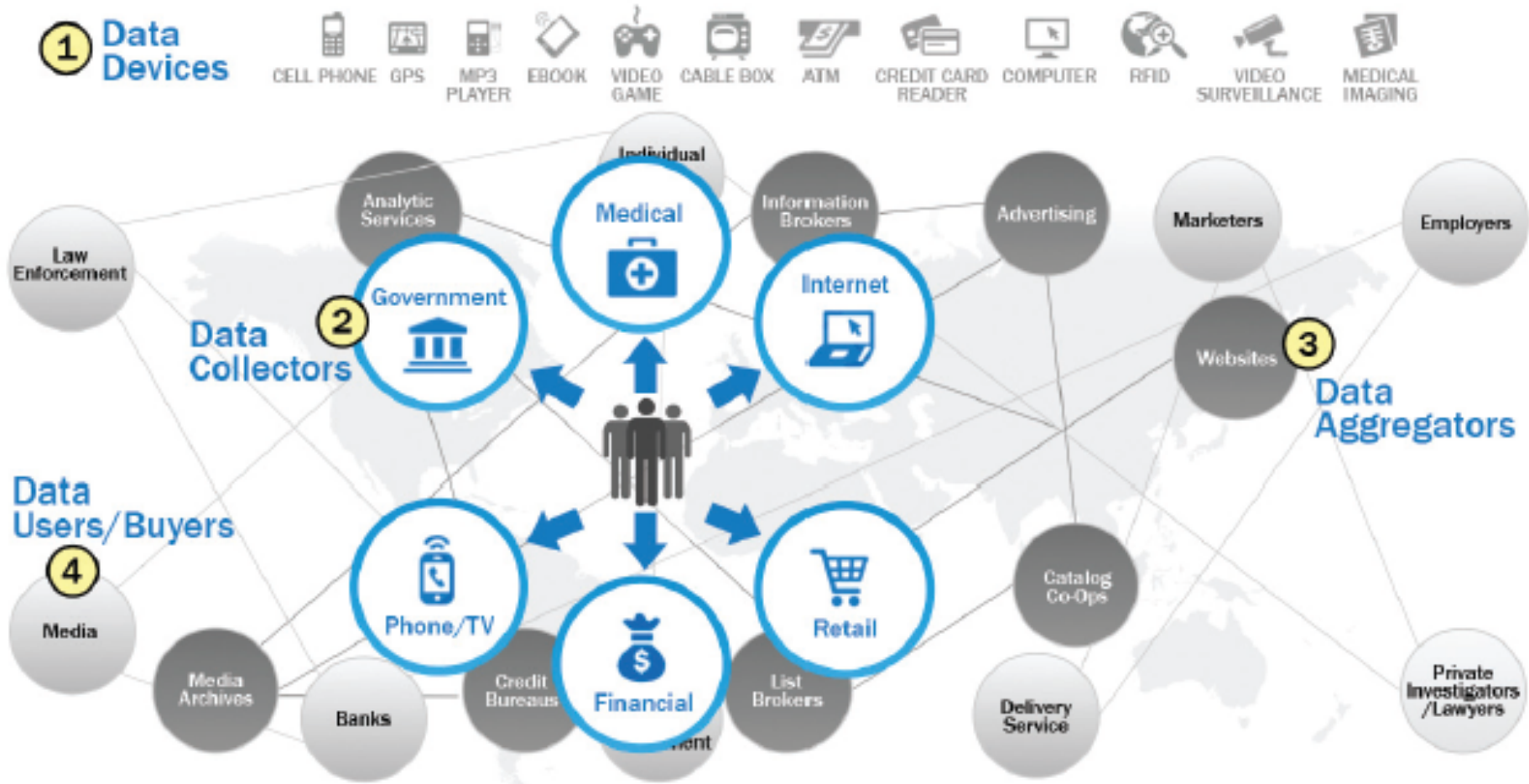
Typical Analytic Architecture



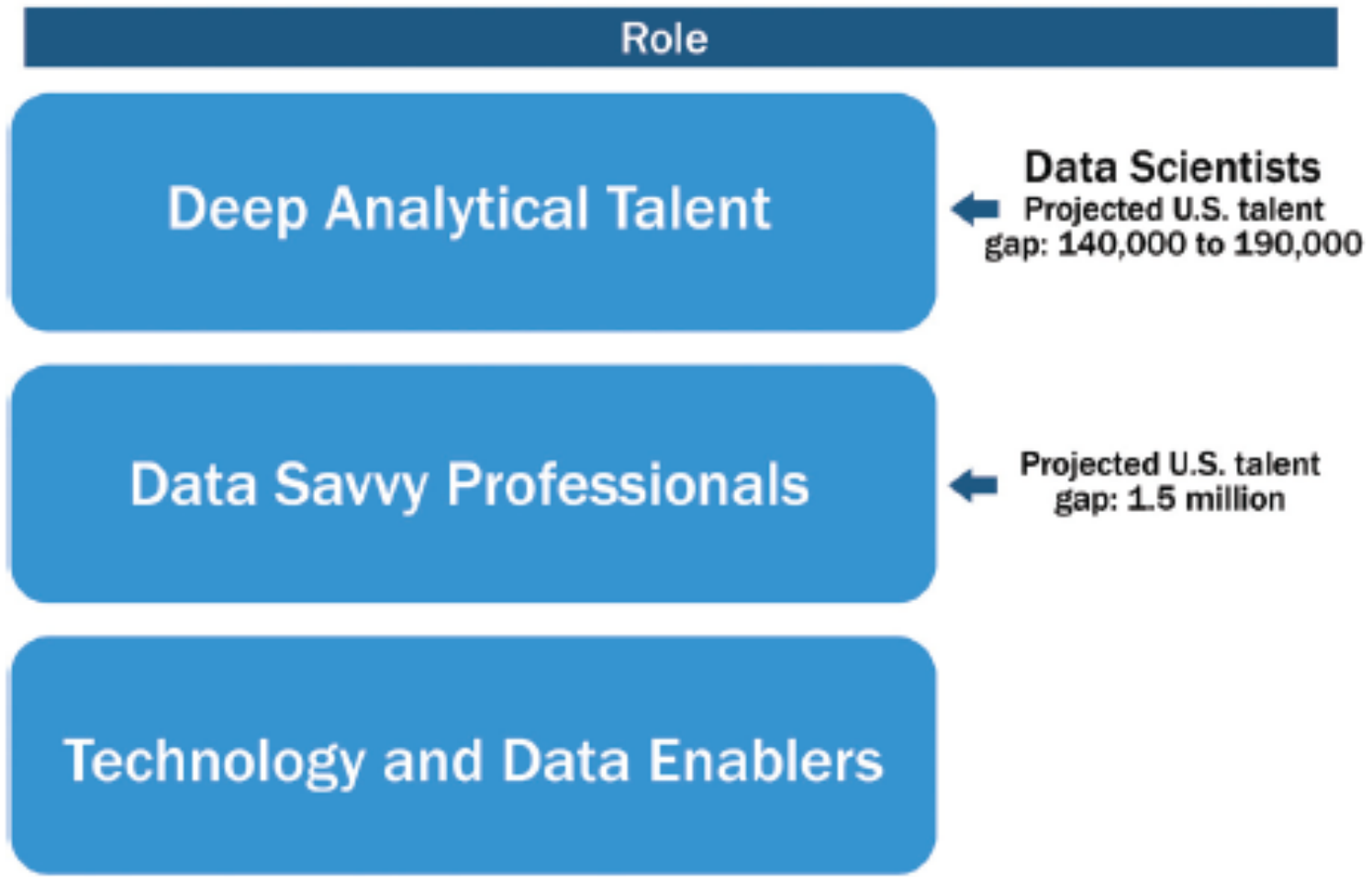
Data Evolution and the Rise of Big Data Sources



Emerging Big Data Ecosystem



Key Roles for the New Big Data Ecosystem

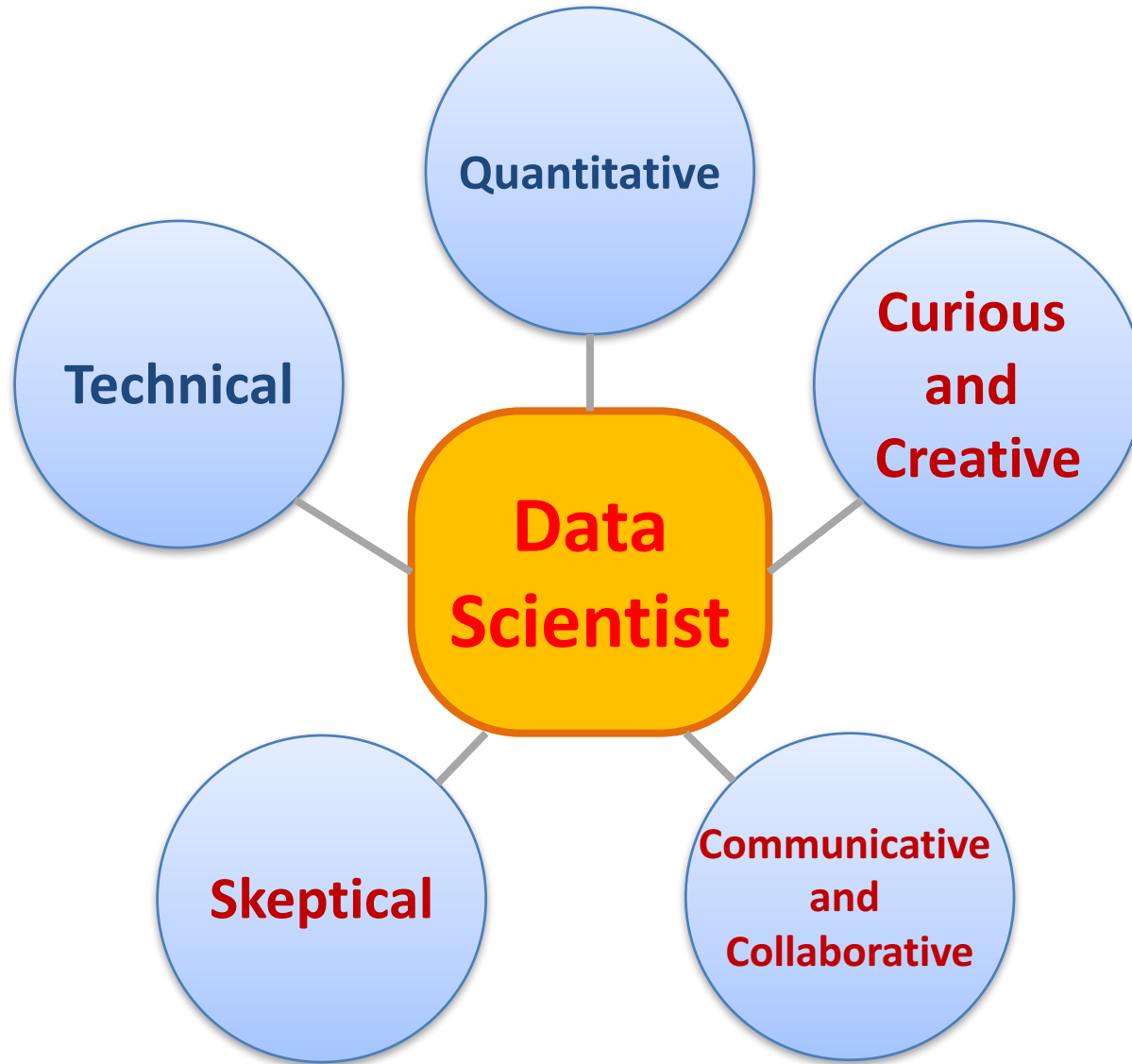


Note: Figures above reflect a projected talent gap in US in 2018, as shown in McKinsey May 2011 article "Big Data: The Next Frontier for Innovation, Competition, and Productivity"

