

Big Data Mining

巨量資料探勘

巨量資料基礎：

MapReduce 典範、Hadoop 與 Spark 生態系統

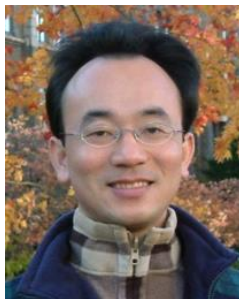
(Fundamental Big Data:

MapReduce Paradigm, Hadoop and Spark Ecosystem)

1042DM02

MI4 (M2244) (3094)

Tue, 3, 4 (10:10-12:00) (B216)



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/>

2016-02-23



課程大綱 (Syllabus)

週次 (Week)	日期 (Date)	內容 (Subject/Topics)
1	2016/02/16	巨量資料探勘課程介紹 (Course Orientation for Big Data Mining)
2	2016/02/23	巨量資料基礎：MapReduce典範、Hadoop與Spark生態系統 (Fundamental Big Data: MapReduce Paradigm, Hadoop and Spark Ecosystem)
3	2016/03/01	關連分析 (Association Analysis)
4	2016/03/08	分類與預測 (Classification and Prediction)
5	2016/03/15	分群分析 (Cluster Analysis)
6	2016/03/22	個案分析與實作一 (SAS EM 分群分析)： Case Study 1 (Cluster Analysis – K-Means using SAS EM)
7	2016/03/29	個案分析與實作二 (SAS EM 關連分析)： Case Study 2 (Association Analysis using SAS EM)

課程大綱 (Syllabus)

週次 (Week)	日期 (Date)	內容 (Subject/Topics)
8	2016/04/05	教學行政觀摩日 (Off-campus study)
9	2016/04/12	期中報告 (Midterm Project Presentation)
10	2016/04/19	期中考試週 (Midterm Exam)
11	2016/04/26	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	2016/05/03	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	2016/05/10	Google TensorFlow 深度學習 (Deep Learning with Google TensorFlow)
14	2016/05/17	期末報告 (Final Project Presentation)
15	2016/05/24	畢業班考試 (Final Exam)

2016/02/23

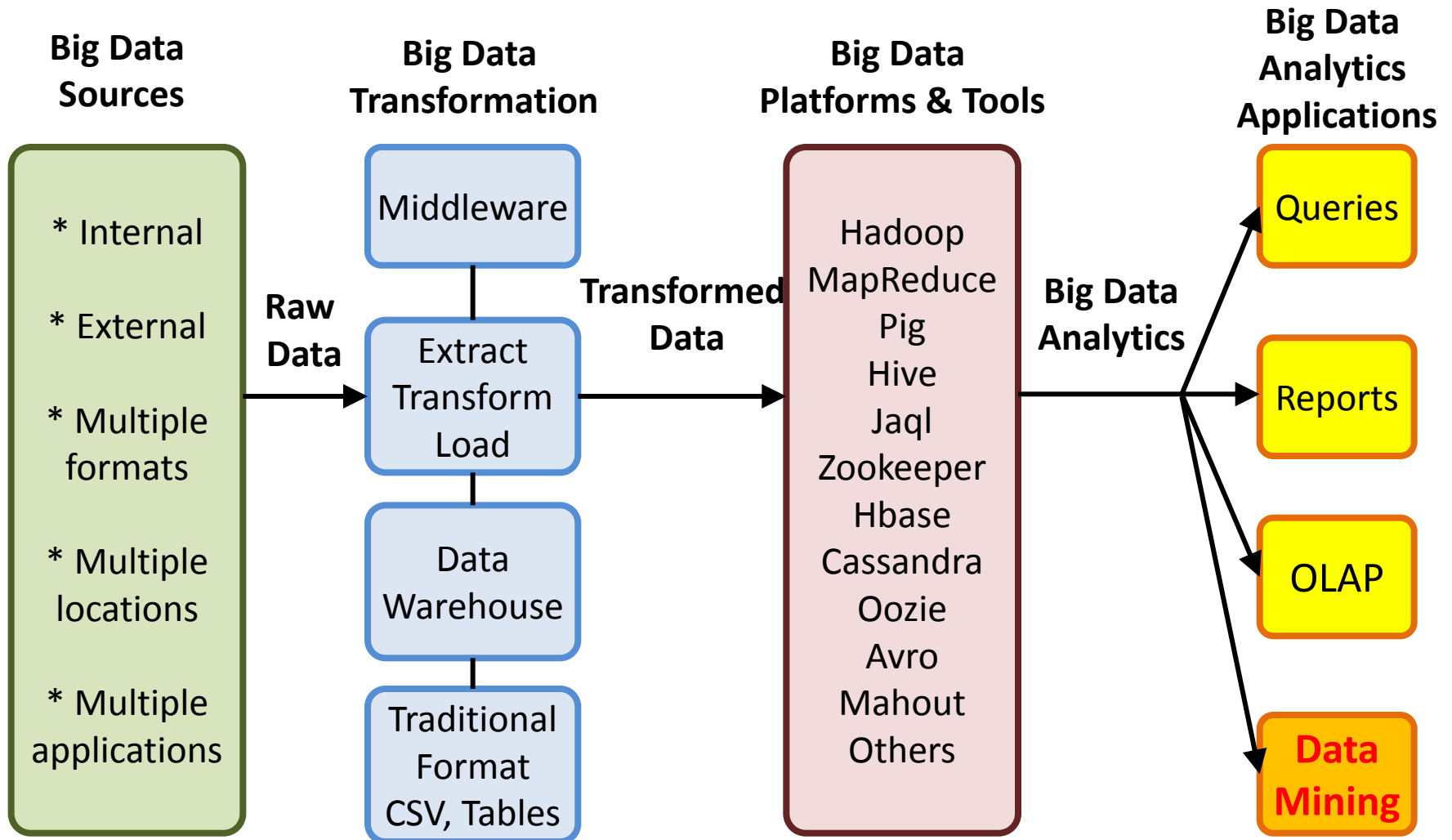
巨量資料基礎：

**MapReduce典範、
Hadoop與Spark生態系統**

(Fundamental Big Data:

**MapReduce Paradigm,
Hadoop and Spark Ecosystem)**

Architecture of Big Data Analytics



Architecture of Big Data Analytics



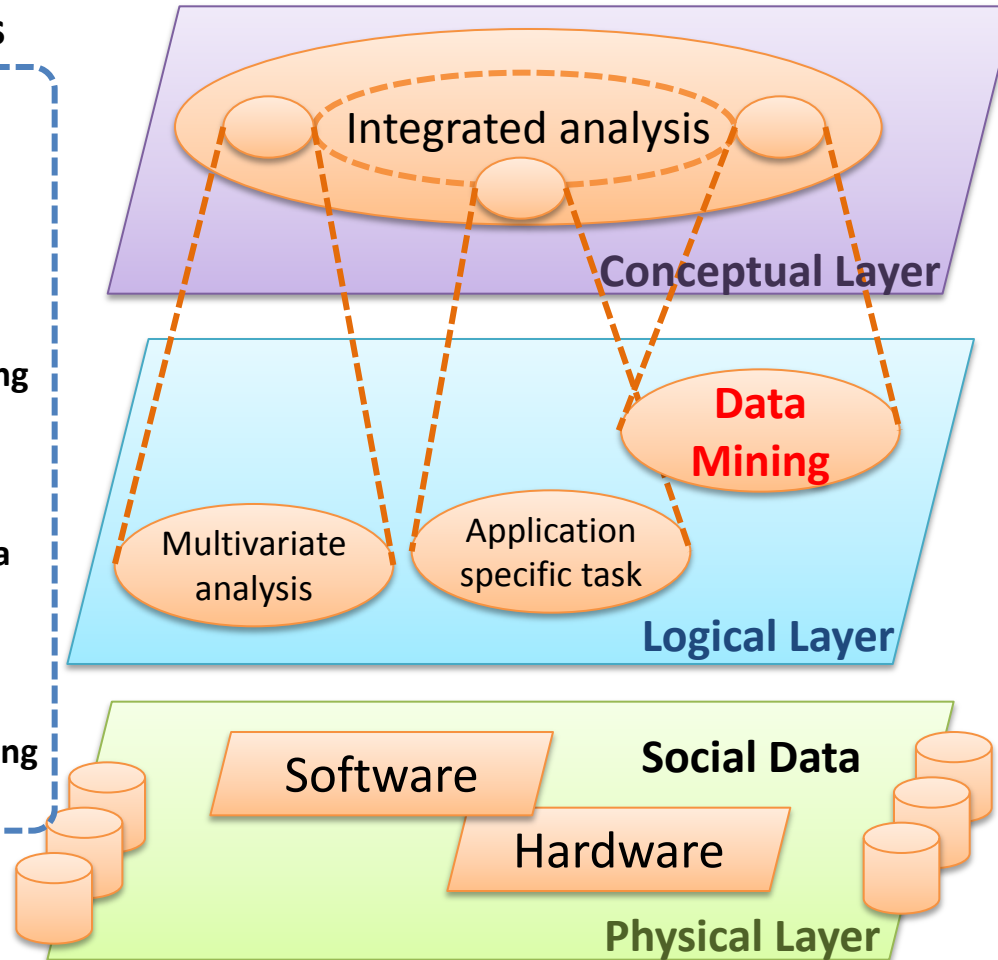
Source: Stephan Kudyba (2014), Big Data, Mining, and Analytics: Components of Strategic Decision Making, Auerbach Publications

Architecture for Social Big Data Mining

(Hiroshi Ishikawa, 2015)

Enabling Technologies

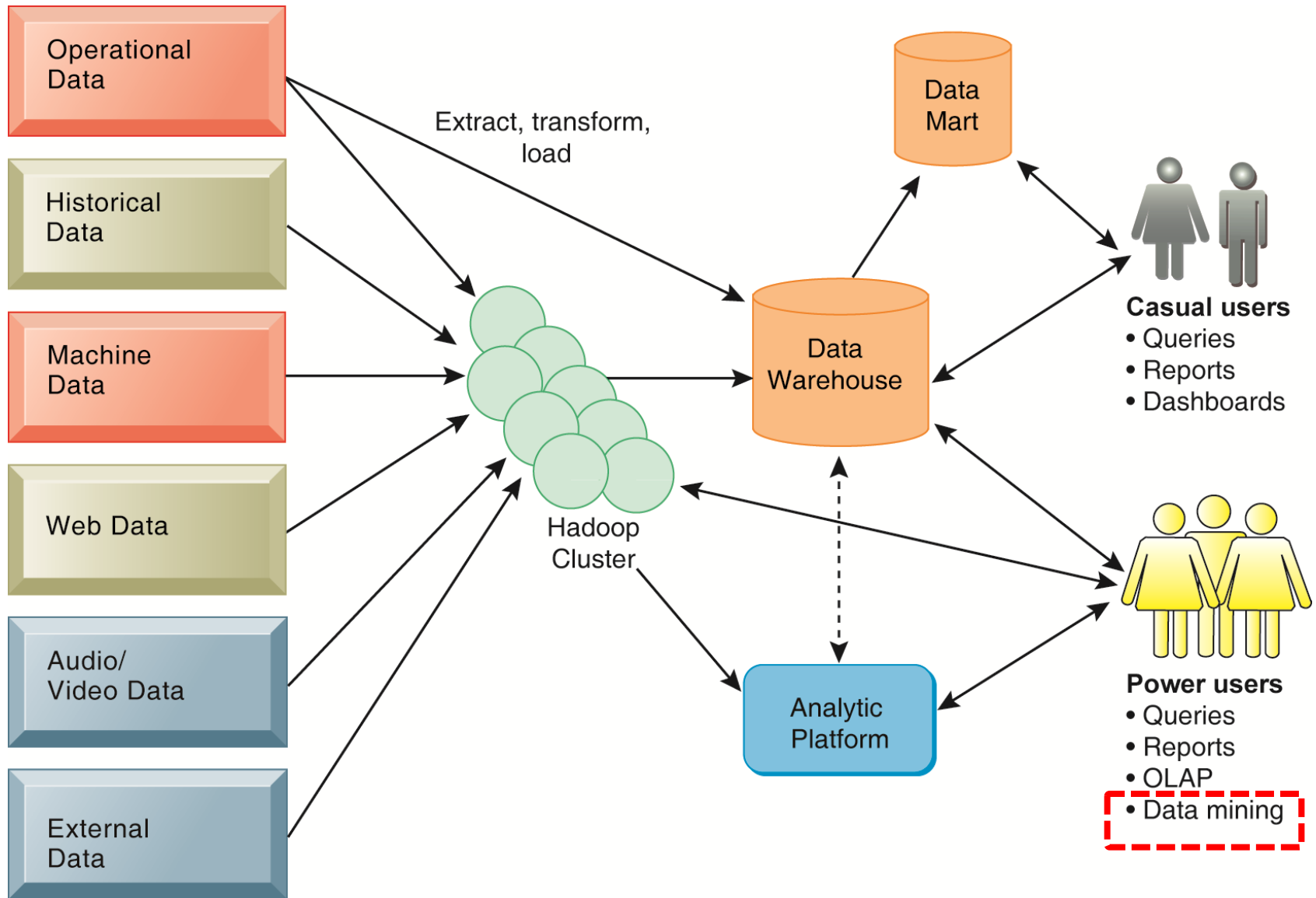
- Integrated analysis model
- Natural Language Processing
- Information Extraction
- Anomaly Detection
- Discovery of relationships among heterogeneous data
- Large-scale visualization
- Parallel distributed processing



Analysts

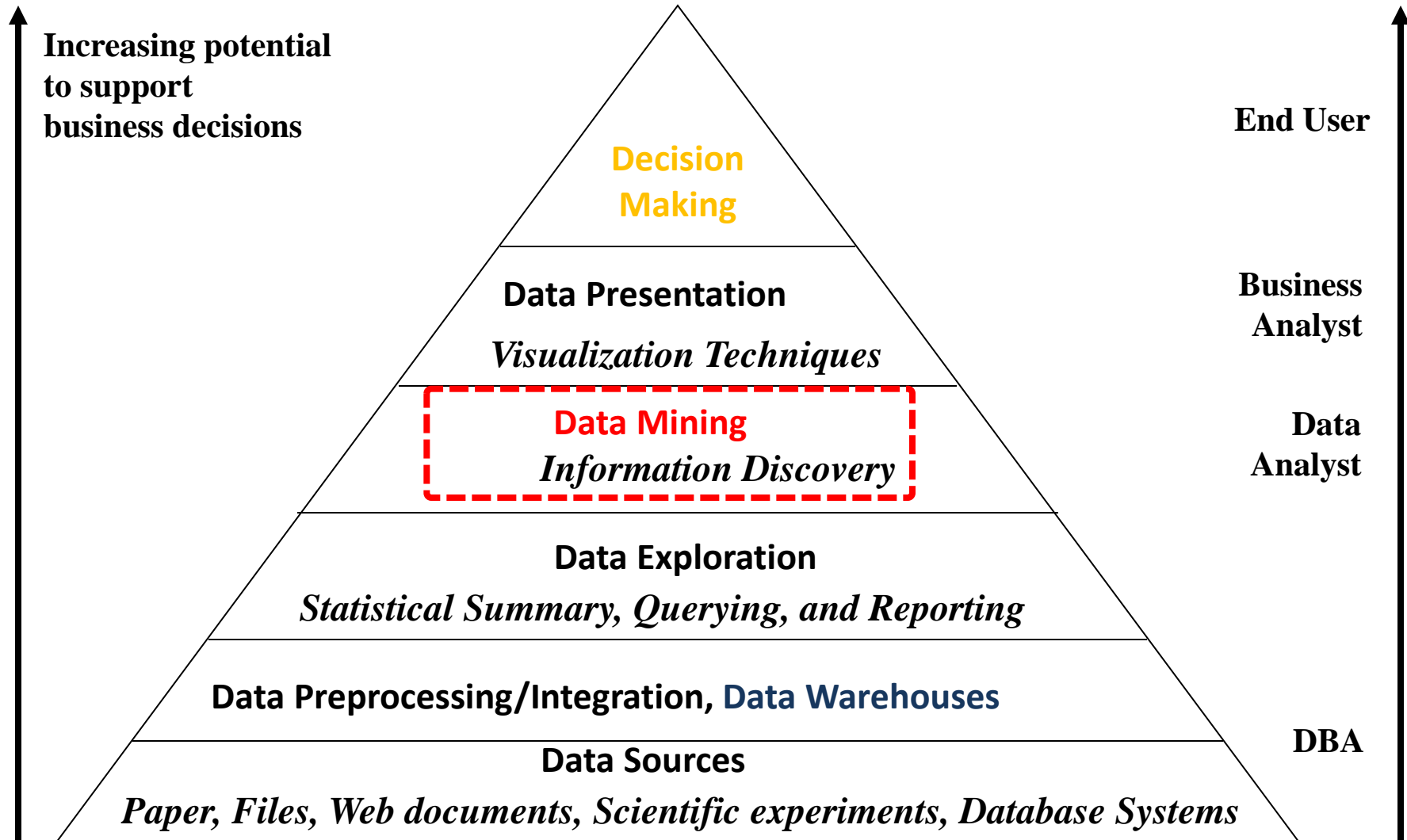
- Model Construction
- Explanation by Model
- Construction and confirmation of individual hypothesis
- Description and execution of application-specific task