Profile of a Data Scientist

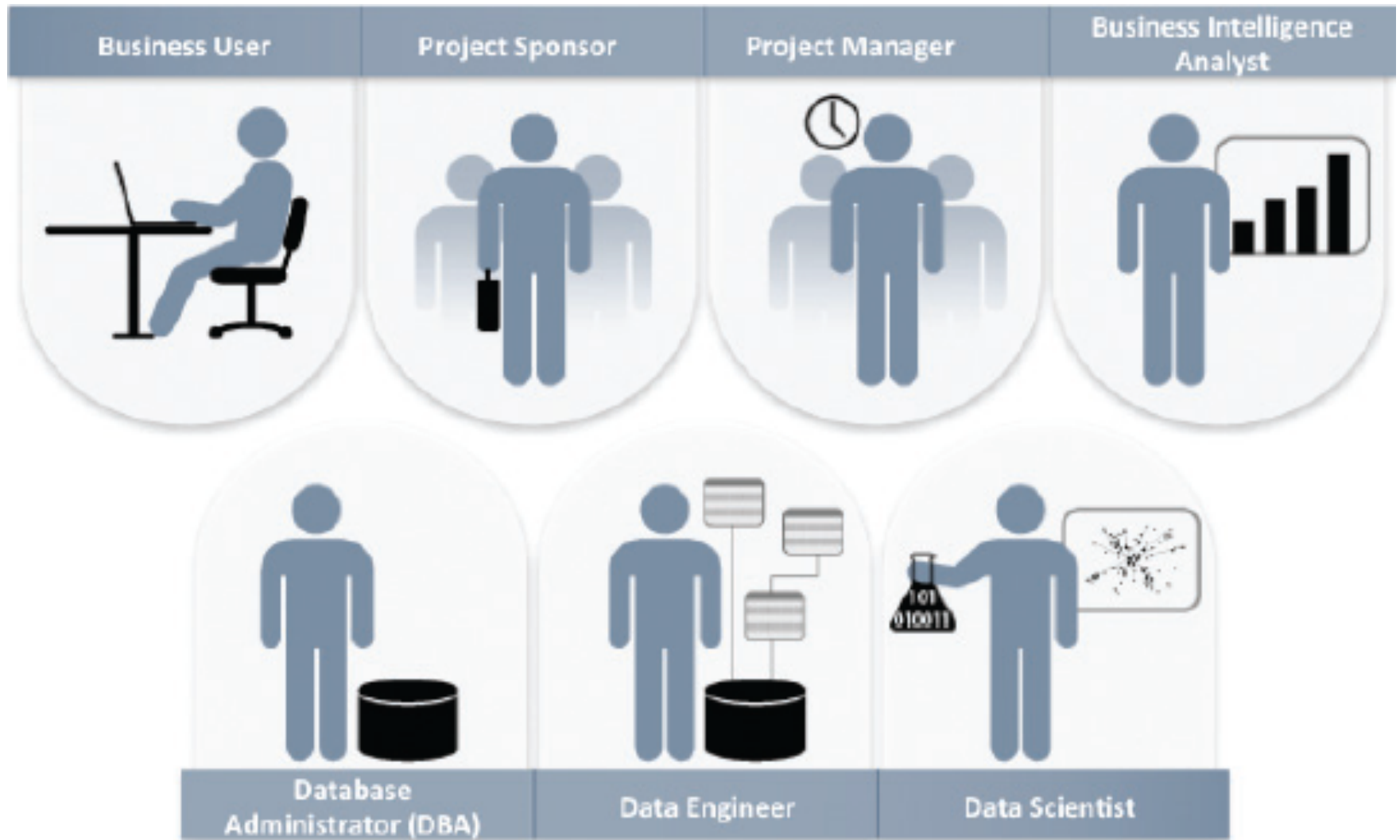
- **Quantitative**
 - mathematics or statistics
- **Technical**
 - software engineering,
machine learning,
and programming skills
- **Skeptical mind-set** and **critical thinking**
- **Curious** and **creative**
- **Communicative** and **collaborative**