Business Intelligence (BI) Infrastructure



Data Warehouse

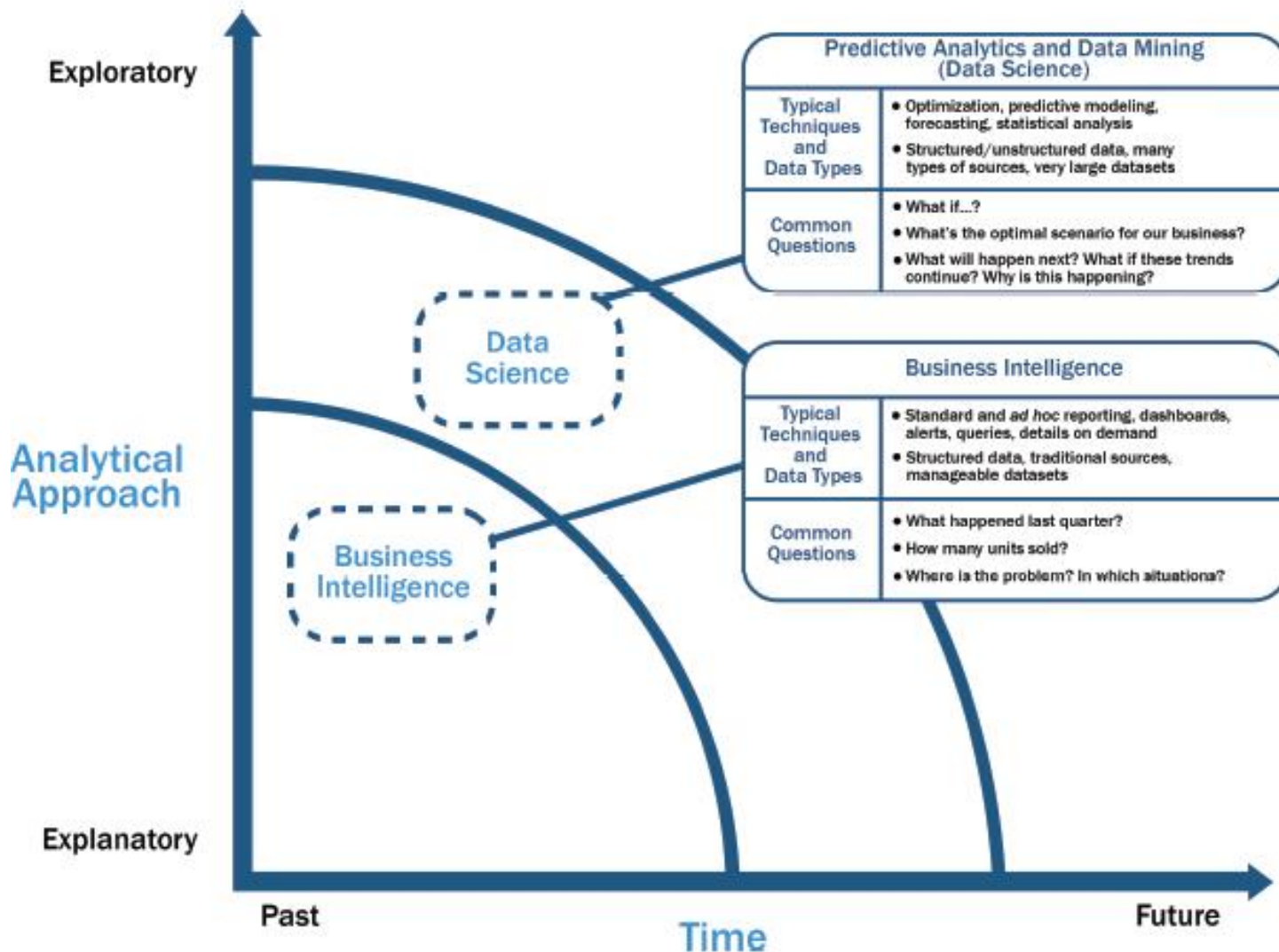
Data Mining and Business Intelligence



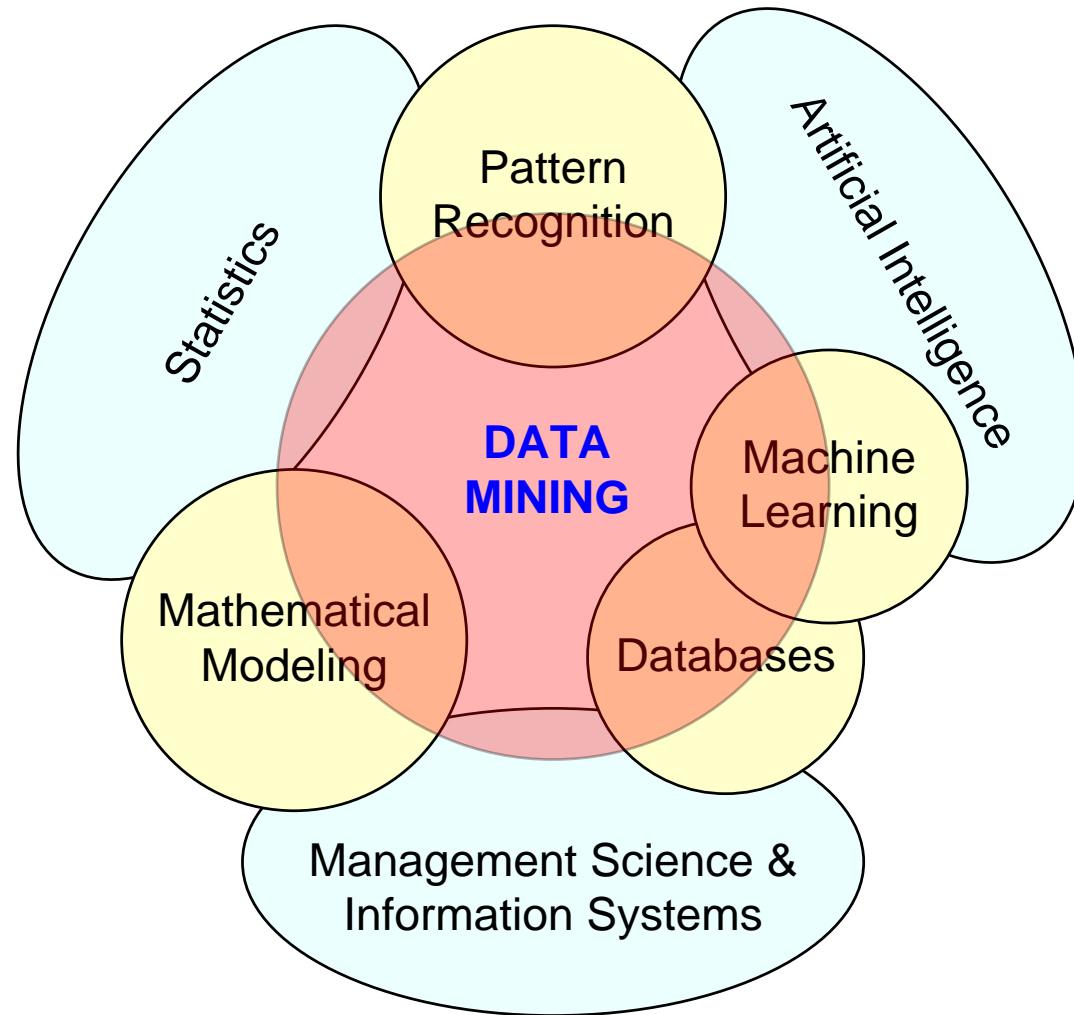
The Evolution of BI Capabilities



Data Science and Business Intelligence

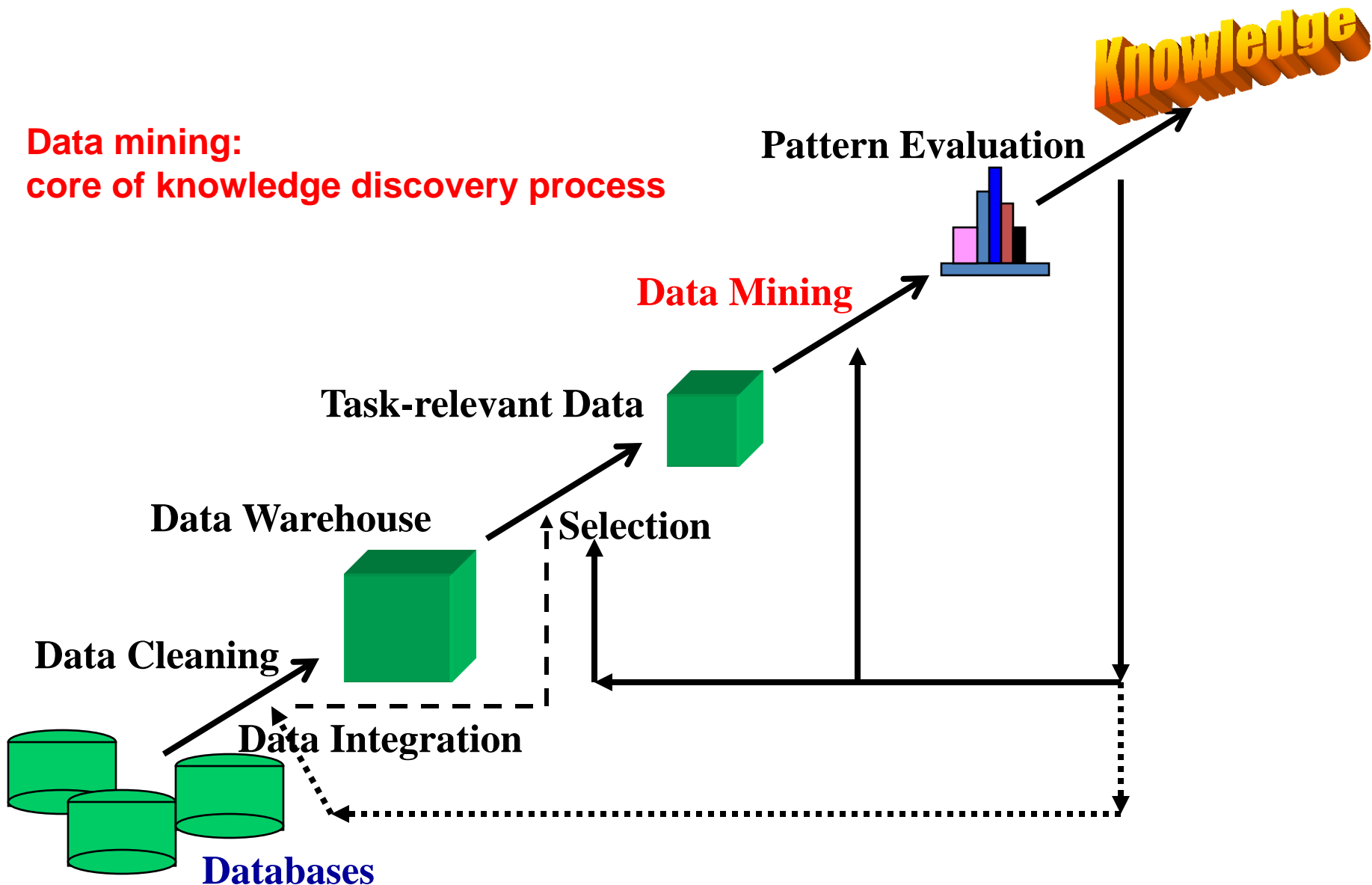


Data Mining at the Intersection of Many Disciplines

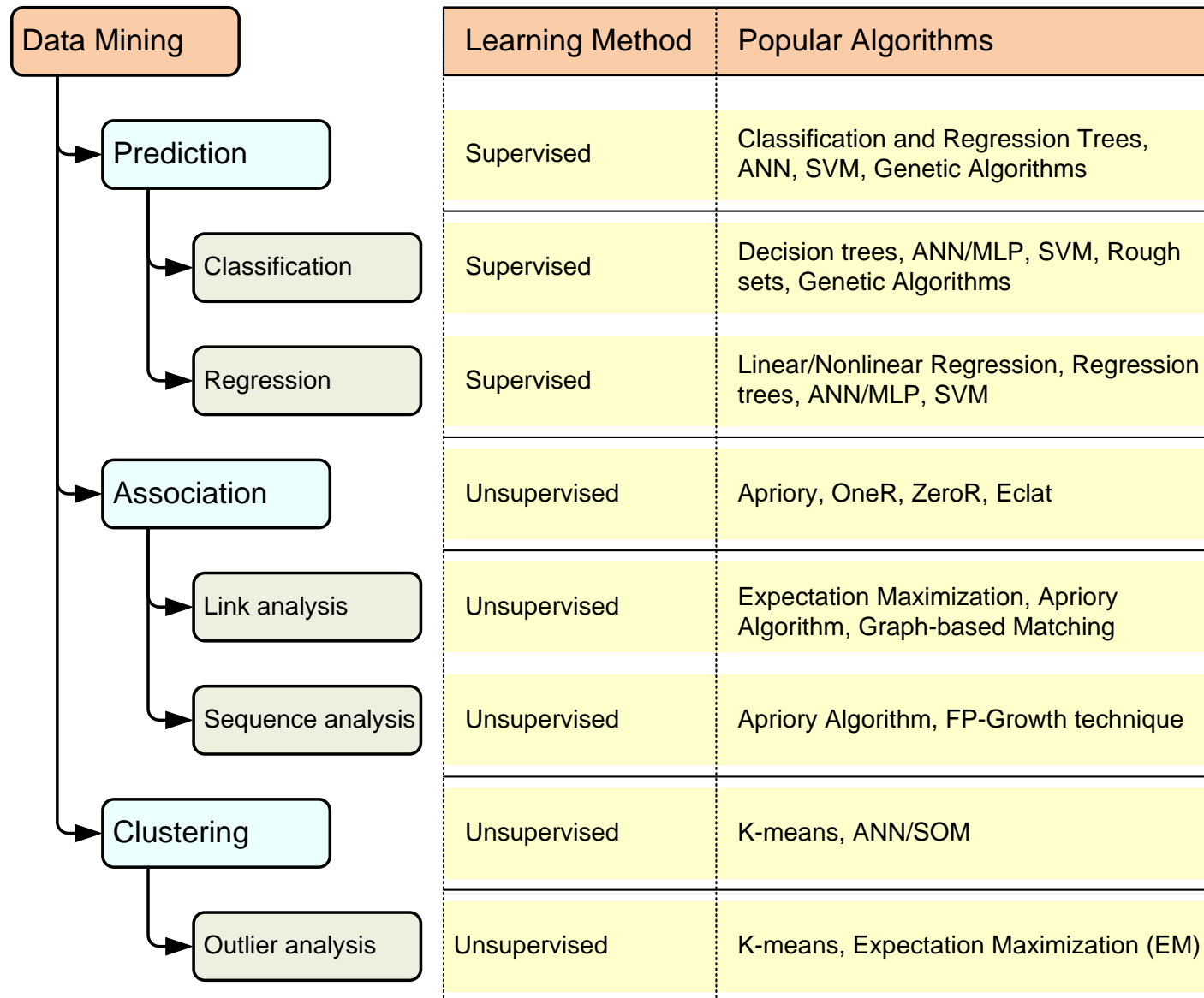


Knowledge Discovery (KDD) Process

Data mining:
core of knowledge discovery process



A Taxonomy for Data Mining Tasks