Data Scientist Profile

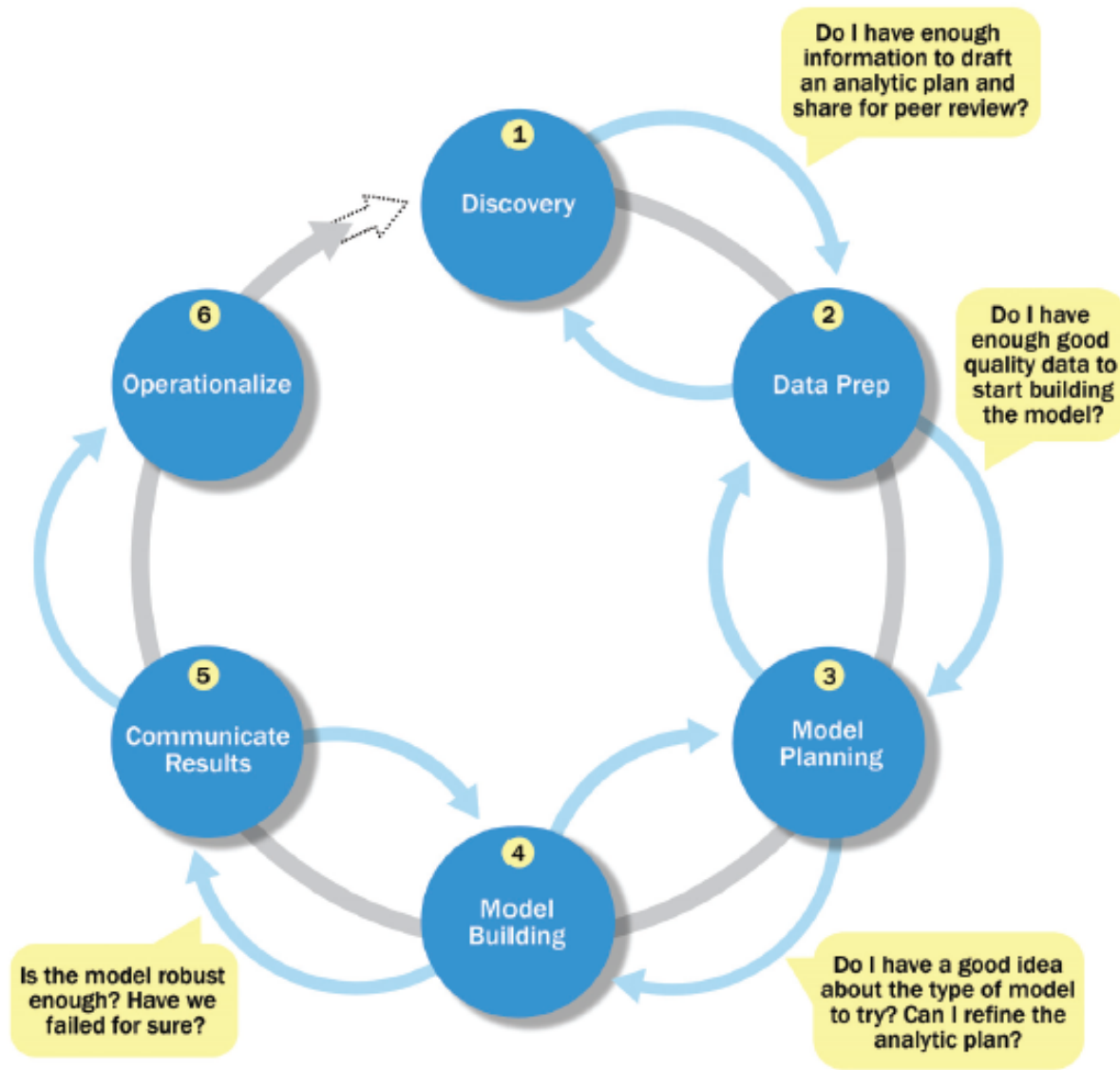


Big Data Analytics Lifecycle

Key Roles for a Successful Analytics Project



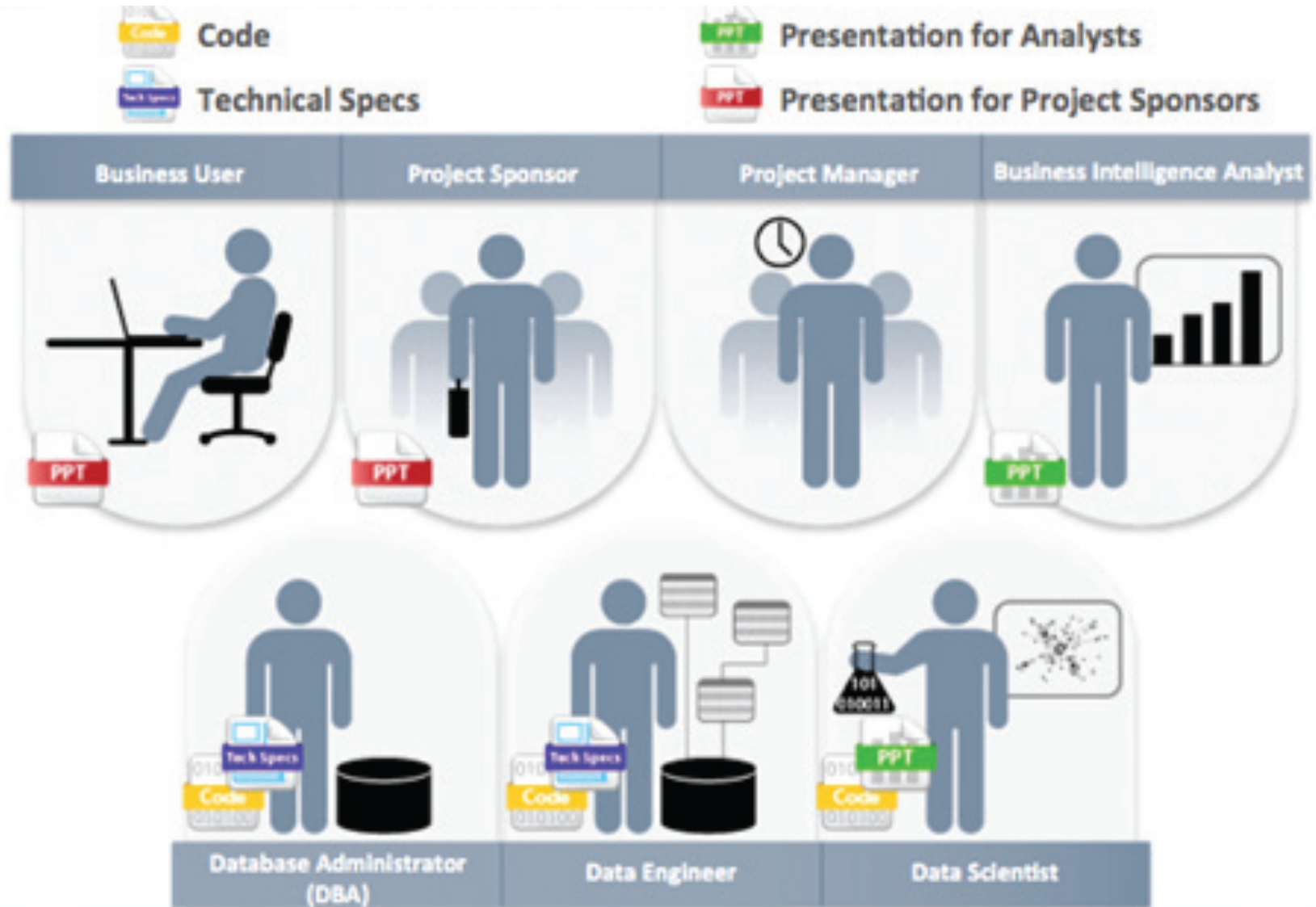
Overview of Data Analytics Lifecycle



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

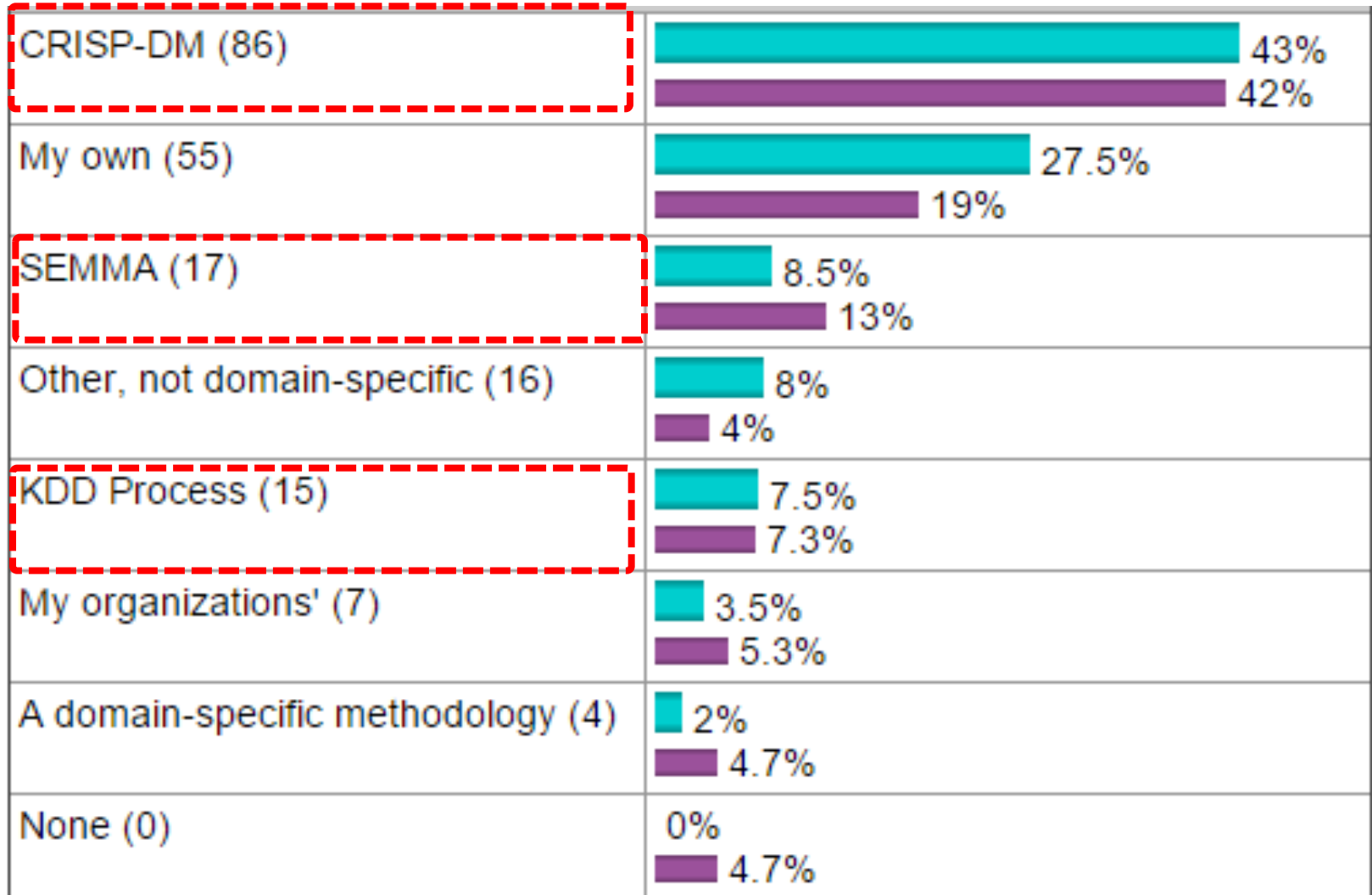
Data Mining Process

- 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



Data Mining:

Core **Analytics** Process

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

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.

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

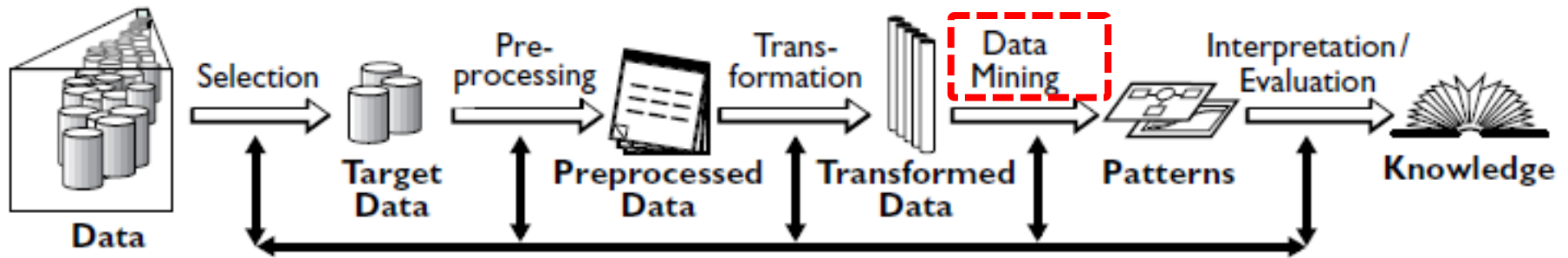


TEHRAN UNIVERSITY

Data Mining

Knowledge Discovery in Databases (KDD) Process

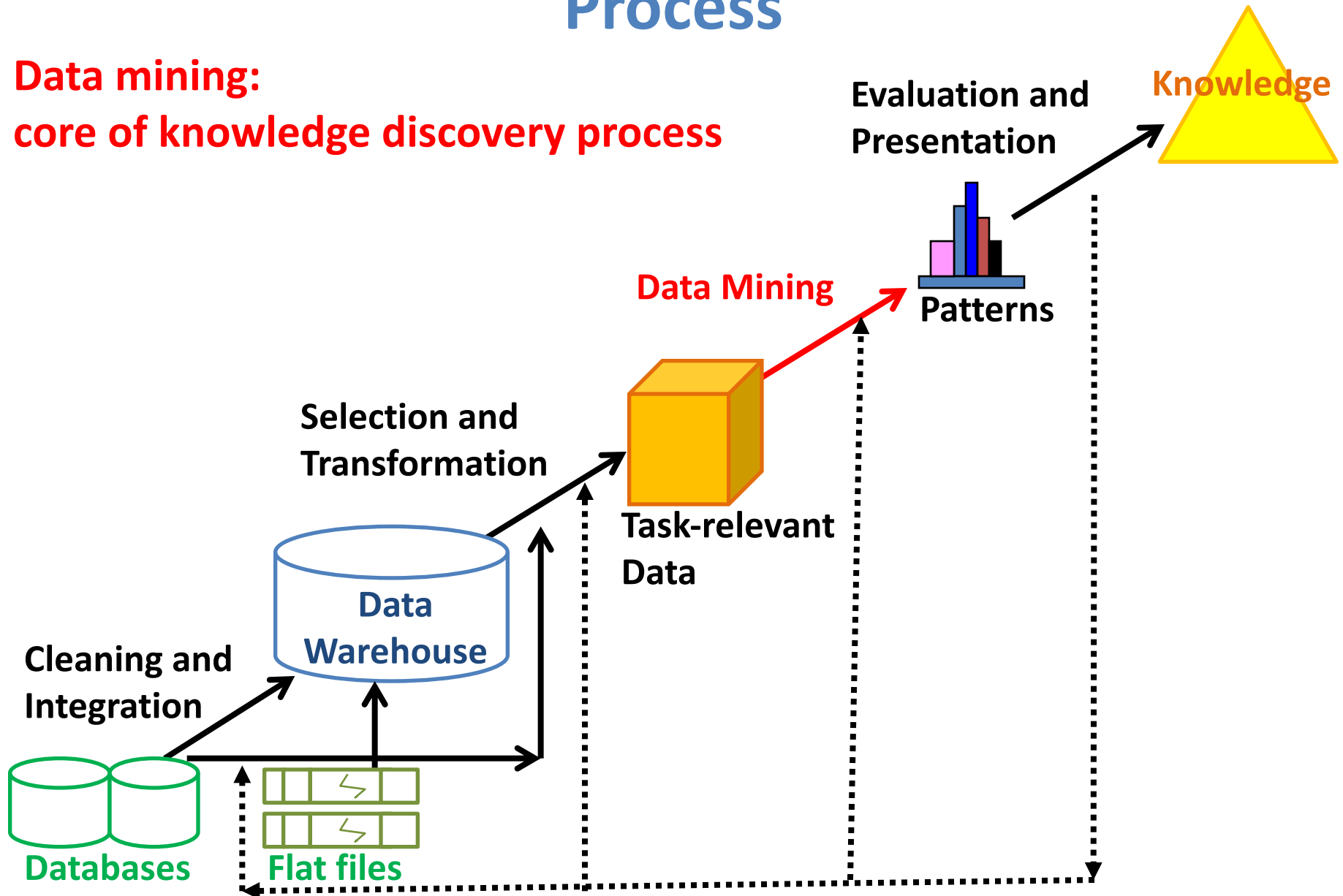
(Fayyad et al., 1996)



Knowledge Discovery in Databases (KDD)

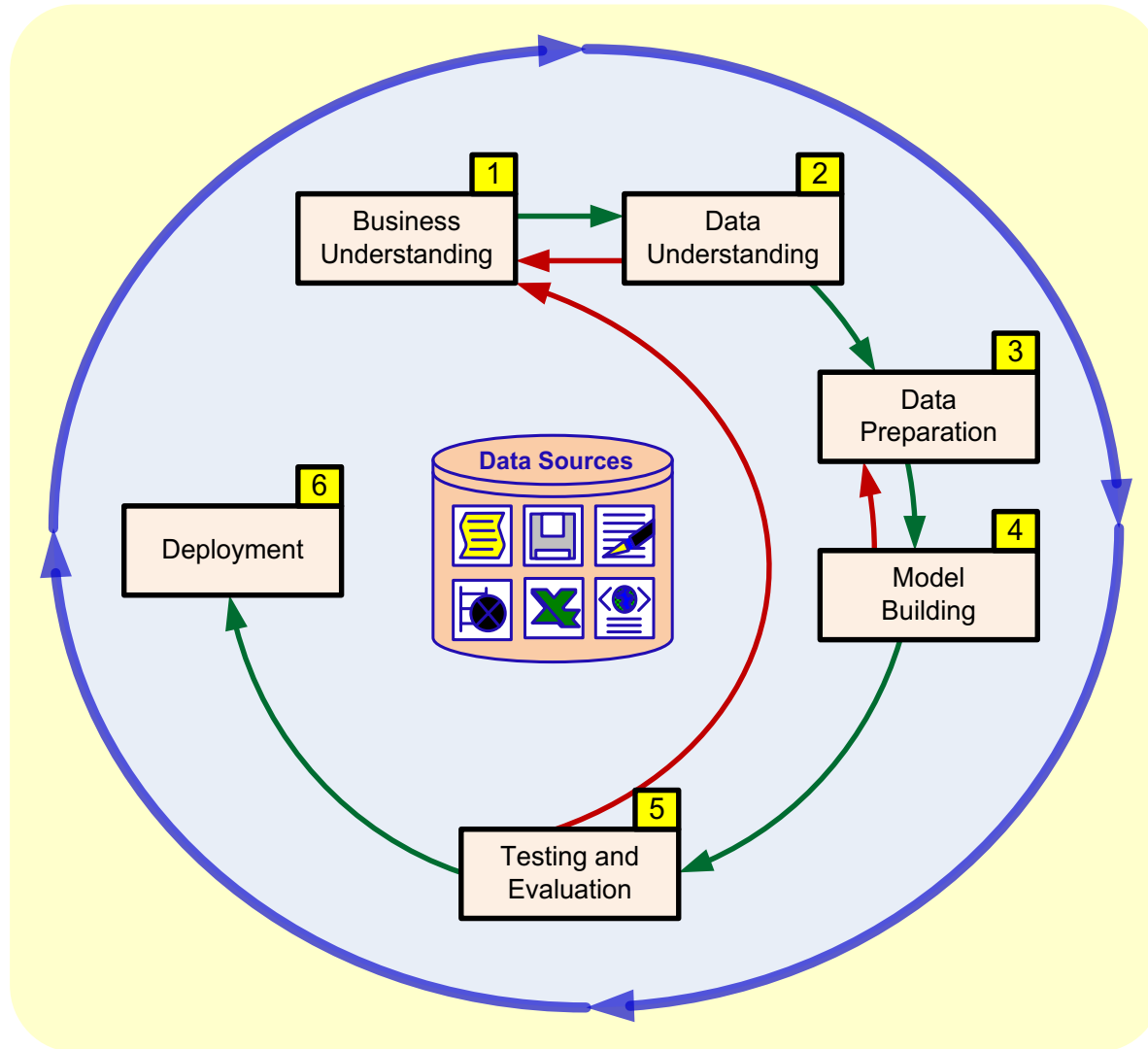
Process

Data mining:
core of knowledge discovery process



Data Mining Process:

CRISP-DM



Data Mining Process:

CRISP-DM

Step 1: Business Understanding

Step 2: Data Understanding

Step 3: Data Preparation (!)