Deep Learning

Intelligence from Big Data



Big Data



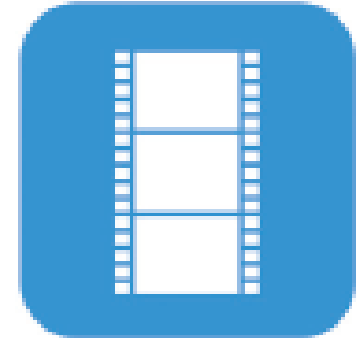
**Mobile
Sensors**



**Social
Media**



**Video
Surveillance**



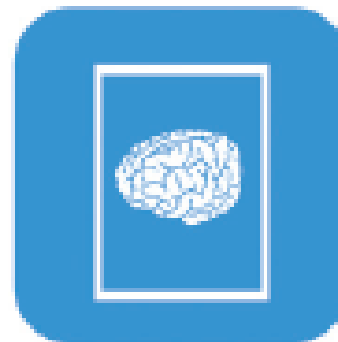
**Video
Rendering**



**Smart
Grids**



**Geophysical
Exploration**

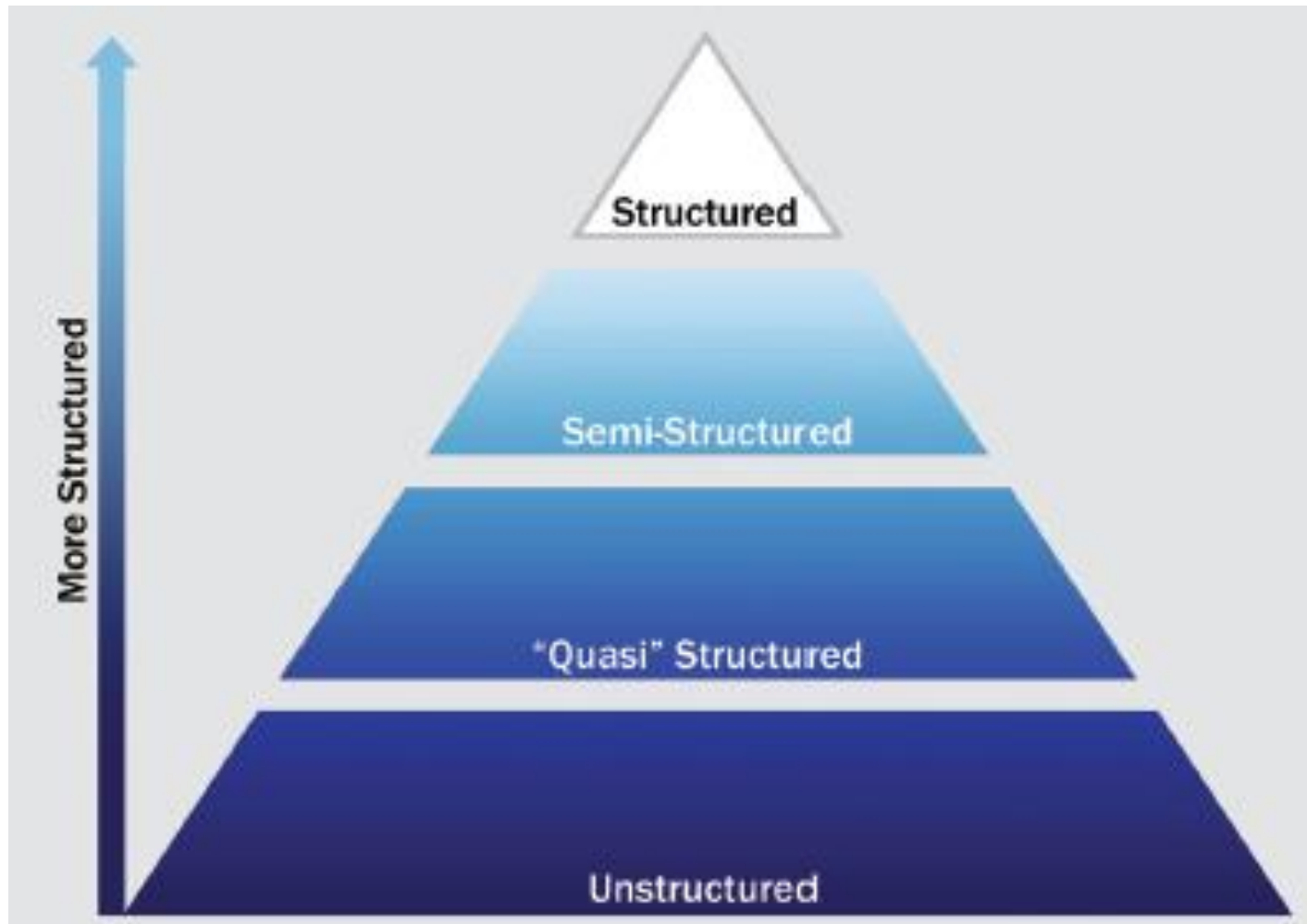


**Medical
Imaging**

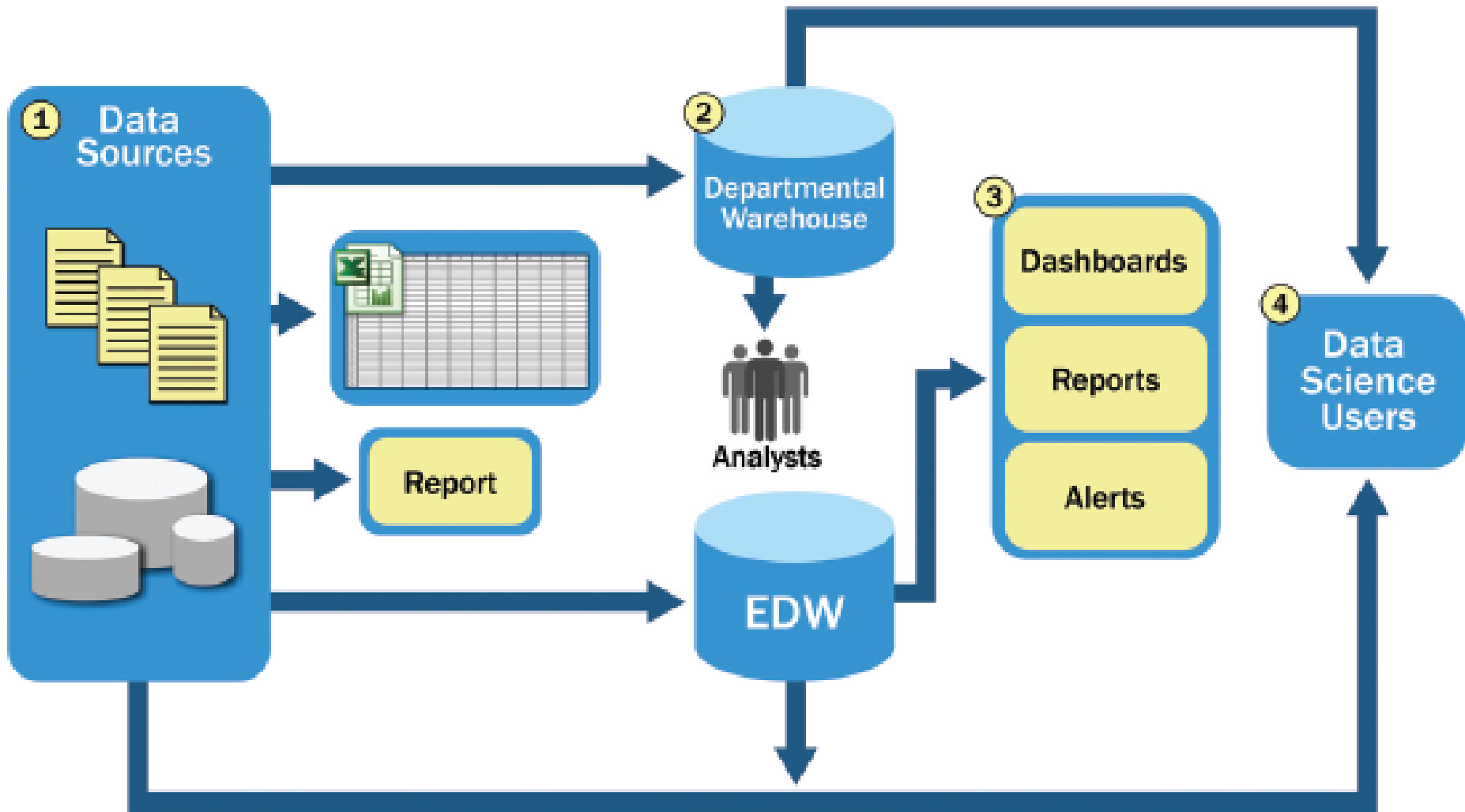


**Gene
Sequencing**

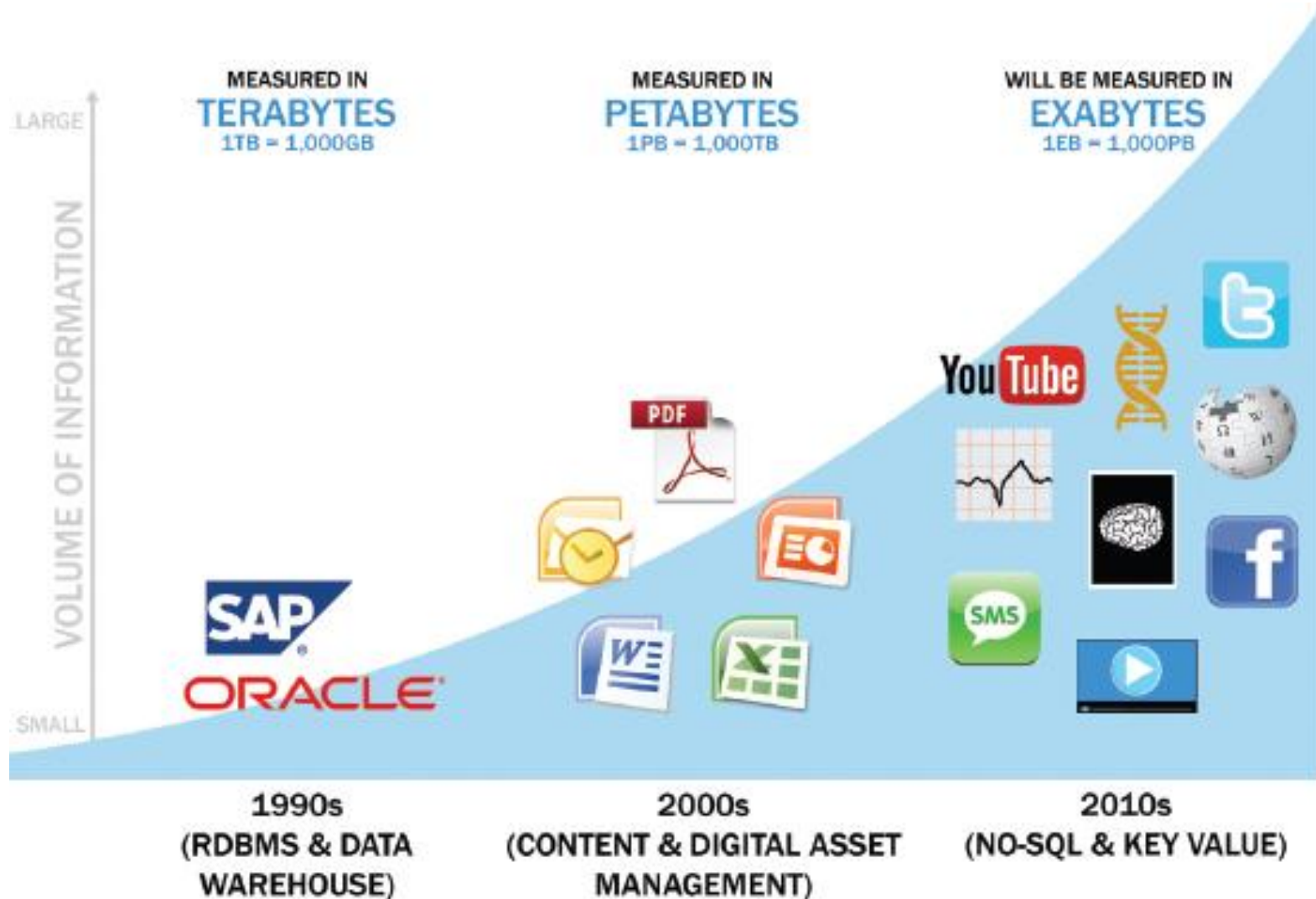
Big Data Growth is increasingly **unstructured**



Typical Analytic Architecture



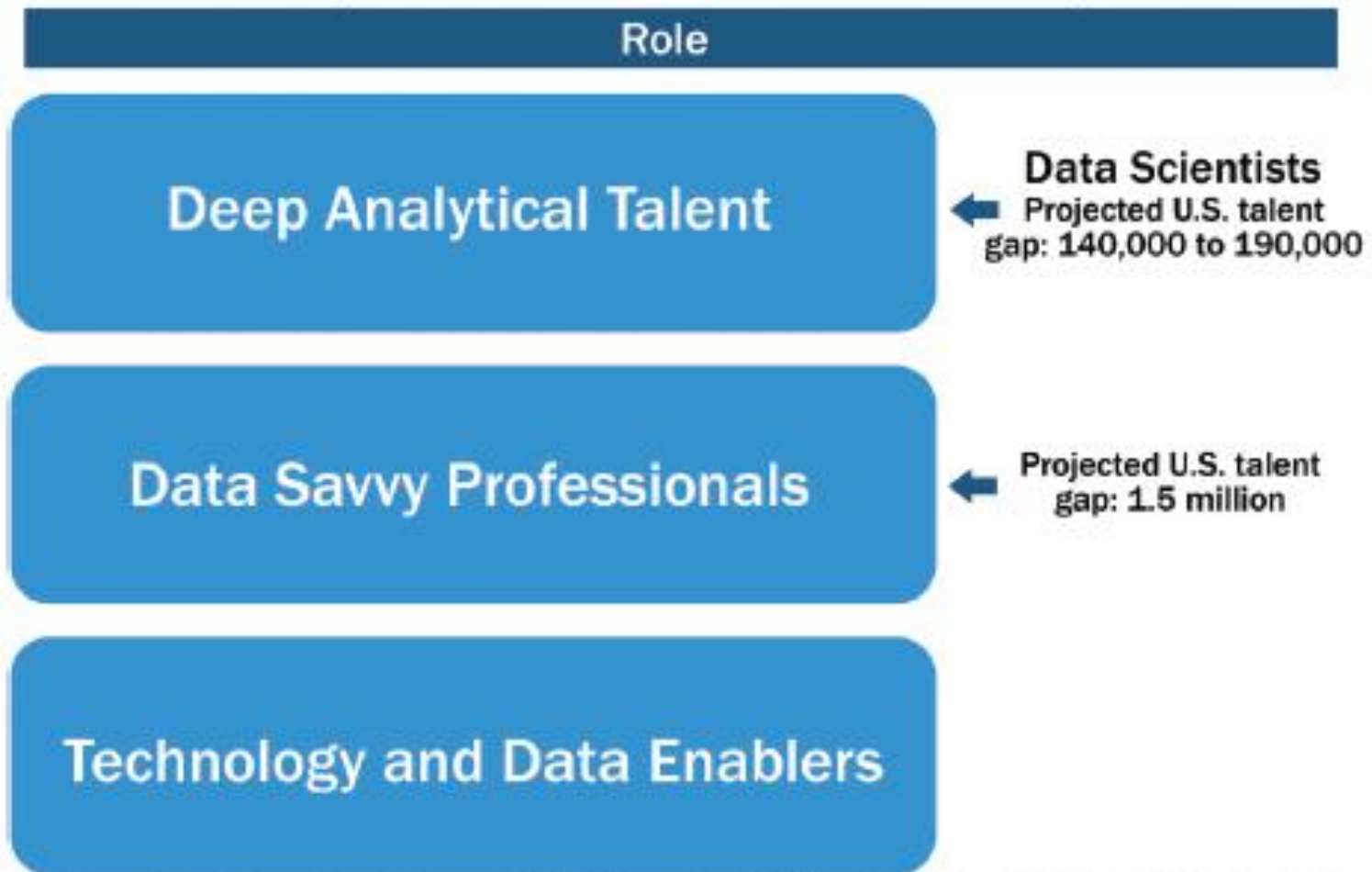
Data Evolution and the Rise of Big Data Sources



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

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
R²I²

DETECT

Anomalies
Communities
Typologies

PREDICT

Trending
Real-time
Prediction

DECIDE

Simulation
Optimization



BIG DATA – Batch



BIG DATA – Real Time



Complex by nature

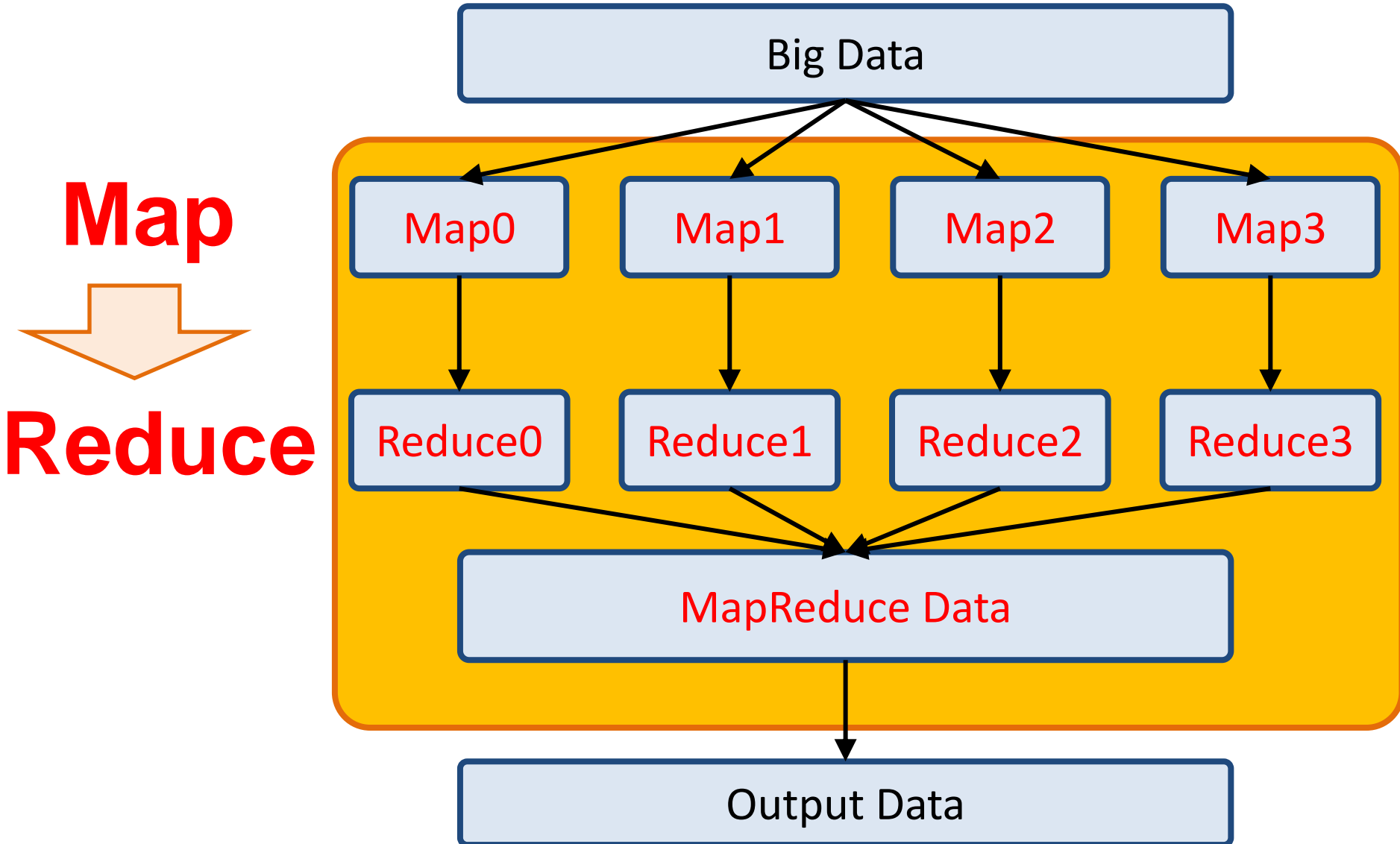


DATA

Complex by structure



MapReduce Paradigm





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



MapReduce

Processing

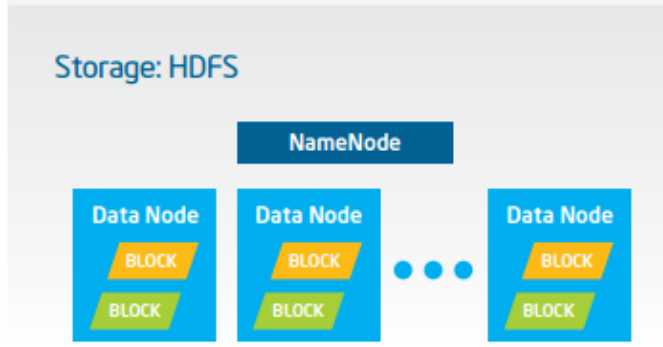
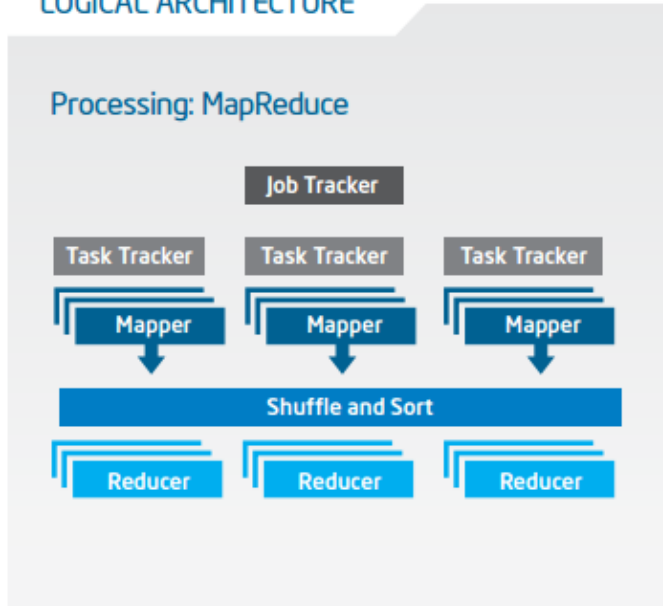