Step 4: Model Building

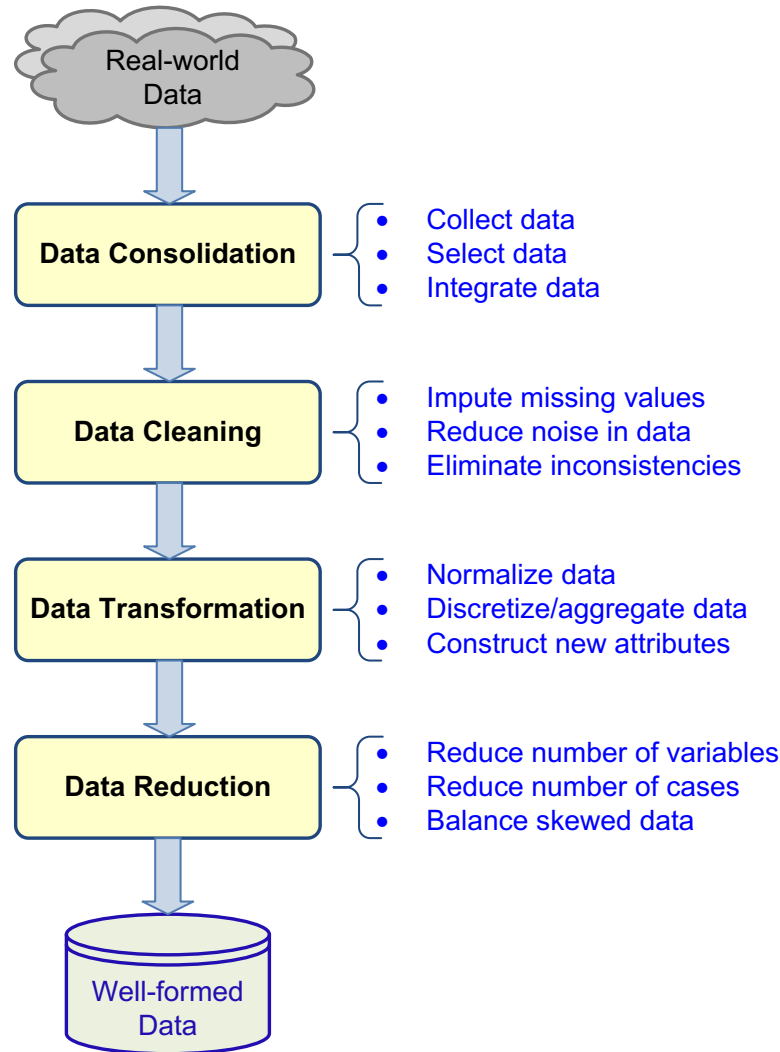
Step 5: Testing and Evaluation

Step 6: Deployment

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

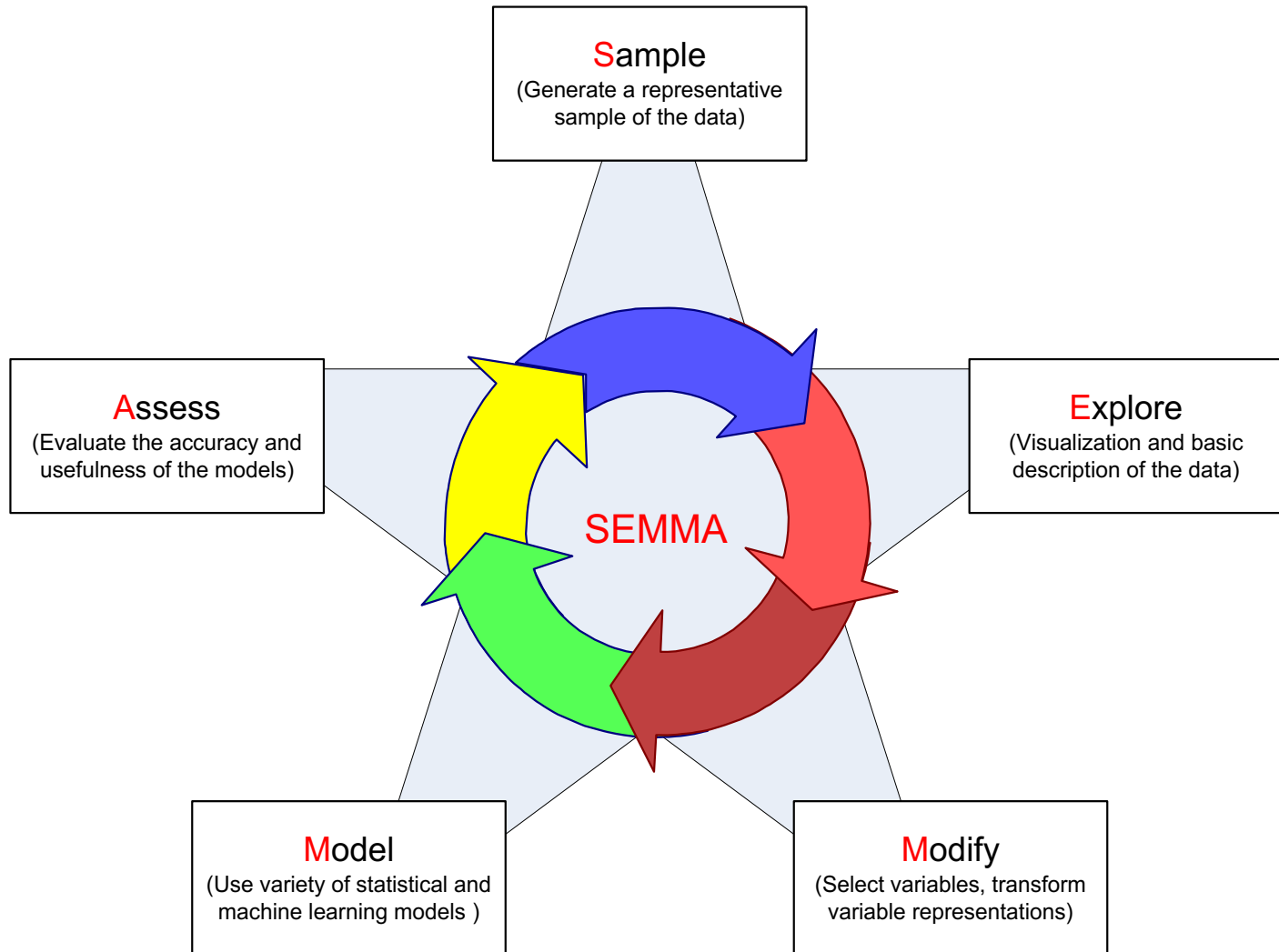
Accounts for
~85% of total
project time

Data Preparation – A Critical DM Task



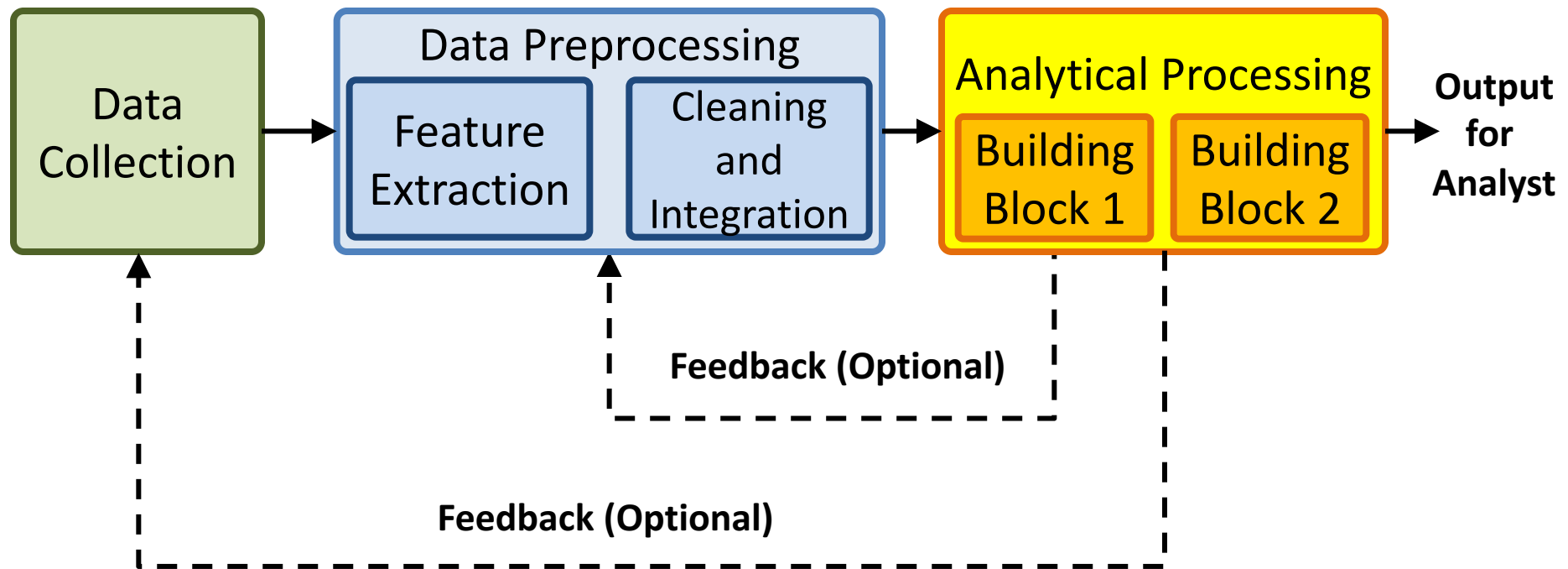
Data Mining Process:

SEMMA

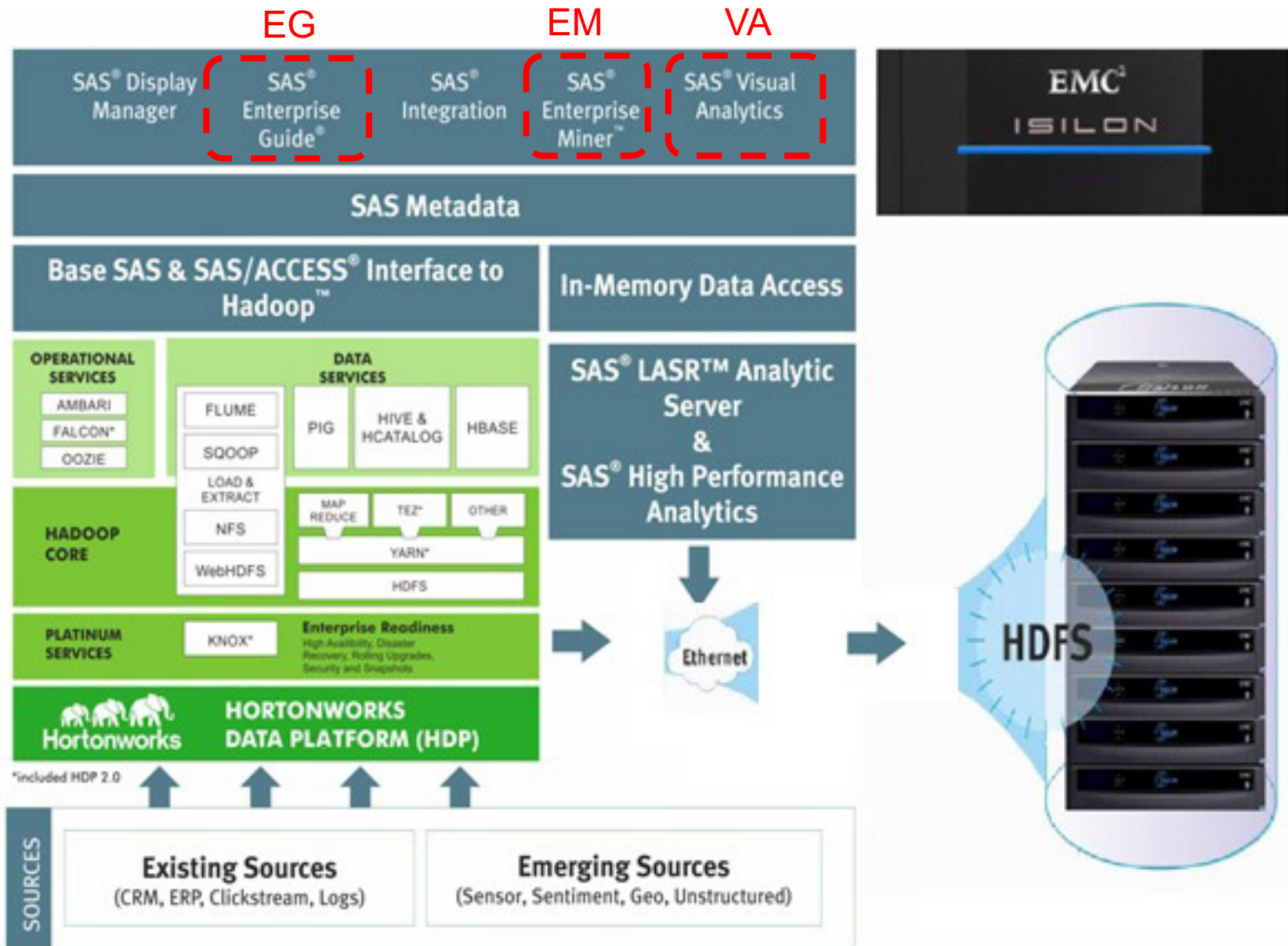


Data Mining Processing Pipeline

(Charu Aggarwal, 2015)