HDFS

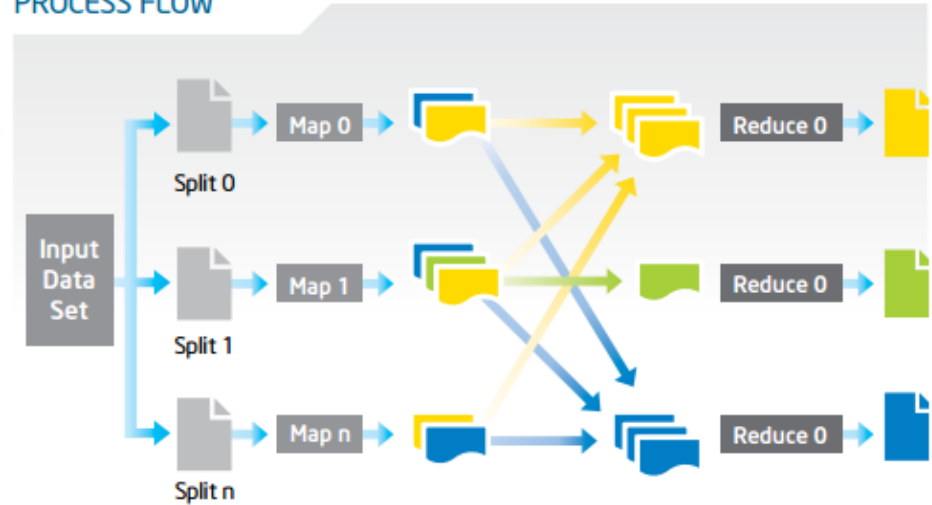
Storage

Big Data with Hadoop Architecture

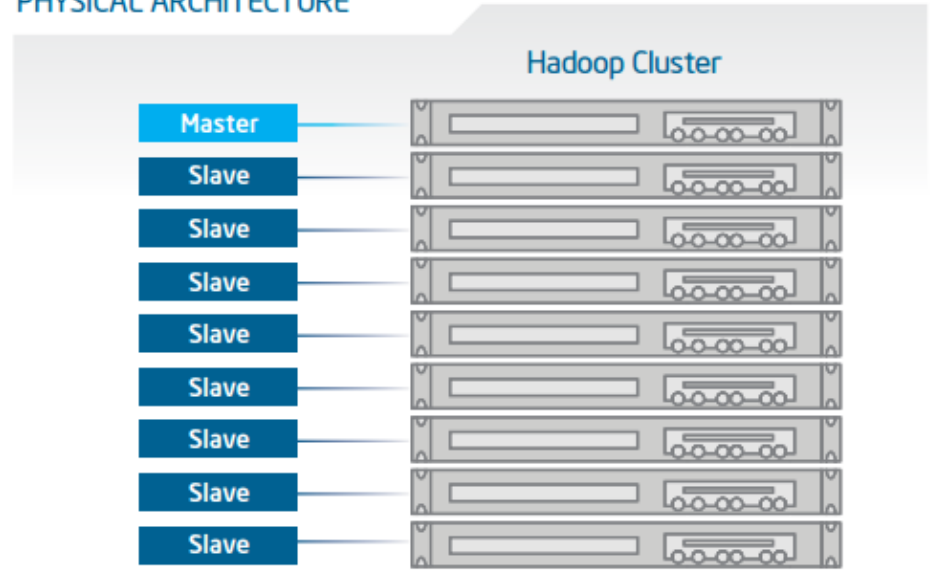
LOGICAL ARCHITECTURE



PROCESS FLOW



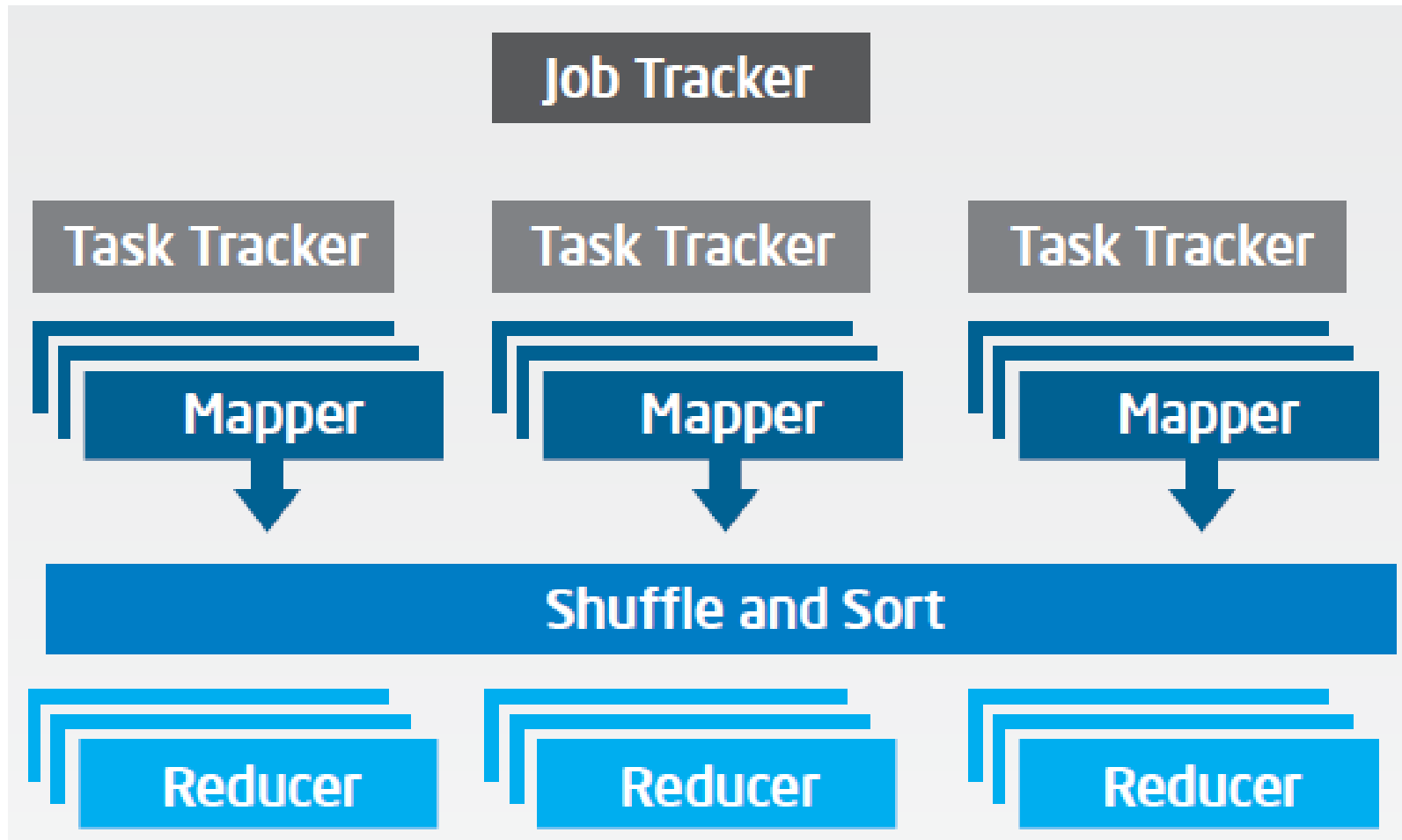
PHYSICAL ARCHITECTURE



Big Data with Hadoop Architecture

Logical Architecture

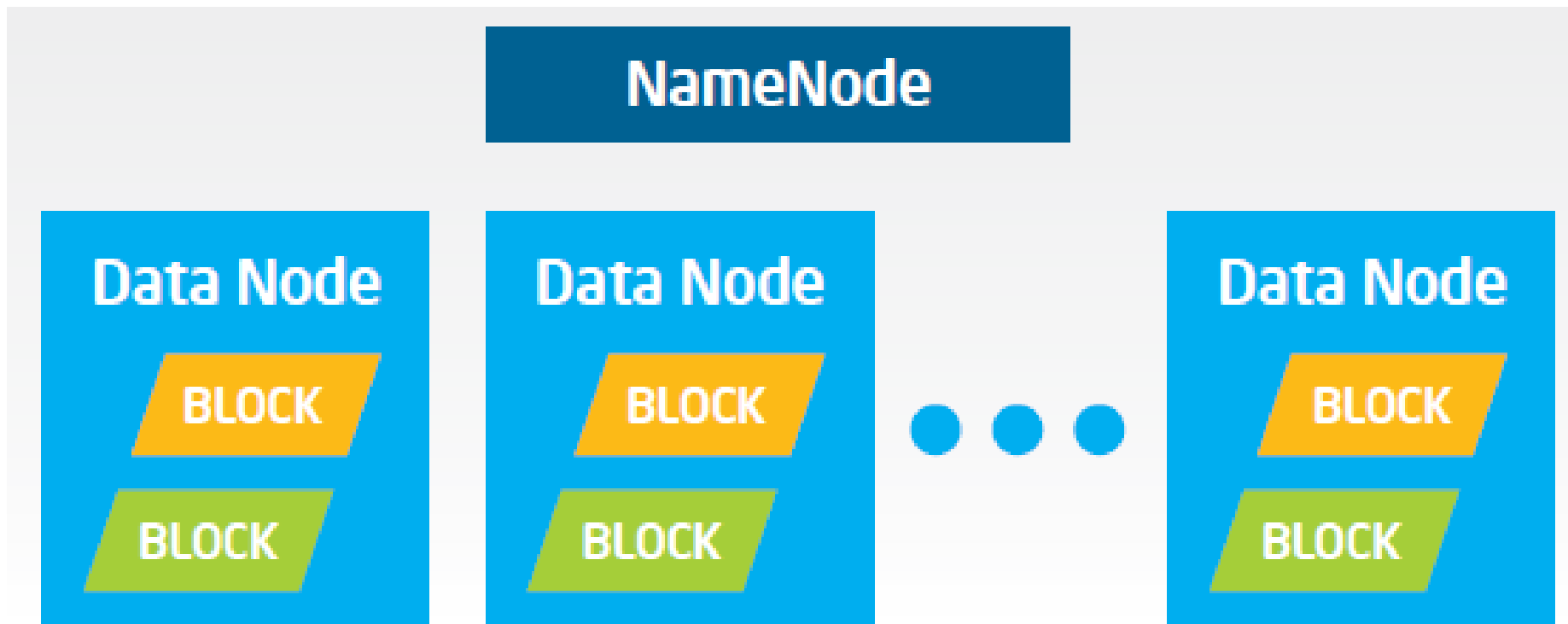
Processing: MapReduce



Big Data with Hadoop Architecture

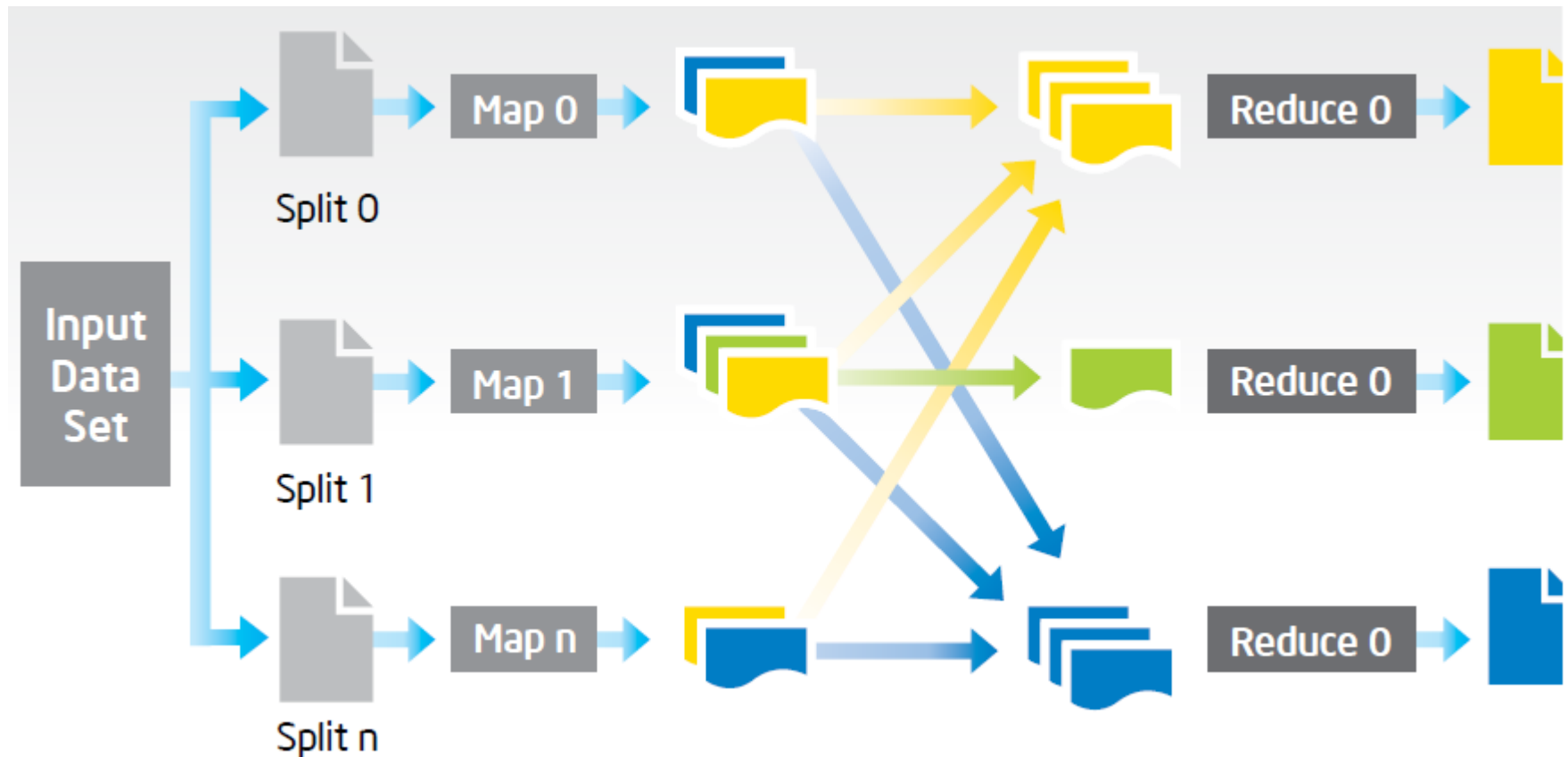