Big Data Solution





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



Complex by nature













DATA

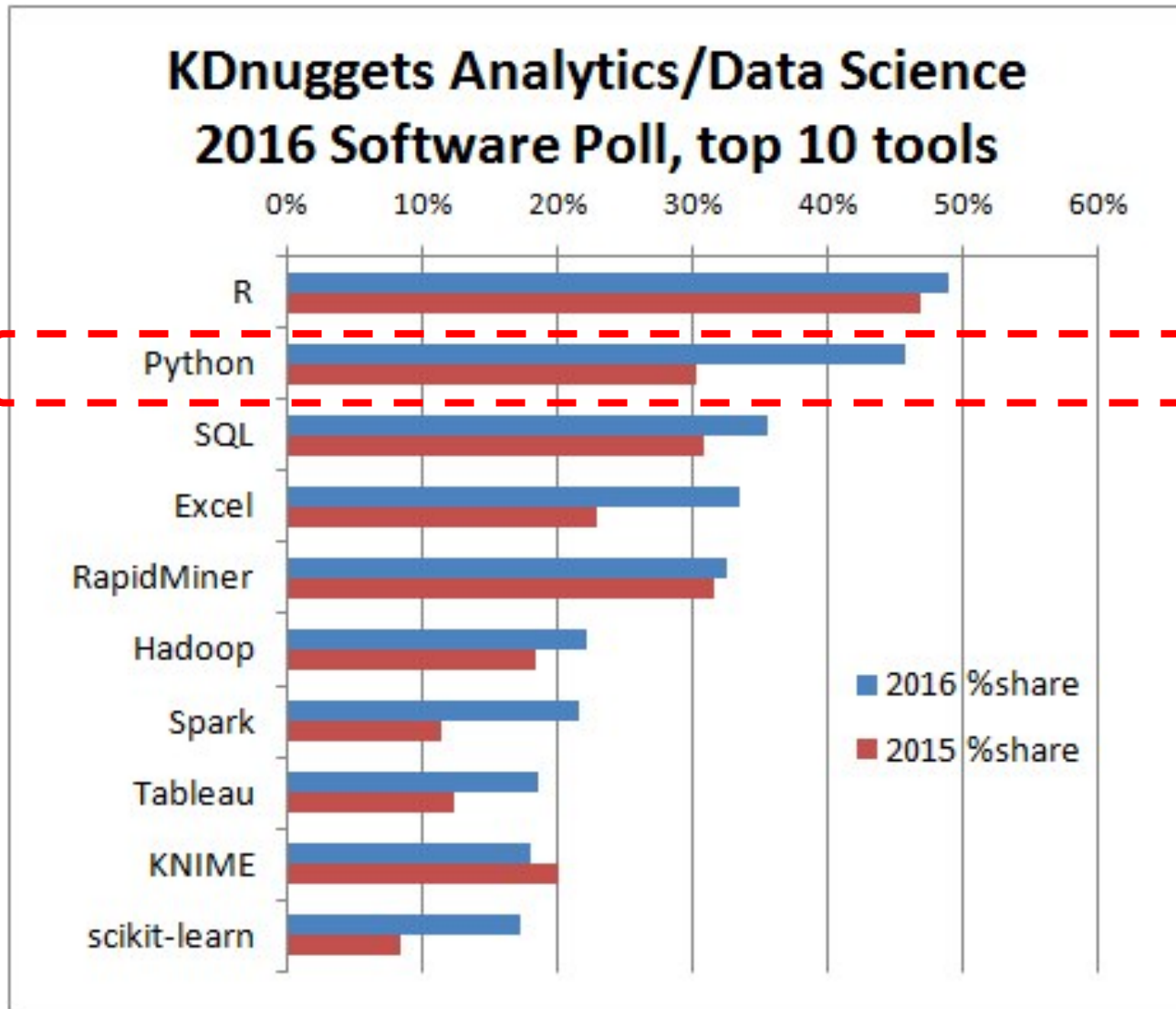
Complex by structure



Python for Big Data Analytics

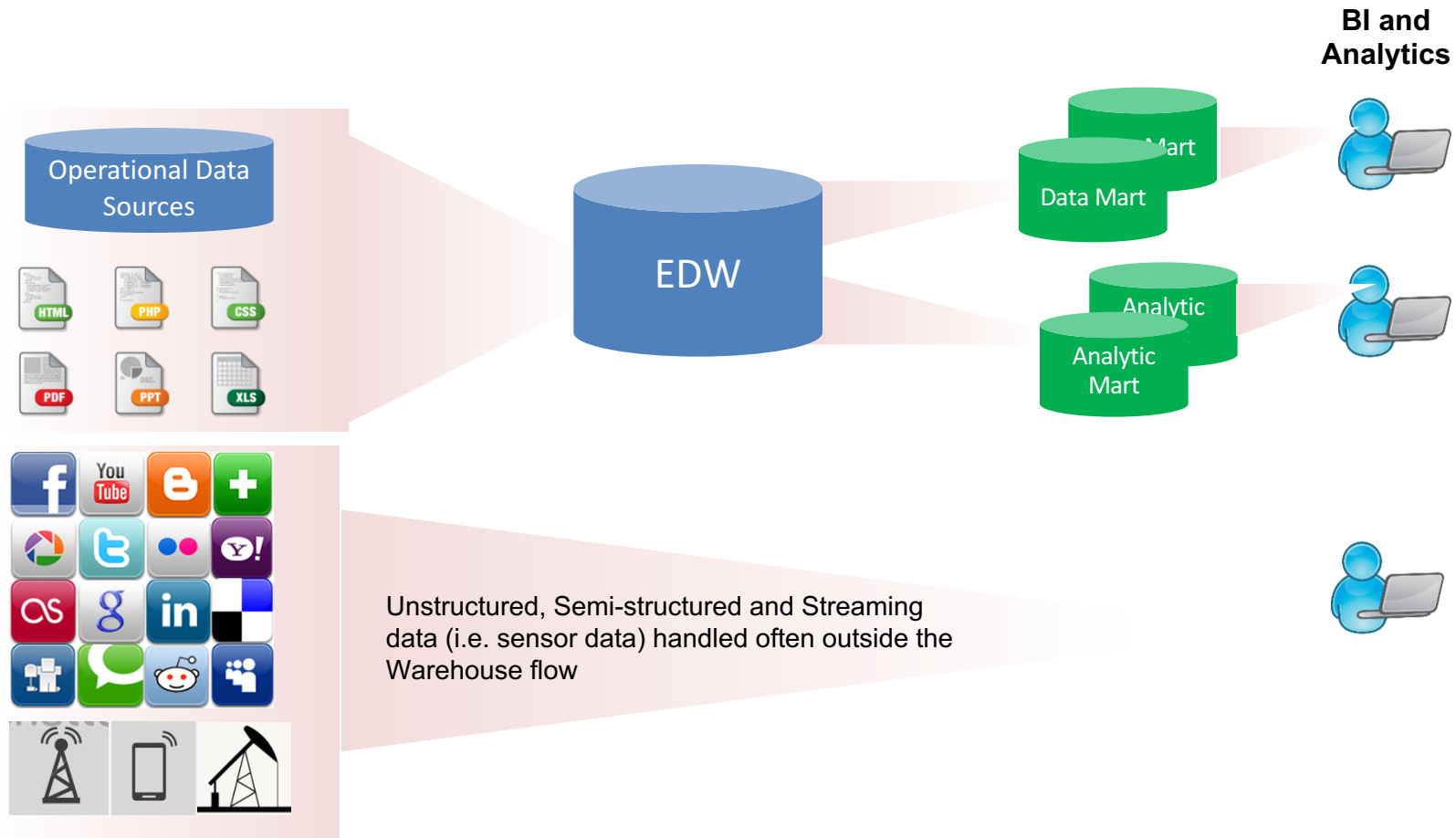
Language Rank	Types	Spectrum Ranking
1. C		100.0
2. Java		98.1
3. Python		98.0
4. C++		95.9
5. R		87.9
6. C#		86.7
7. PHP		82.8
8. JavaScript		82.2
9. Ruby		74.5
10. Go		71.9