Logical Architecture

Storage: HDFS



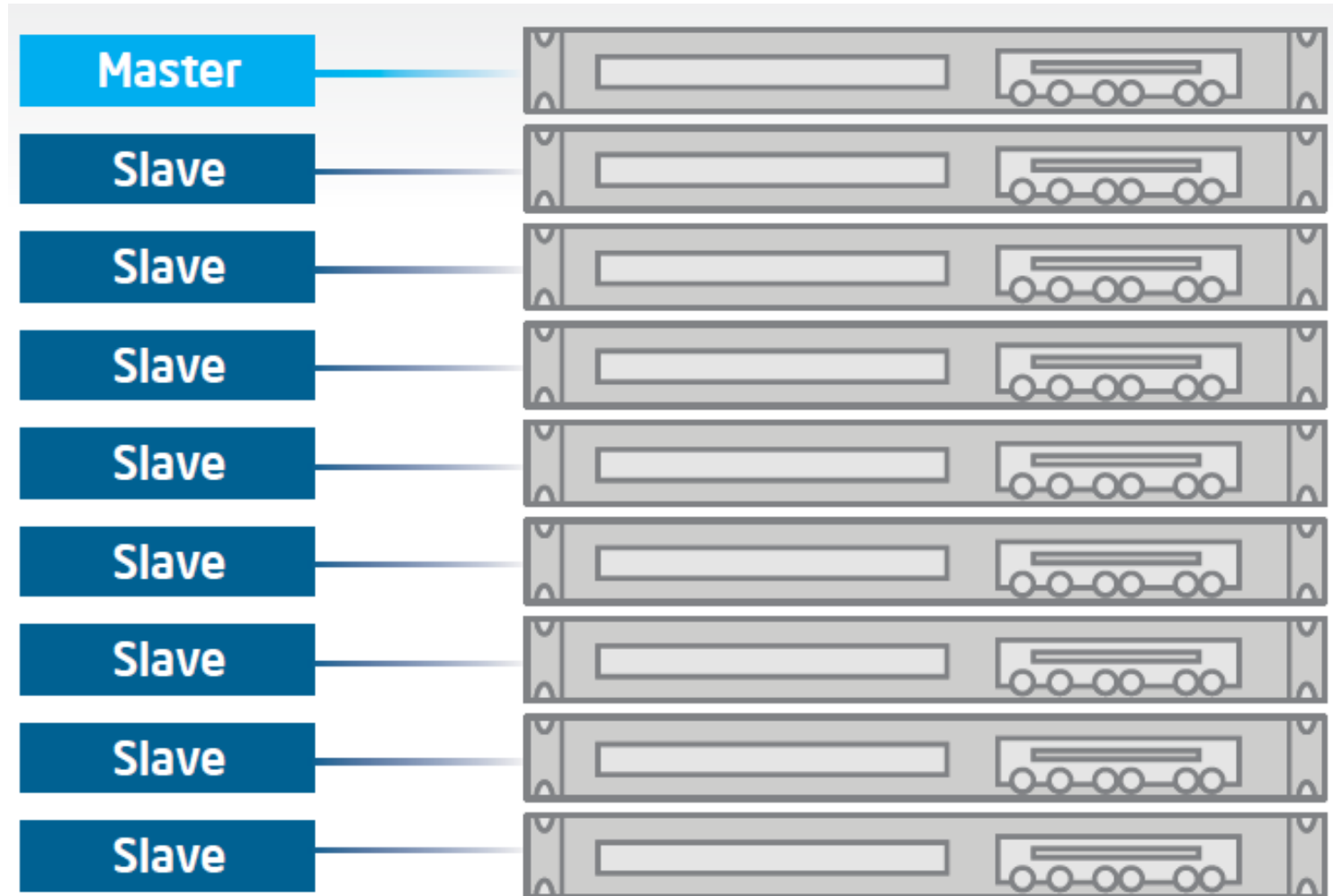
Big Data with Hadoop Architecture

Process Flow

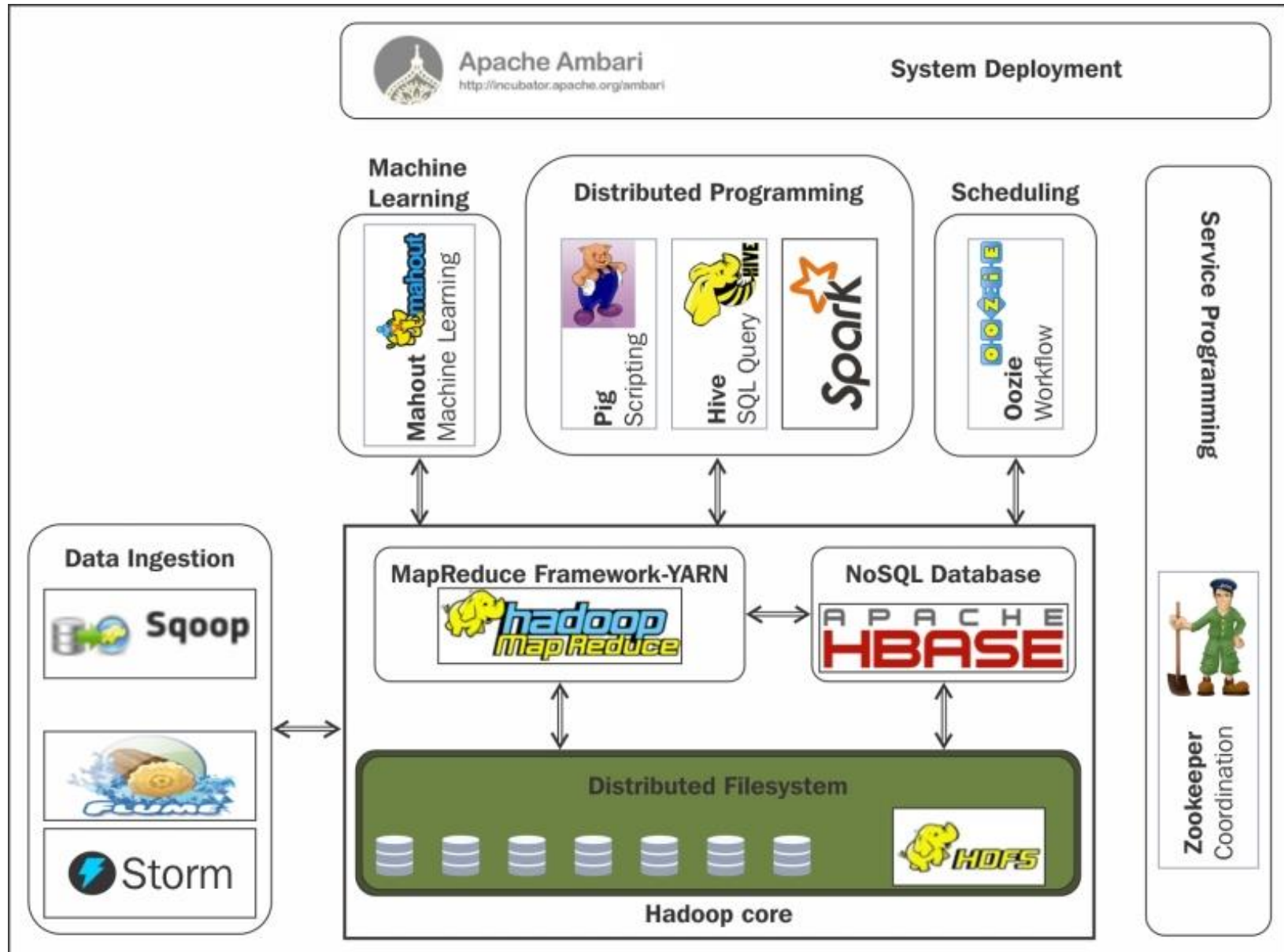


Big Data with Hadoop Architecture

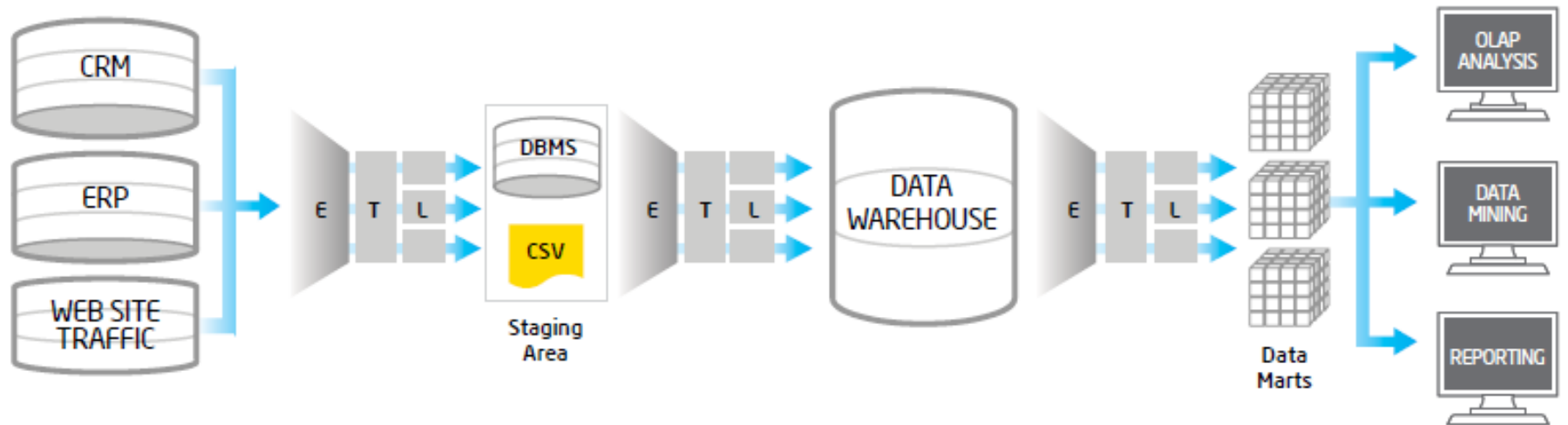
Hadoop Cluster



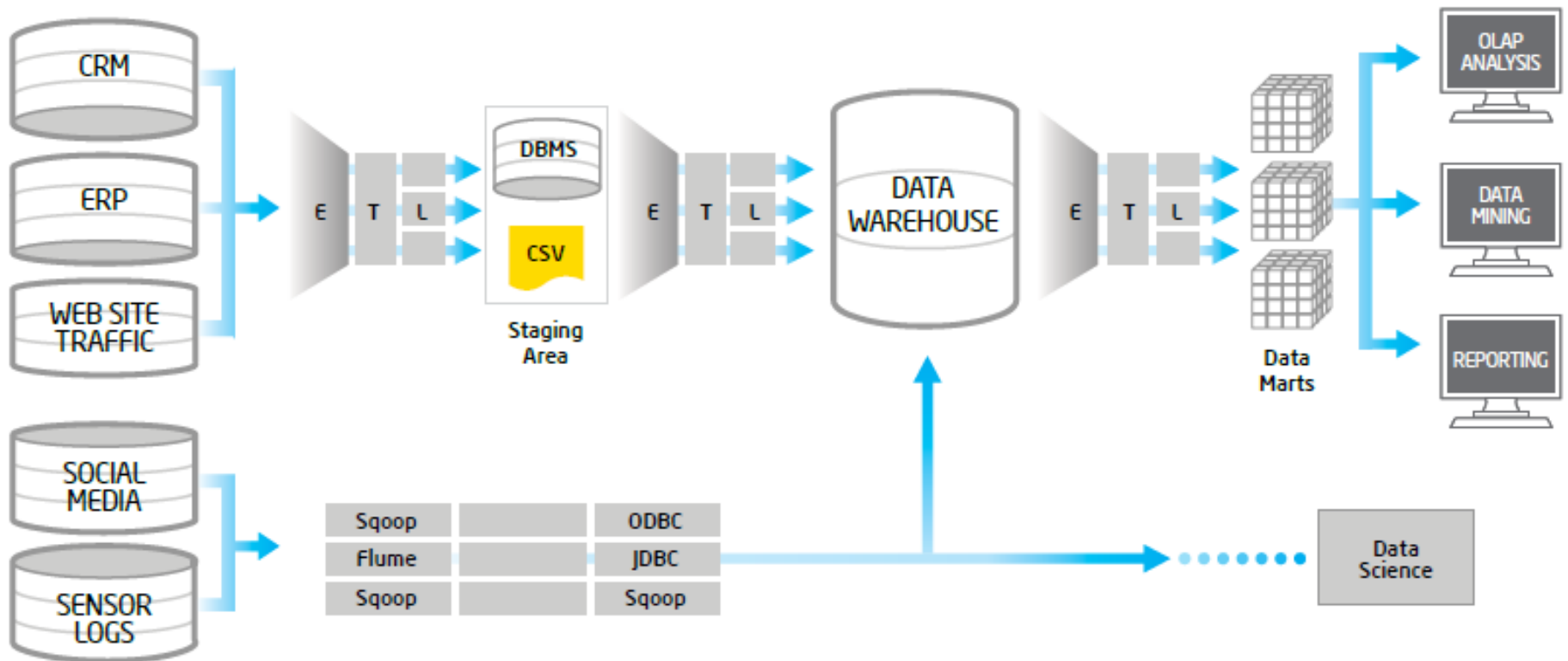
Hadoop Ecosystem



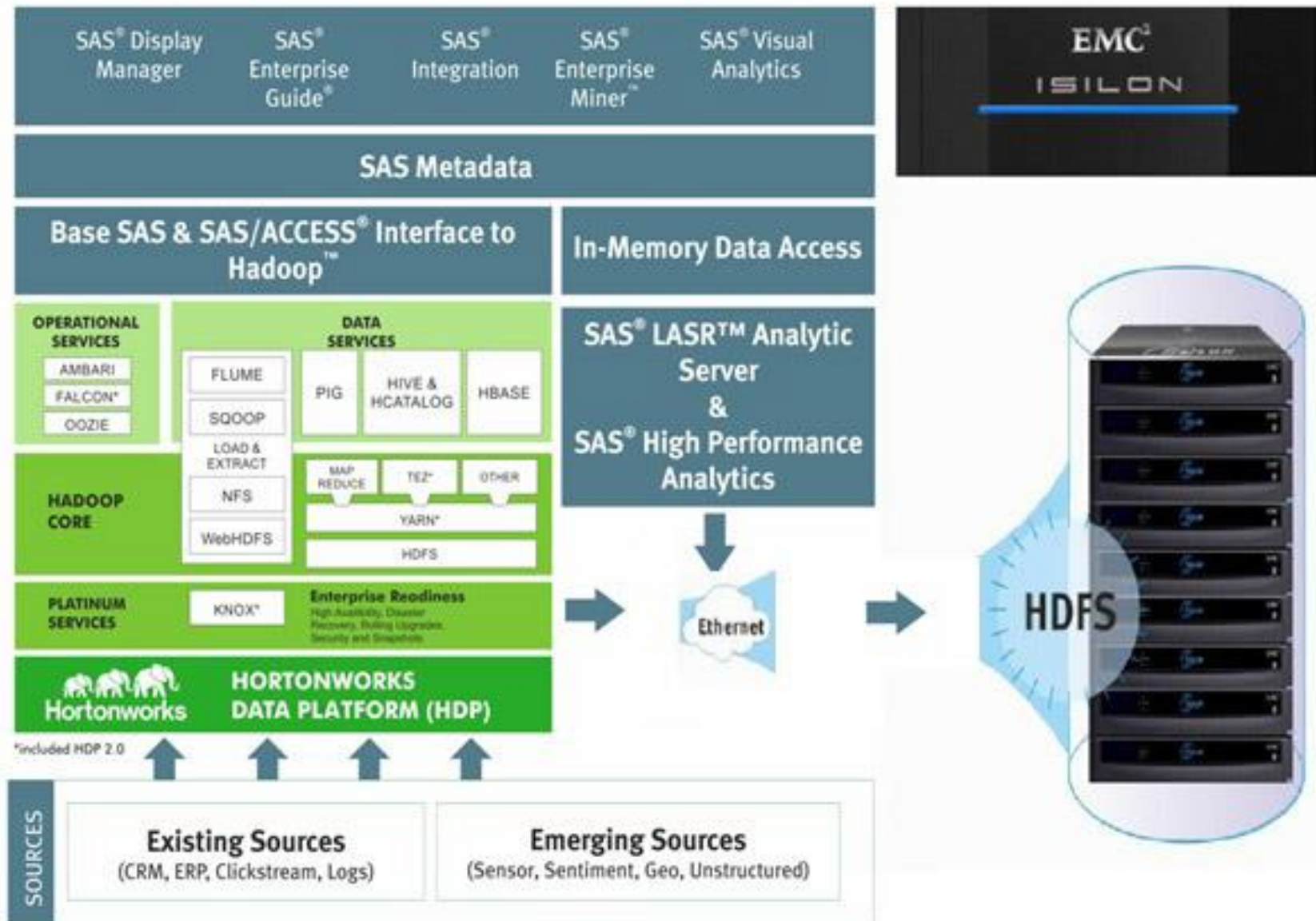
Traditional ETL Architecture



Offload ETL with Hadoop (Big Data Architecture)