Python: Analytics and Data Science Software

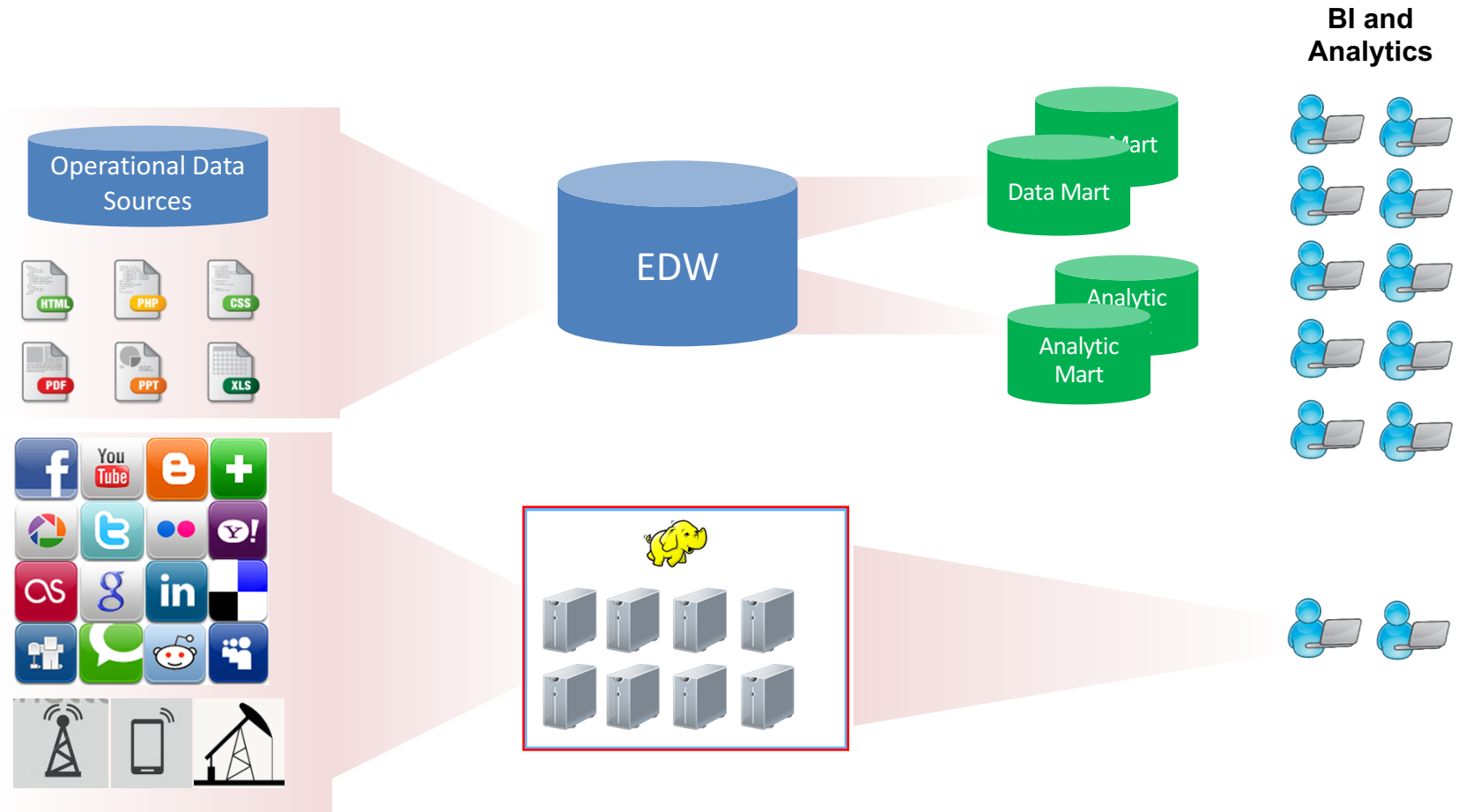


Architectures of Big Data Analytics

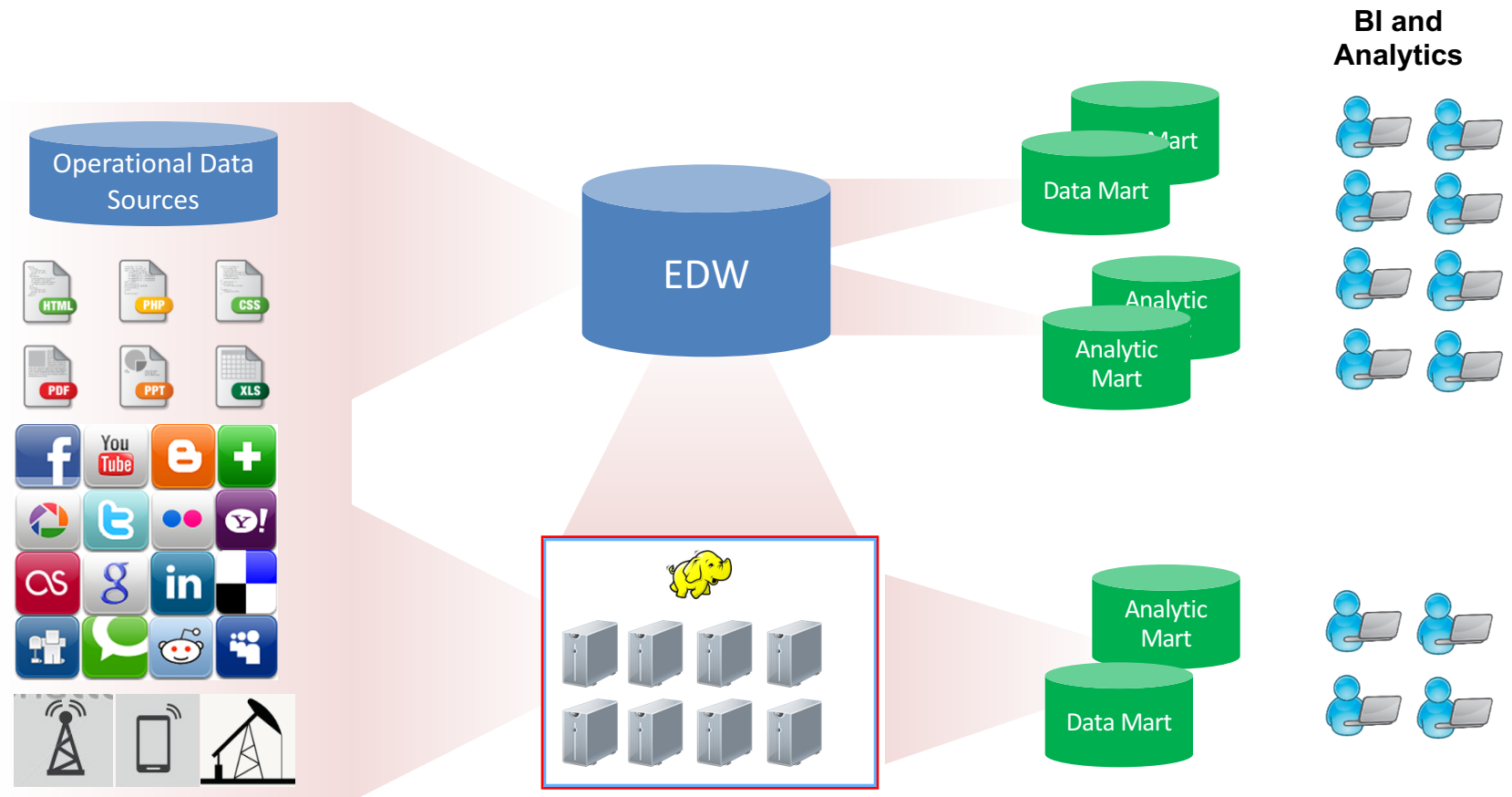
Traditional Analytics



Hadoop as a “new data” Store

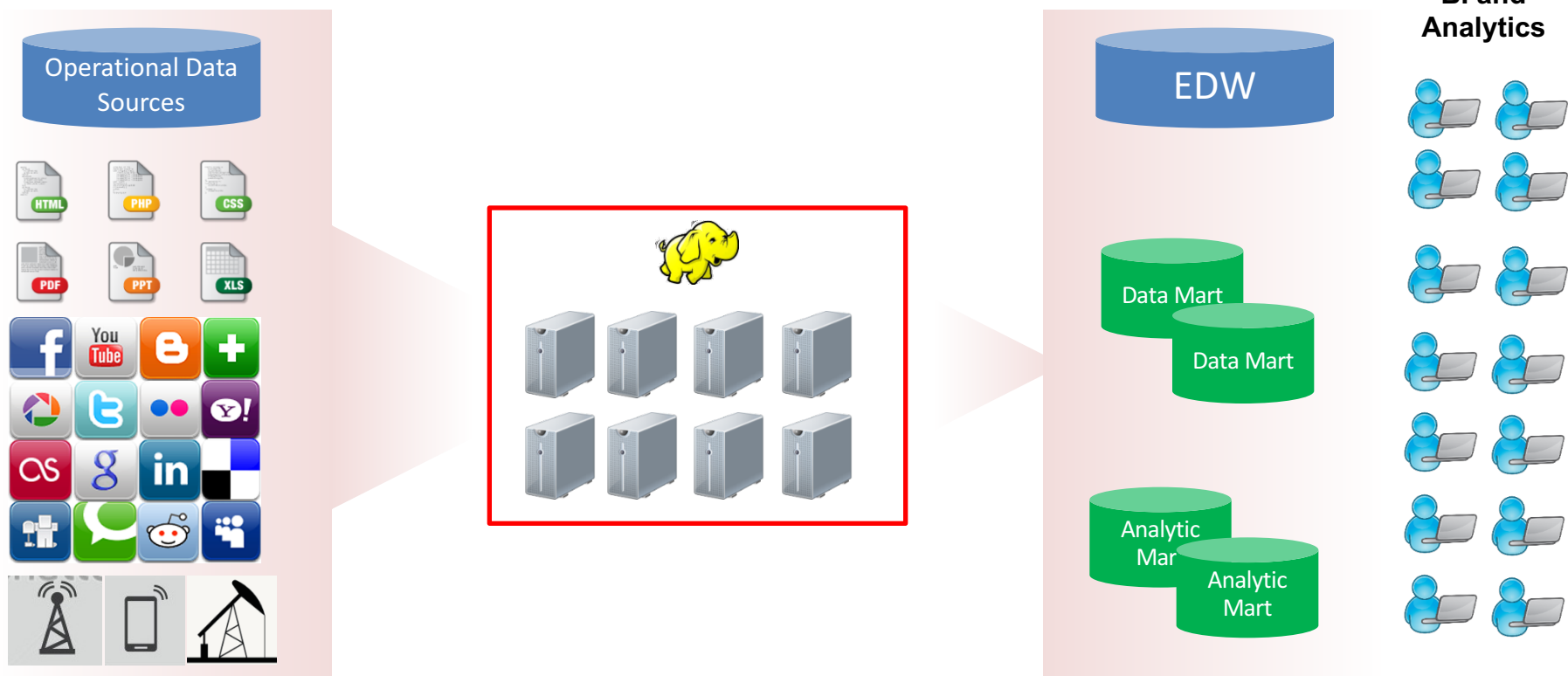


Hadoop as an additional input to the EDW



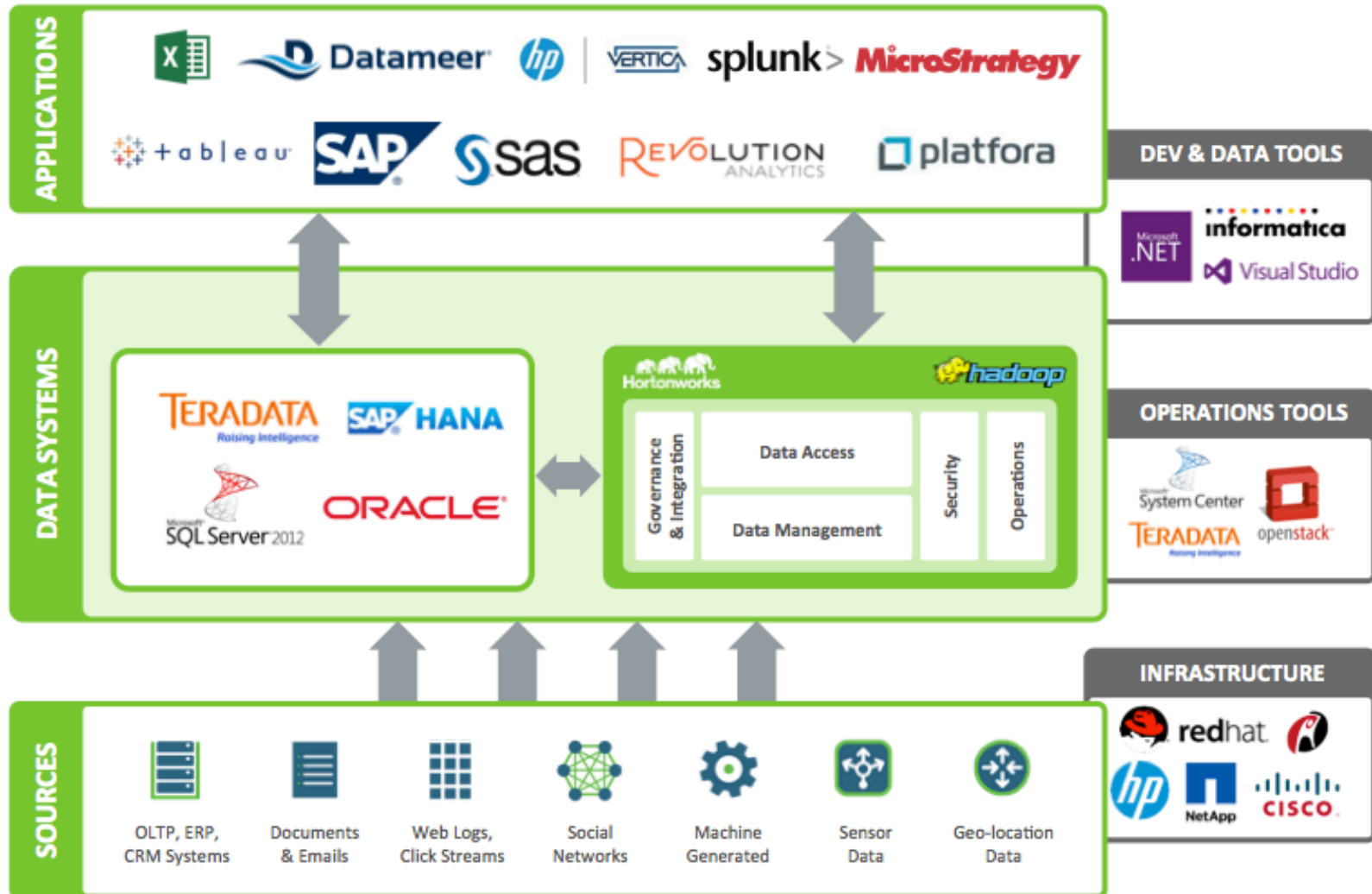
Hadoop Data Platform As a “staging Layer” as part of a “data Lake”