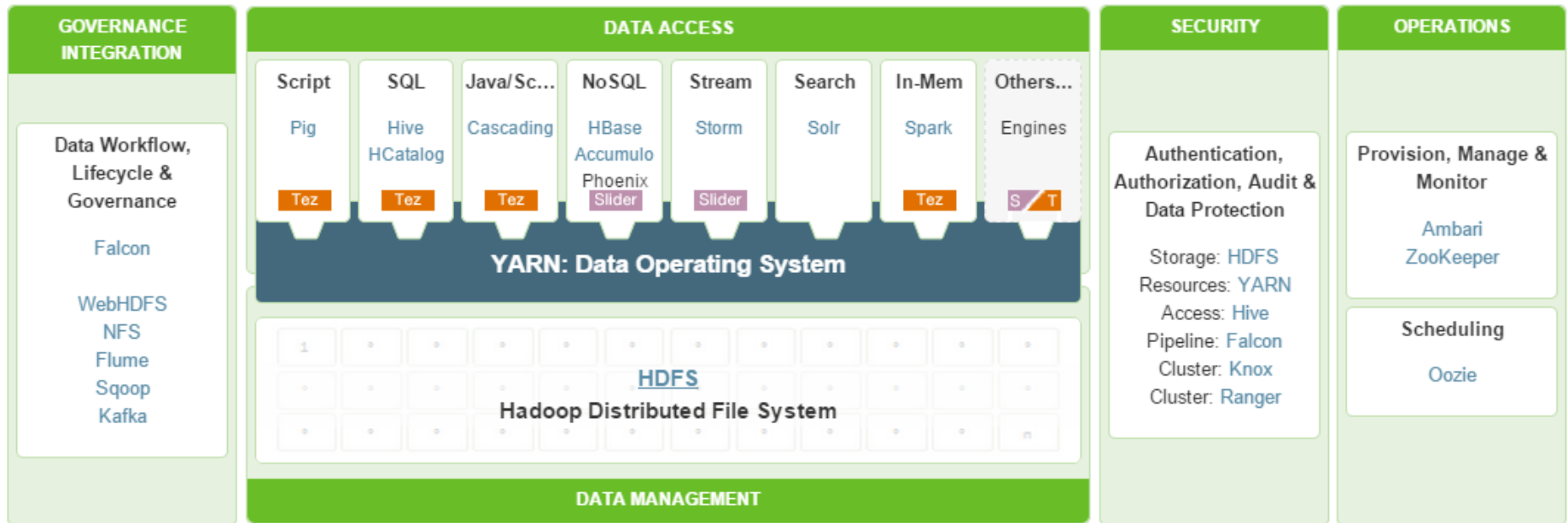


Big Data Solution



HDP

A Complete Enterprise Hadoop Data Platform



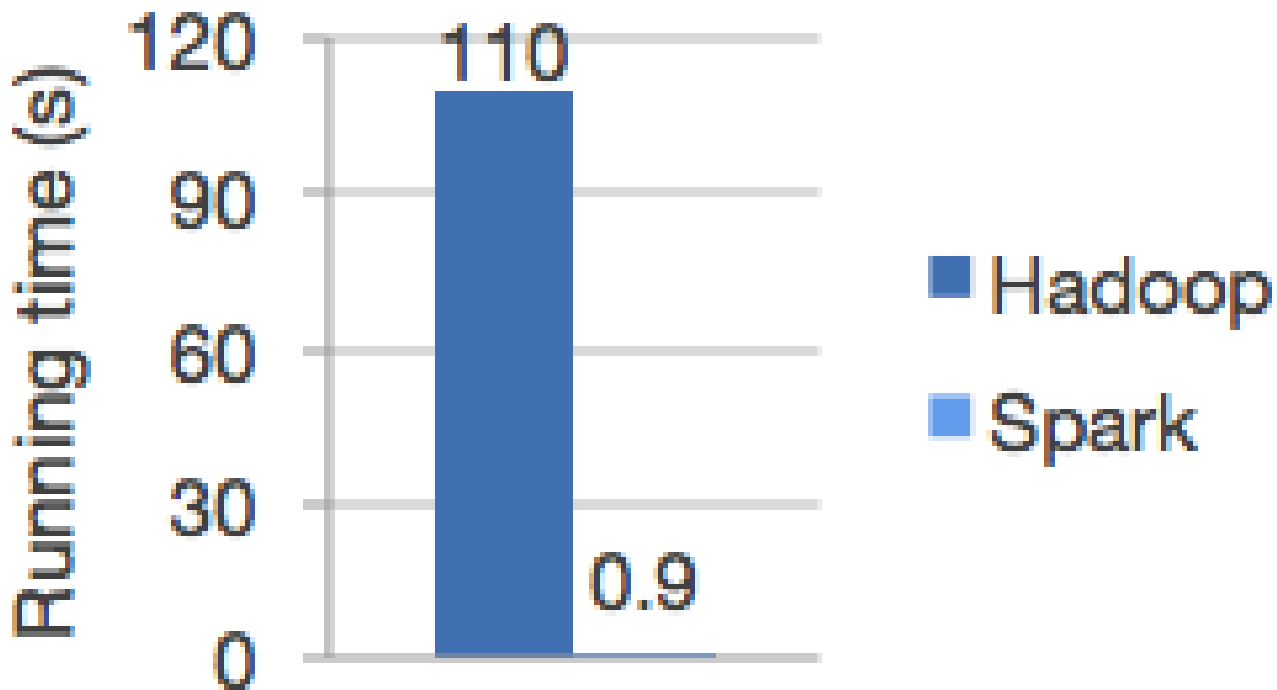


Lightning-fast cluster computing

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

```
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

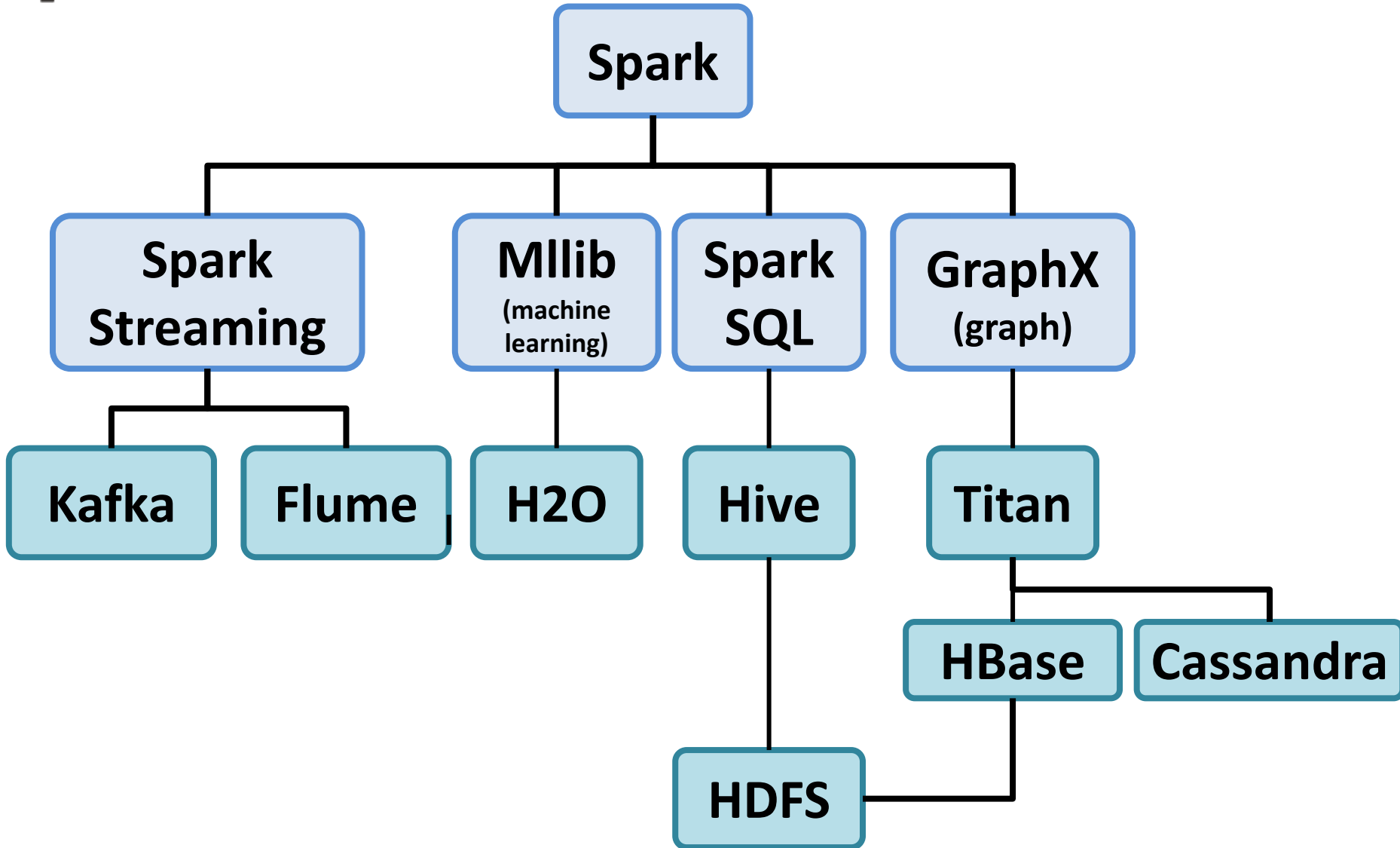
MLlib
(machine
learning)

GraphX
(graph)

Apache Spark













Spark Ecosystem

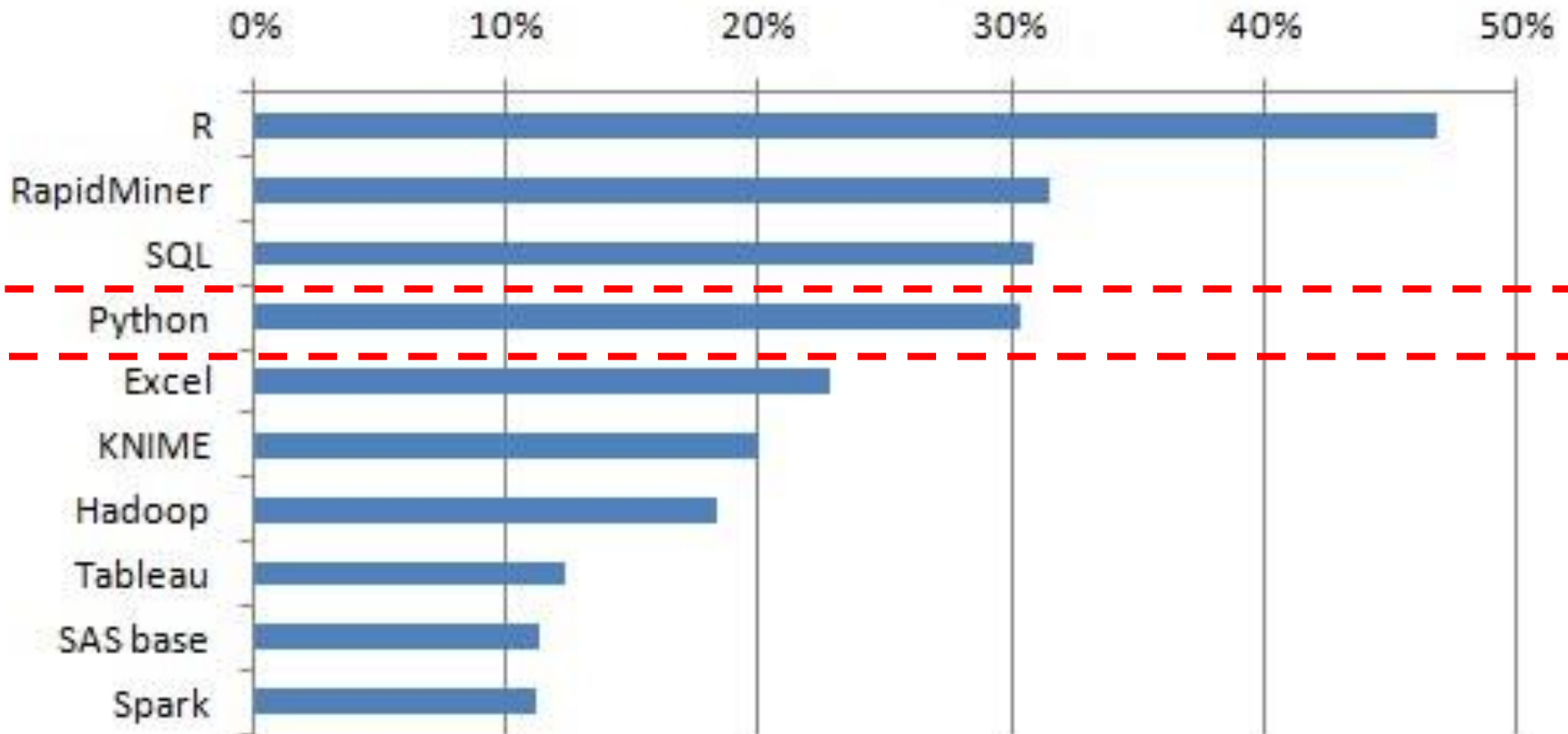


Python for Big Data Analytics

(The column on the left is the 2015 ranking; the column on the right is the 2014 ranking for comparison)

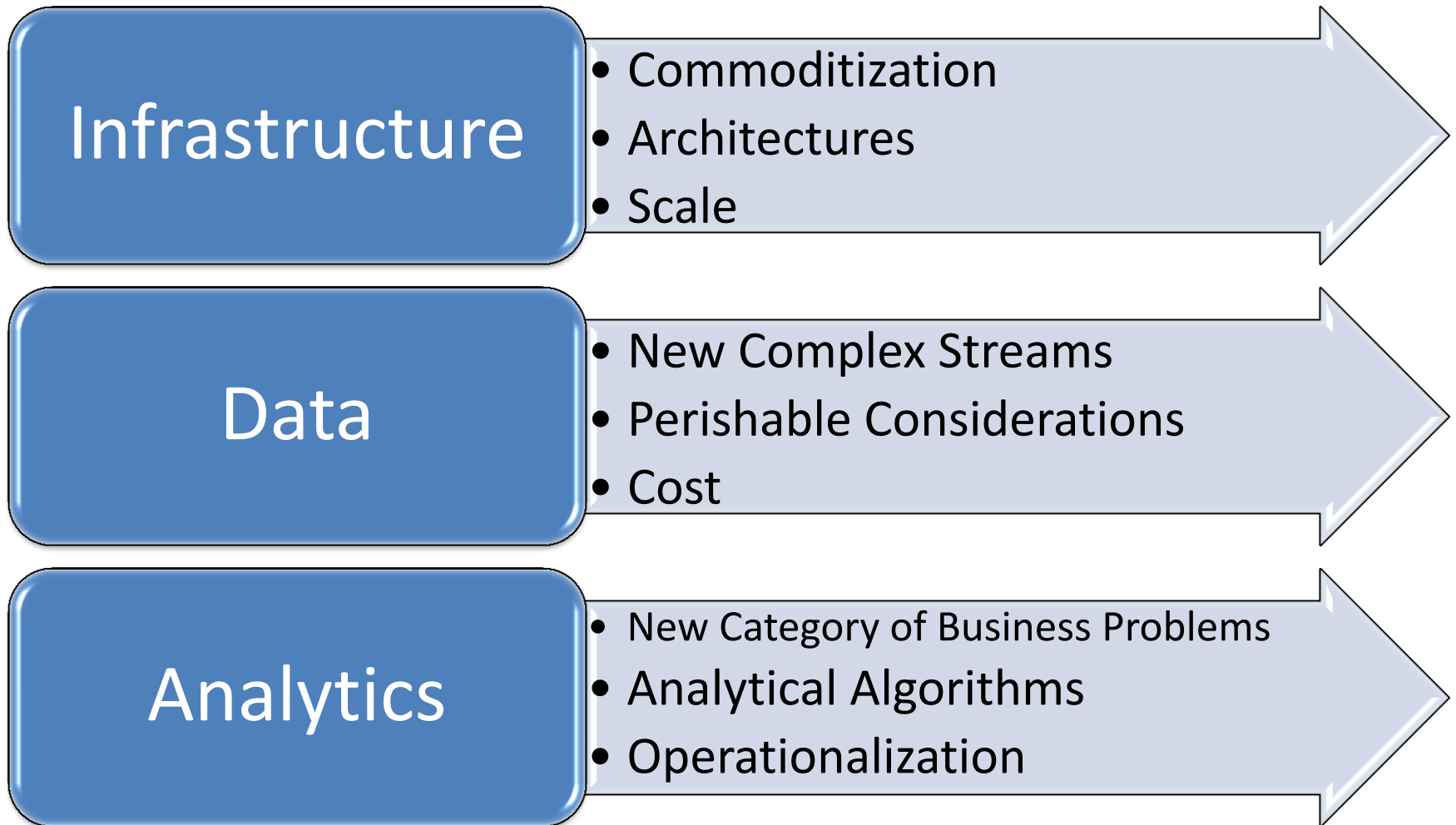
Language Rank	Types	2015 Spectrum Ranking	2014 Spectrum Ranking
1. Java		100.0	100.0
2. C		99.9	99.3
3. C++		99.4	95.5
4. Python		96.5	93.5
5. C#		91.3	92.4
6. R		84.8	84.8
7. PHP		84.5	84.5
8. JavaScript		83.0	78.9
9. Ruby		76.2	74.3
10. Matlab		72.4	72.8

Top Analytics, Data Mining, Data Science software used, 2015



SAS & Hadoop

Modern Reality



ADVANCED ANALYTICS

TEXT ANALYTICS