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



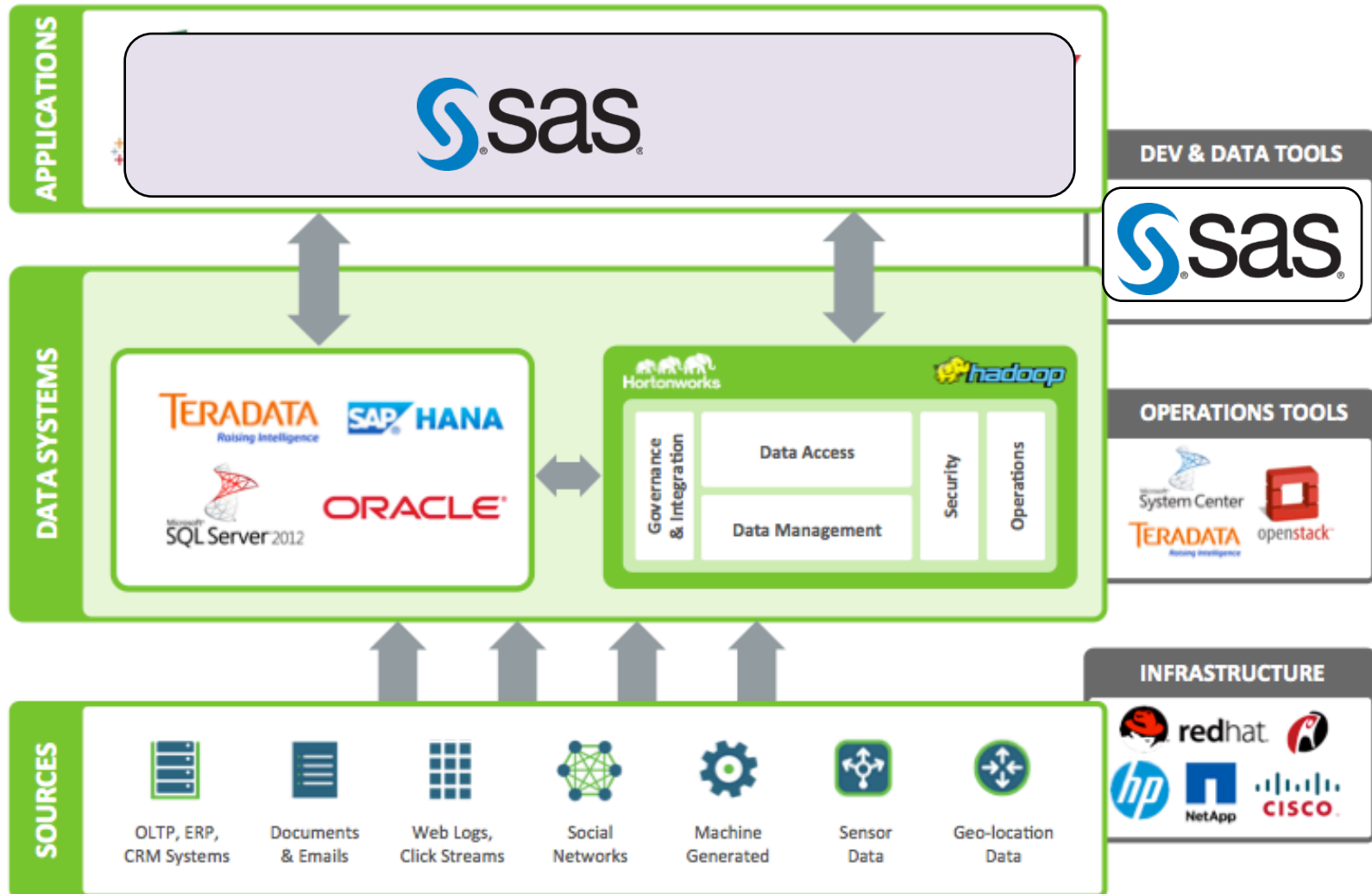
SAS Big data Strategy

– SAS areas

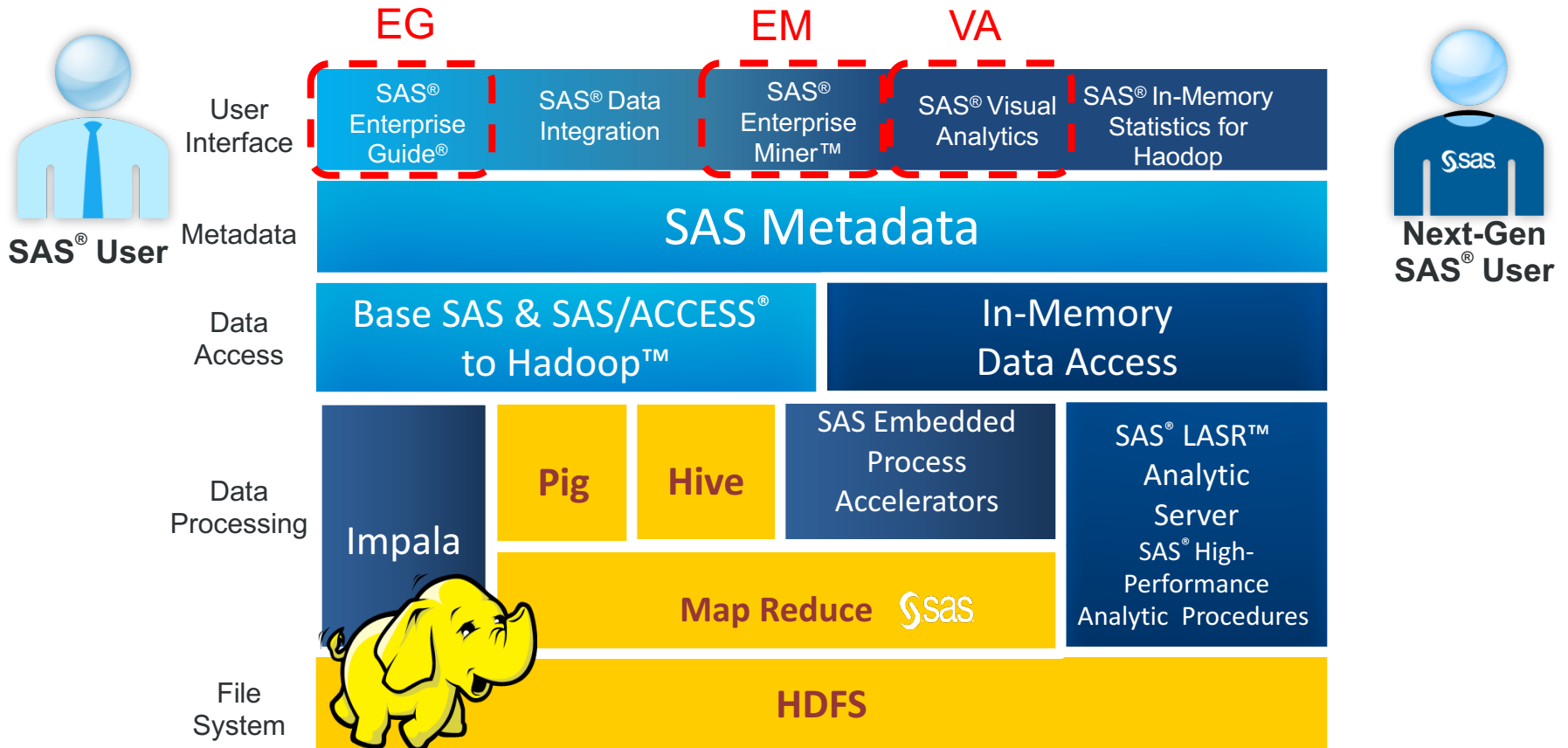


SAS Big data Strategy

– SAS areas

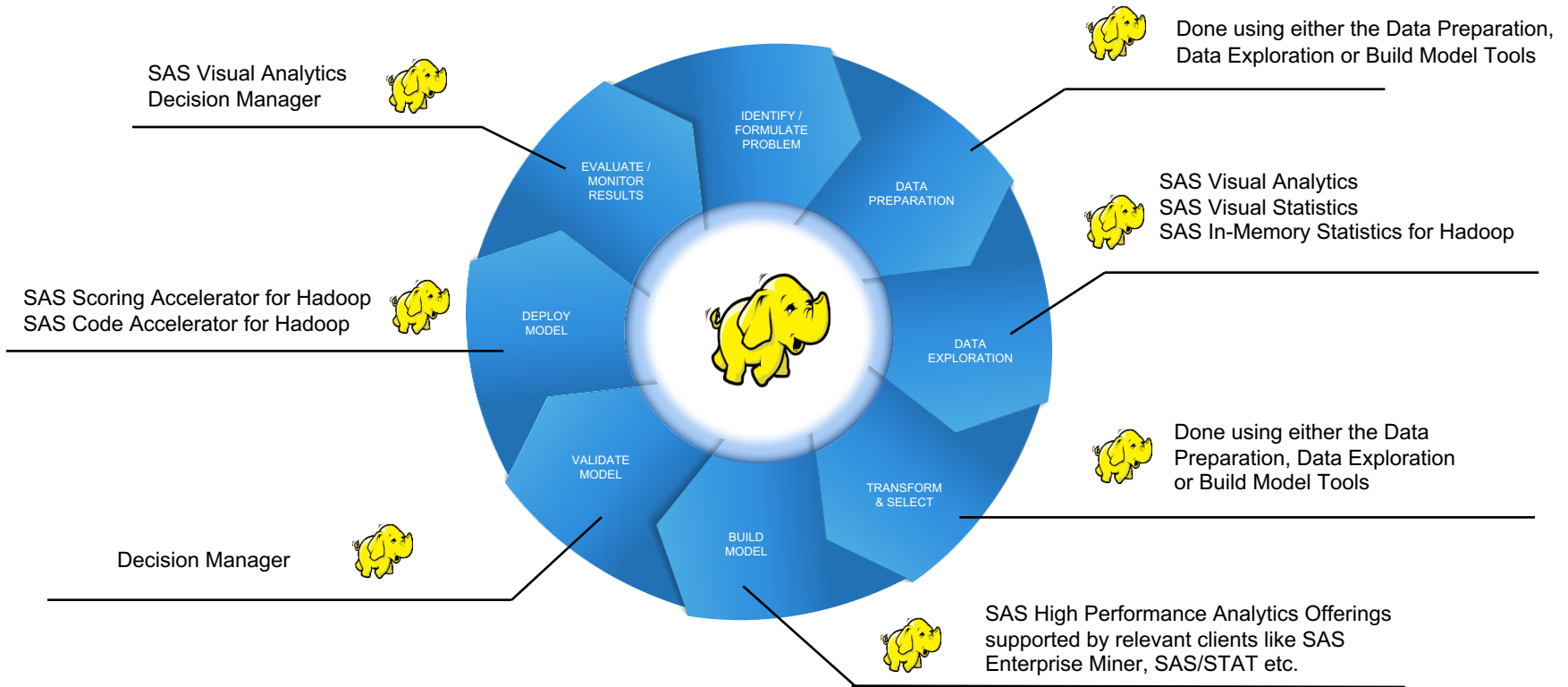


SAS® Within the HADOOP ECOSYSTEM



SAS enables the entire lifecycle around HADOOP

SAS enableS the entire lifecycle around HADOOP



SAS[®] VISUAL ANALYTICS

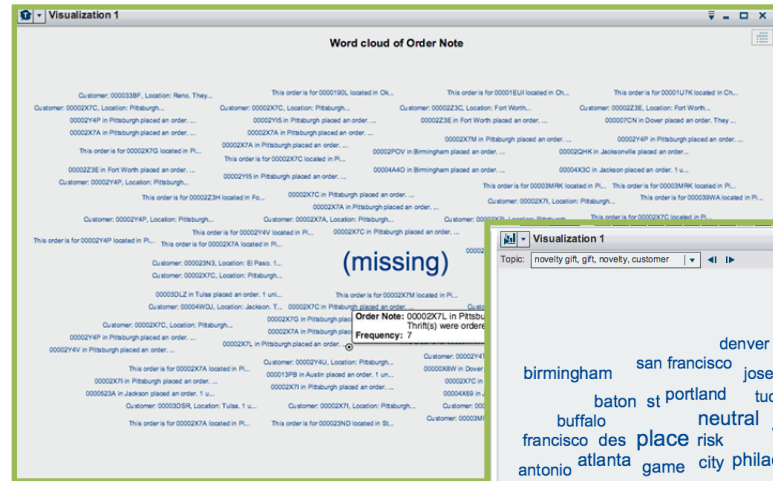
**A Single solution for
Data Discovery,
Visualization, analytics and
reporting**

SAS® VISUAL ANALYTICS

Example: text analysis gives you insight to customer experience and opinion

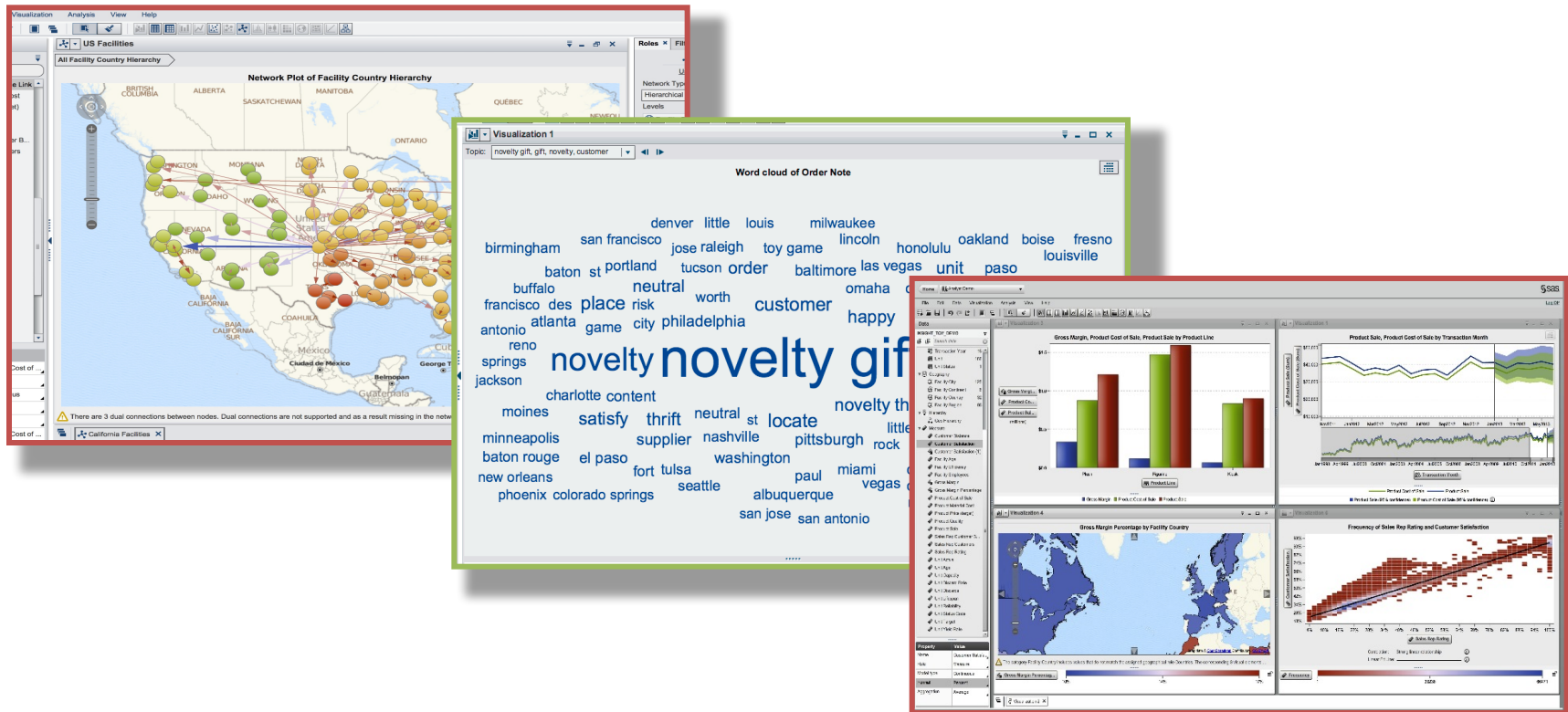


ANALYTICS POWERED



Analytics applied to text provides real MEANING

Visualization



Gephi



[Download](#) [Blog](#) [Wiki](#) [Forum](#) [Support](#) [Bug tracker](#)

[Home](#) [Features](#) [Learn](#) [Develop](#) [Plugins](#) [Services](#) [Consortium](#)

The Open Graph Viz Platform

Gephi is the leading visualization and exploration software for all kinds of graphs and networks. Gephi is open-source and free.

Runs on Windows, Mac OS X and Linux.

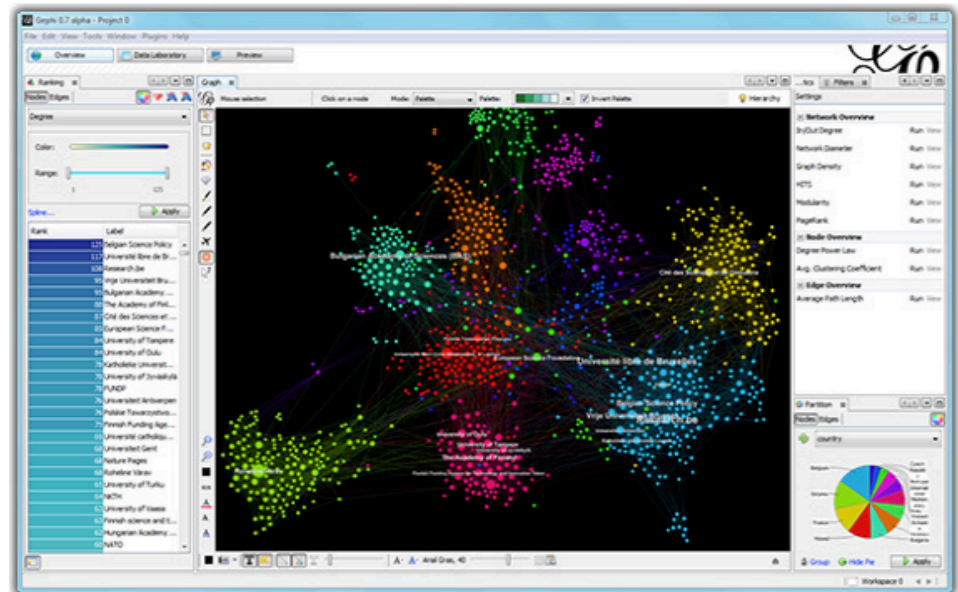
[Learn More on Gephi Platform >](#)



[Release Notes](#) | [System Requirements](#)

► [Features](#)
► [Quick start](#)

► [Screenshots](#)
► [Videos](#)



Support us! We are non-profit. Help us to innovate and empower the community by donating only 8€:

Donate



APPLICATIONS

- ✓ **Exploratory Data Analysis:** intuition-oriented analysis by networks manipulations in real time.
- ✓ **Link Analysis:** revealing the underlying structures of associations between objects.
- ✓ **Social Network Analysis:** easy creation of social

Like Photoshop™ for graphs.

— the Community

LATEST NEWS

► [Gephi updates with 0.9.1 version](#)

PAPERS



<https://gephi.org/>

igraph



igraph – The network analysis package

igraph is a collection of network analysis tools with the emphasis on **efficiency**, **portability** and ease of use. igraph is **open source** and free. igraph can be programmed in **R**, **Python** and **C/C++**.

igraph R package

python-igraph

igraph C library

R/igraph 1.0.0

Repositories at Github

R/igraph 0.7.1

C/igraph 0.7.1

R/igraph 0.7.0

python-igraph 0.7.0

C/igraph 0.7.0

R/igraph 0.6.5

Recent news

R/igraph 1.0.0

June 24, 2015

Release Notes

This is a new major release, with a lot of UI changes. We tried to make it easier to use, with short and easy to remember, consistent function names. Unfortunately

<http://igraph.org/redirect.html>

sigma.js

sigma.js

GET STARTED

FEATURES

USE CASES

TUTORIAL

REFERENCES



Sigma is a JavaScript library **dedicated to graph drawing**. It makes easy to publish networks on Web pages, and allows developers to integrate network exploration in rich Web applications.

FORK ON GITHUB

FOLLOW ON TWITTER

DOWNLOAD V1.1.0

SCIENCESPO - MÉDIALAB

<http://sigmajs.org/>

D3.js

[Overview](#) [Examples](#) [Documentation](#) [Source](#)

Data-Driven Documents

Fork me on GitHub



D3.js is a JavaScript library for manipulating documents based on data. **D3** helps you bring data to life using HTML, SVG, and CSS. D3's emphasis on web standards gives you the full capabilities of

See [more examples](#).

<https://d3js.org/>

Google Charts

Google Charts

Search

All Products



HOME

GUIDES

REFERENCE

SUPPORT

Overview

Hello, Charts!

Quickstart

Load the Charts Library

Prepare the Data

Customize the Chart

Draw the Chart

Draw Multiple Charts

Chart Types

[Chart Gallery](#)

Annotation Charts

Area Charts

Bar Charts

Bubble Charts

Calendar Charts

Candlestick Charts

Column Charts

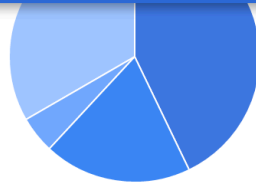
Combo Charts

Diff Charts

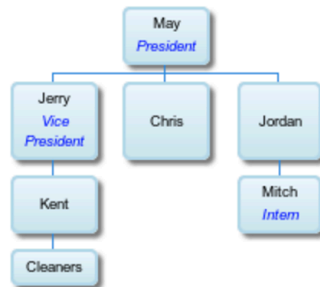
Donut Charts

Gantt Charts

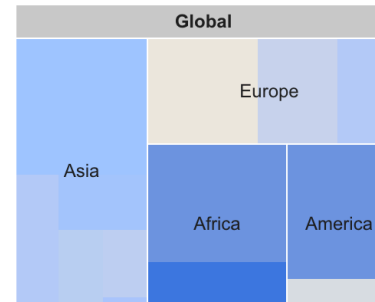
Gauge Charts



Org Chart



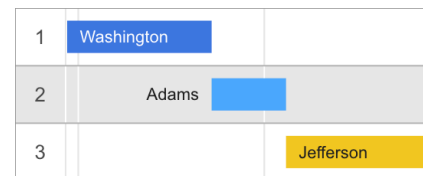
Treemap



Table

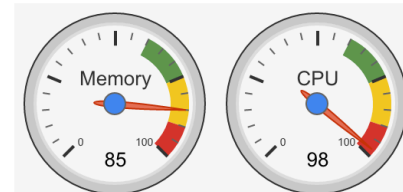
	Name	Salary	Full Time
1	Marie	\$24,700	✓
2	Albert	\$25,200	x
3	Enrico	\$25,700	✓
4	Lise	\$26,600	✓

Timeline

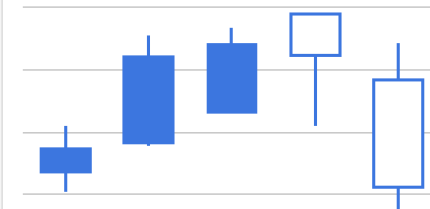


1800

Gauge

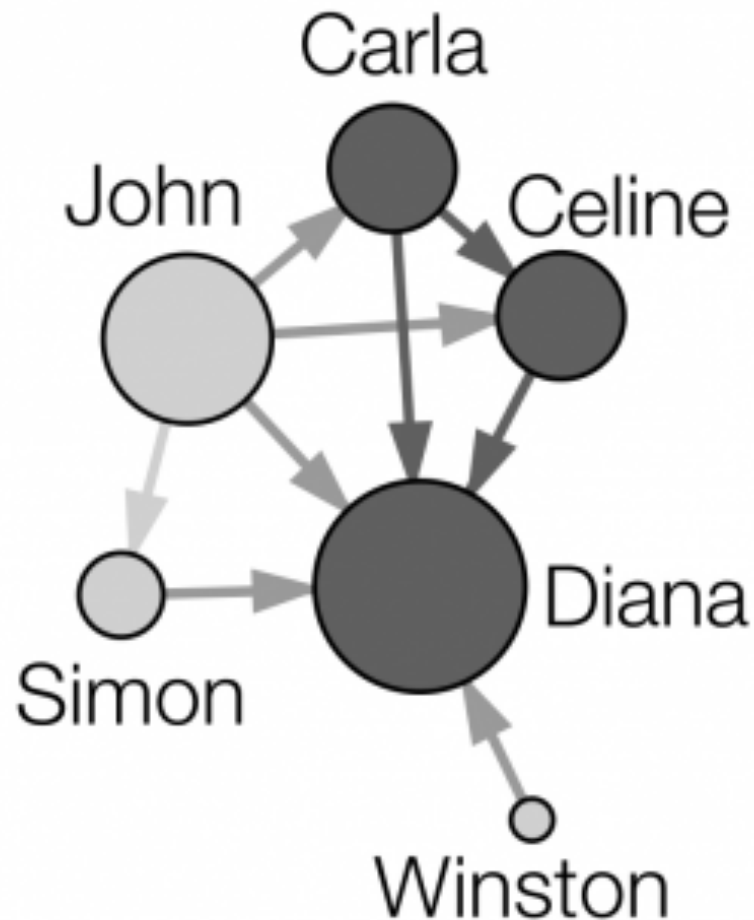


Candlestick Chart



Gephi: Social Network Analysis and Visualization

Network Analysis and Visualization with Gephi



Nodes

Id,Label,Attribute

1,John,1

2,Carla,2

3,Simon,1

4,Celine,2

5,Winston,1

6,Diana,2

Edges

Source,Target

1,2

1,3

1,4

1,6

2,4

2,6

3,6

4,6

5,6

Nodes and Edges

CSV Text Data for Gephi

Nodes1.csv

```
Id,Label,Attribute
1,John,1
2,Carla,2
3,Simon,1
4,Celine,2
5,Winston,1
6,Diana,2
```

Edges1.csv

```
Source,Target
1,2
1,3
1,4
1,6
2,4
2,6
3,6
4,6
5,6
```

Nodes1.csv

×

```
Id,Label,Attribute
1,John,1
2,Carla,2
3,Simon,1
4,Celine,2
5,Winston,1
6,Diana,2
```

Network Visualization with Gephi



Overview Data Laboratory Preview

Workspace 1

Preview Settings Preview

Presets

Default Straight

Settings Manage renderers

opacity	100.0
▼ Node Labels	
Show Labels	<input checked="" type="checkbox"/>
Font	Arial 8 Plain
Proportional size	<input checked="" type="checkbox"/>
Color	custom [0,0...]
Shorten label	<input checked="" type="checkbox"/>
Max characters	14
Outline size	2.0
Outline color	custom [25...]
Outline opacity	80.0
Box	<input type="checkbox"/>
Box color	parent
Box opacity	100.0
▼ Edges	
Show Edges	<input checked="" type="checkbox"/>
Thickness	1.0
Rescale weight	<input type="checkbox"/>
Color	mixed
Opacity	100.0
Curved	<input type="checkbox"/>
Radius	0.0
▼ Edge Arrows	
Size	3.0

Preview ratio: 100%

Refresh

Export: SVG/PDF/PNG

Background Reset zoom - +

```
graph TD; Diana((Diana)) --- Winston((Winston)); Diana --- Celine((Celine)); Diana --- John((John)); Celine --- Diana; Carla((Carla)) --- Diana; John --- Carla; John --- Diana; John --- Simon((Simon)); Simon --- John;
```

References

- EMC Education Services (2015),
Data Science and Big Data Analytics: Discovering,
Analyzing, Visualizing and Presenting Data, Wiley
- SAS Modernization architectures - Big Data Analytics,
<http://www.slideshare.net/deepakramanathan/sas-modernization-architectures-big-data-analytics>