Finding treasures in unstructured data like social media or survey tools that could uncover insights about consumer sentiment

FORECASTING

Leveraging historical data to drive better insight into decision-making for the future



OPTIMIZATION

Analyze massive amounts of data in order to accurately identify areas likely to produce the most profitable results

DATA MINING

Mine transaction databases for data of spending patterns that indicate a stolen card..

STATISTICS

Trends in Analytics



Complex Business Problems Are Driving Analytics Innovation



Speed Will Be Of Essence



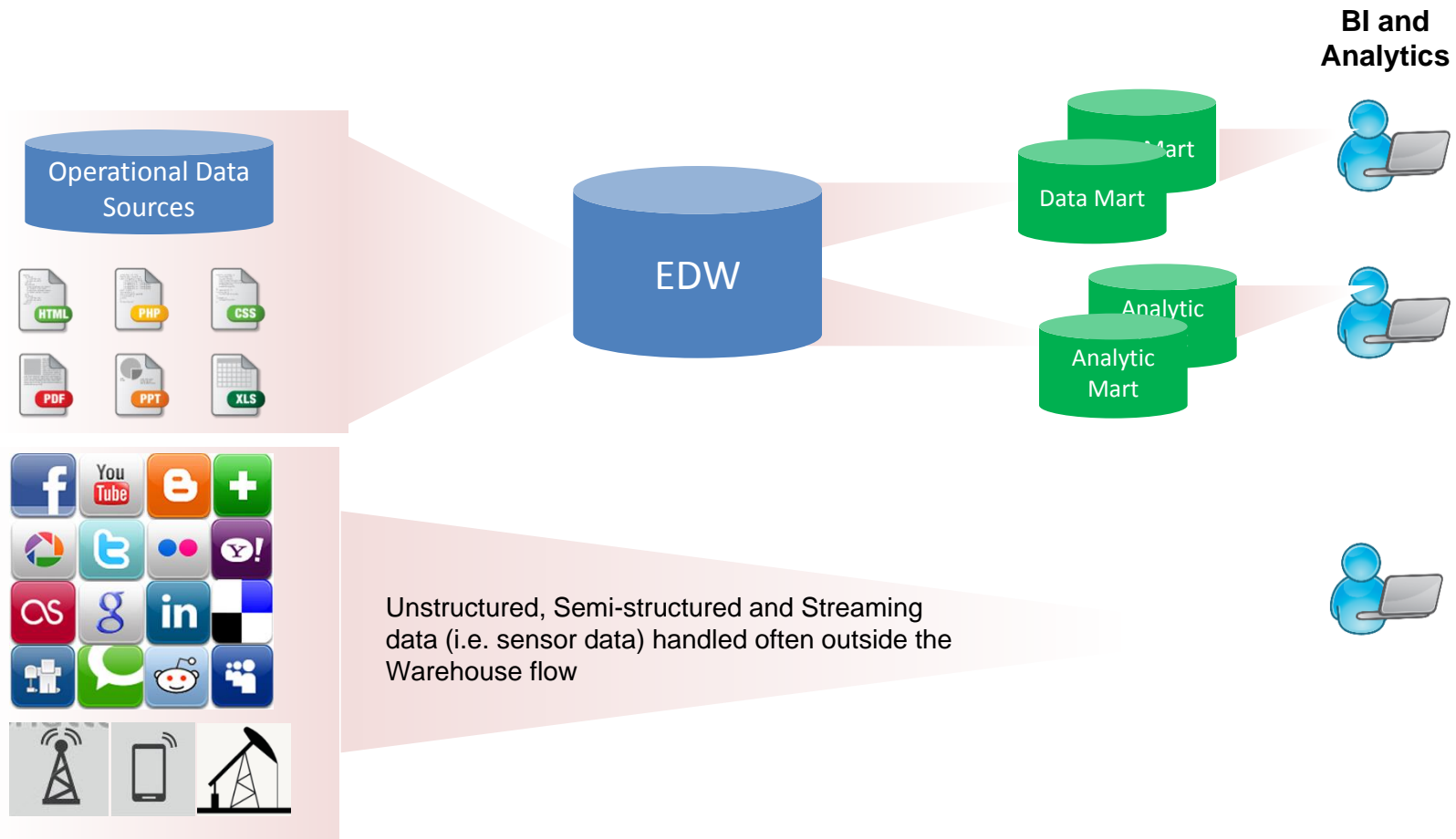
Leverage Analytics To Unlock The Information Contained In Unstructured Data



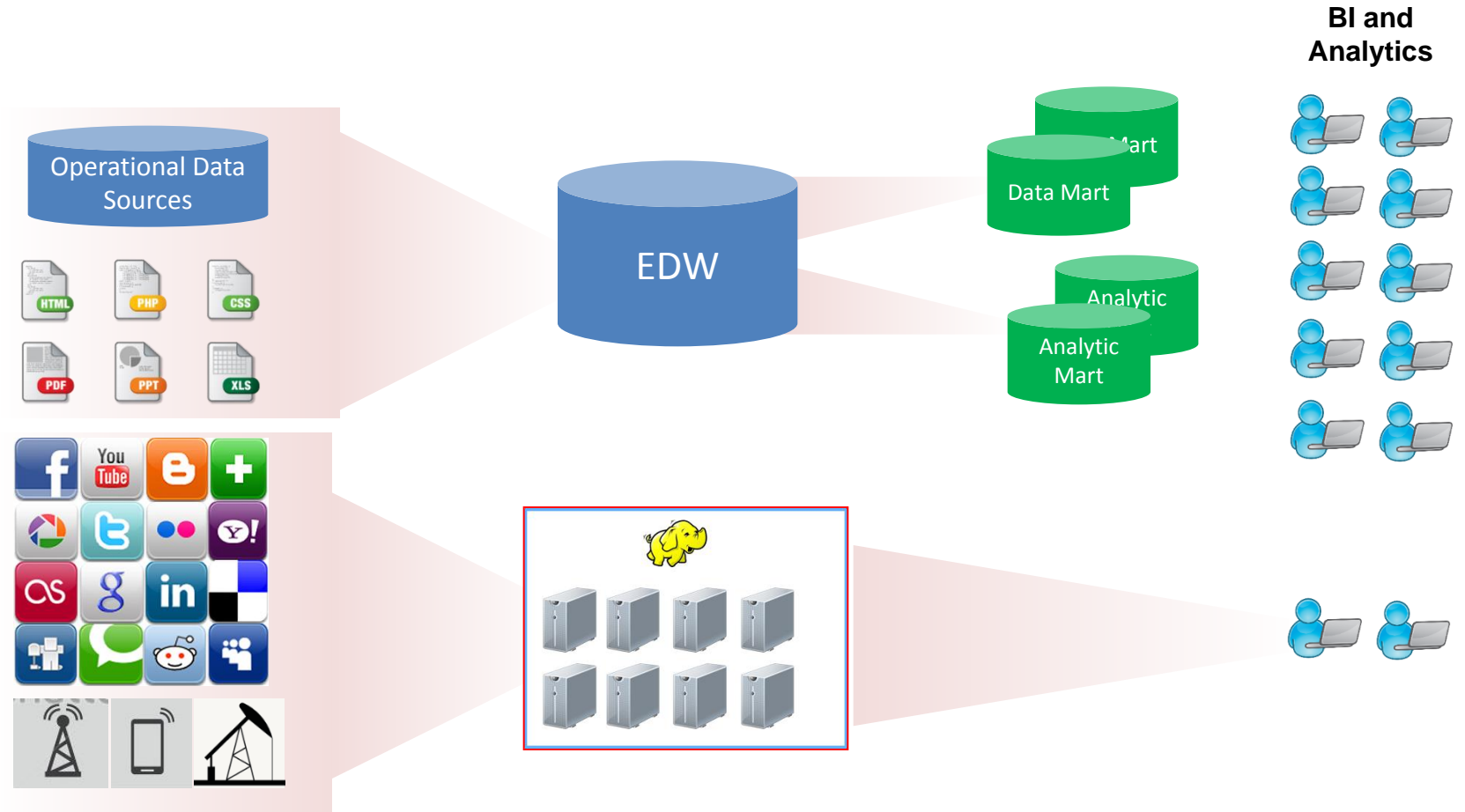
Operationalizing Analytics

Architectures of Big Data Analytics

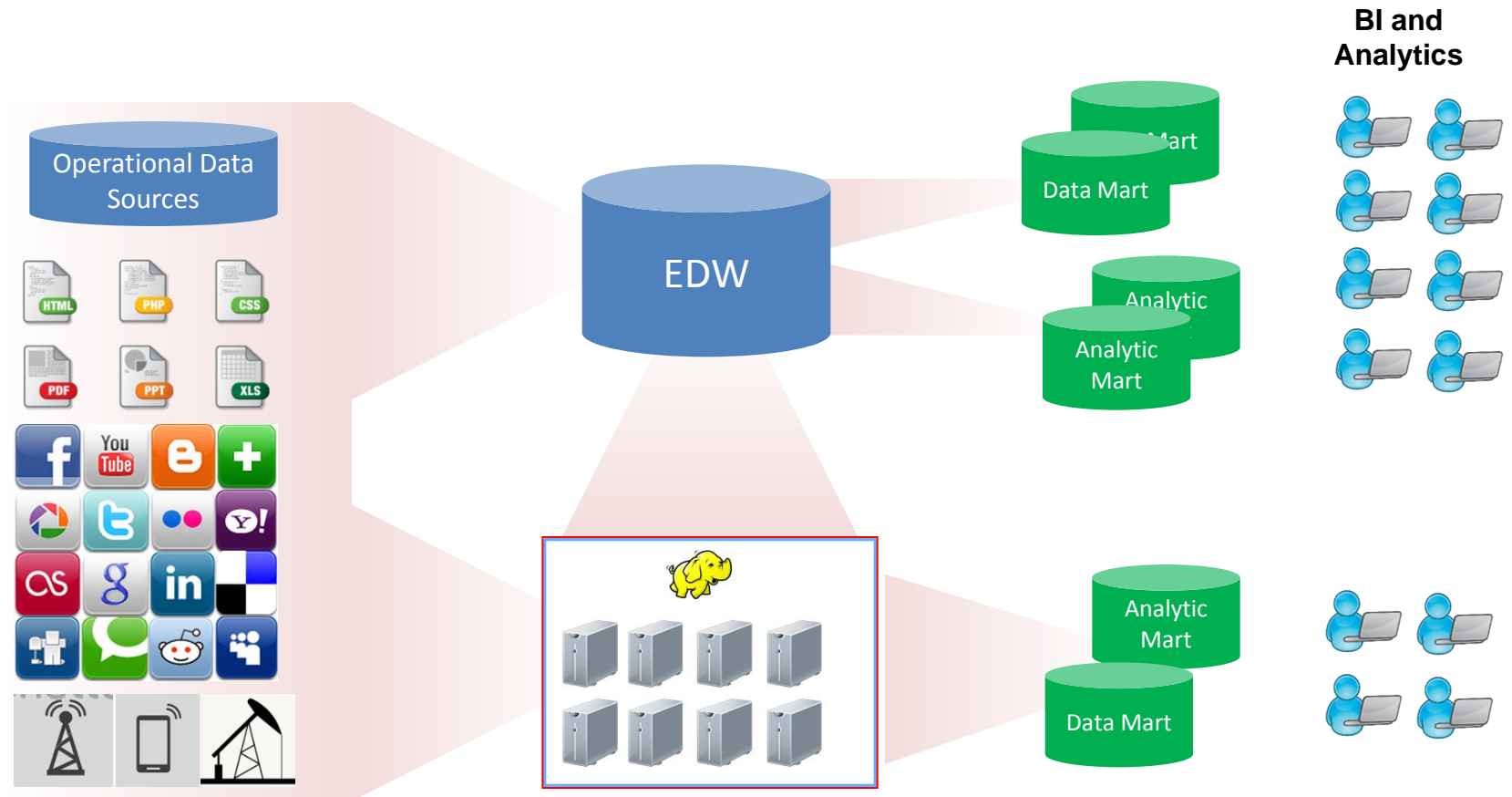
Traditional Analytics



Hadoop as a “new data” Store

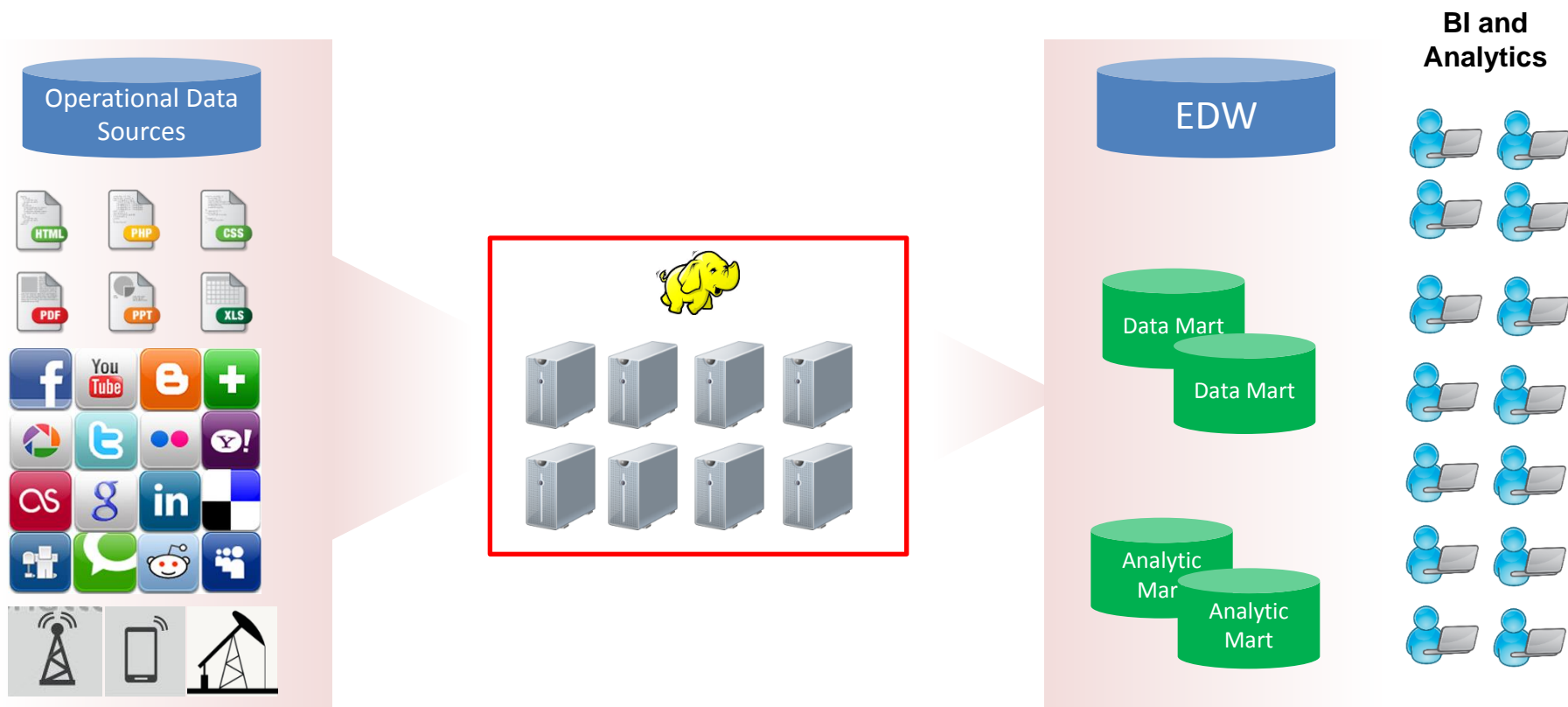


Hadoop as an additional input to the EDW



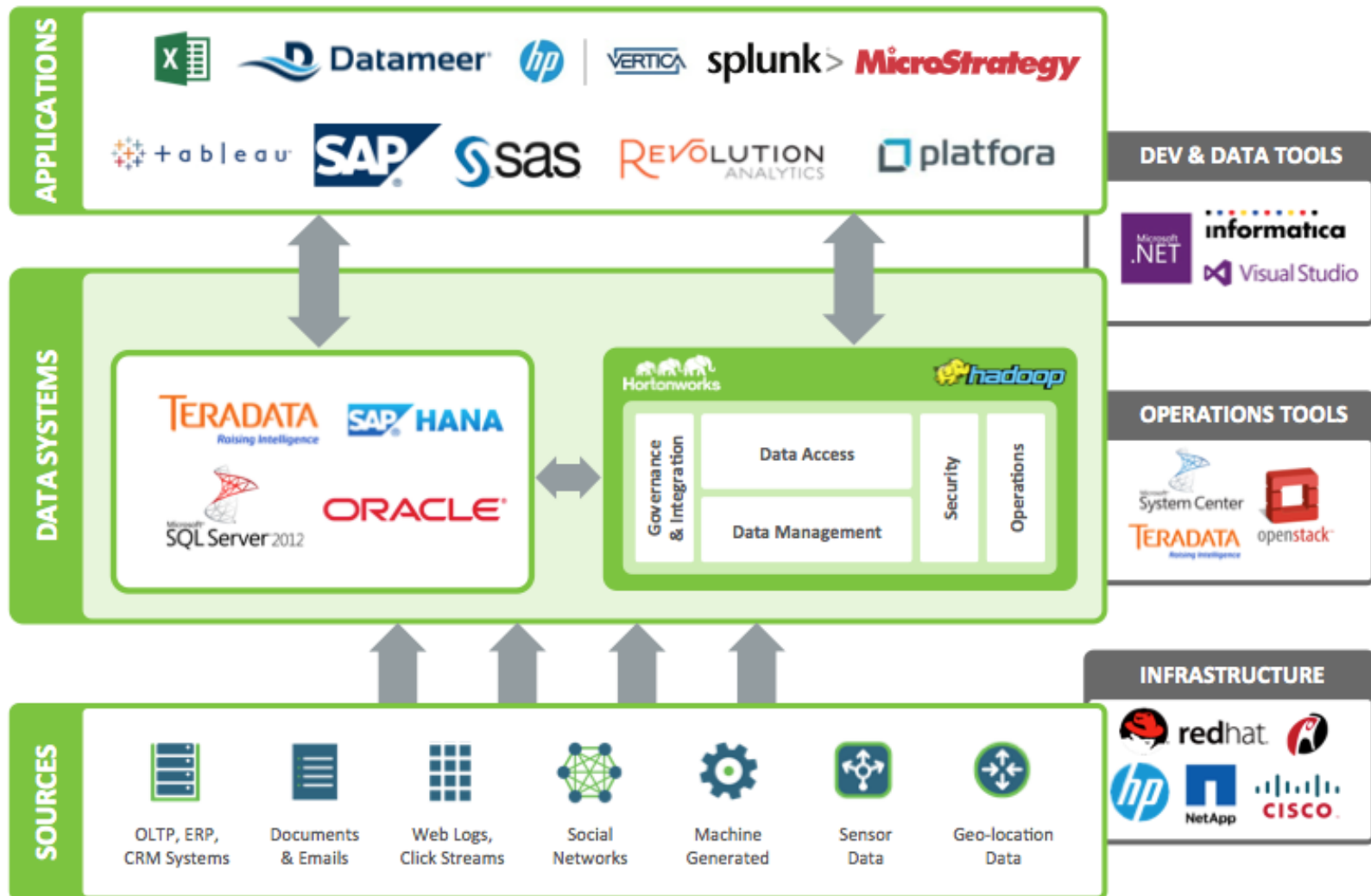
Hadoop Data Platform As a “staging Layer” as part of a “data Lake”

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



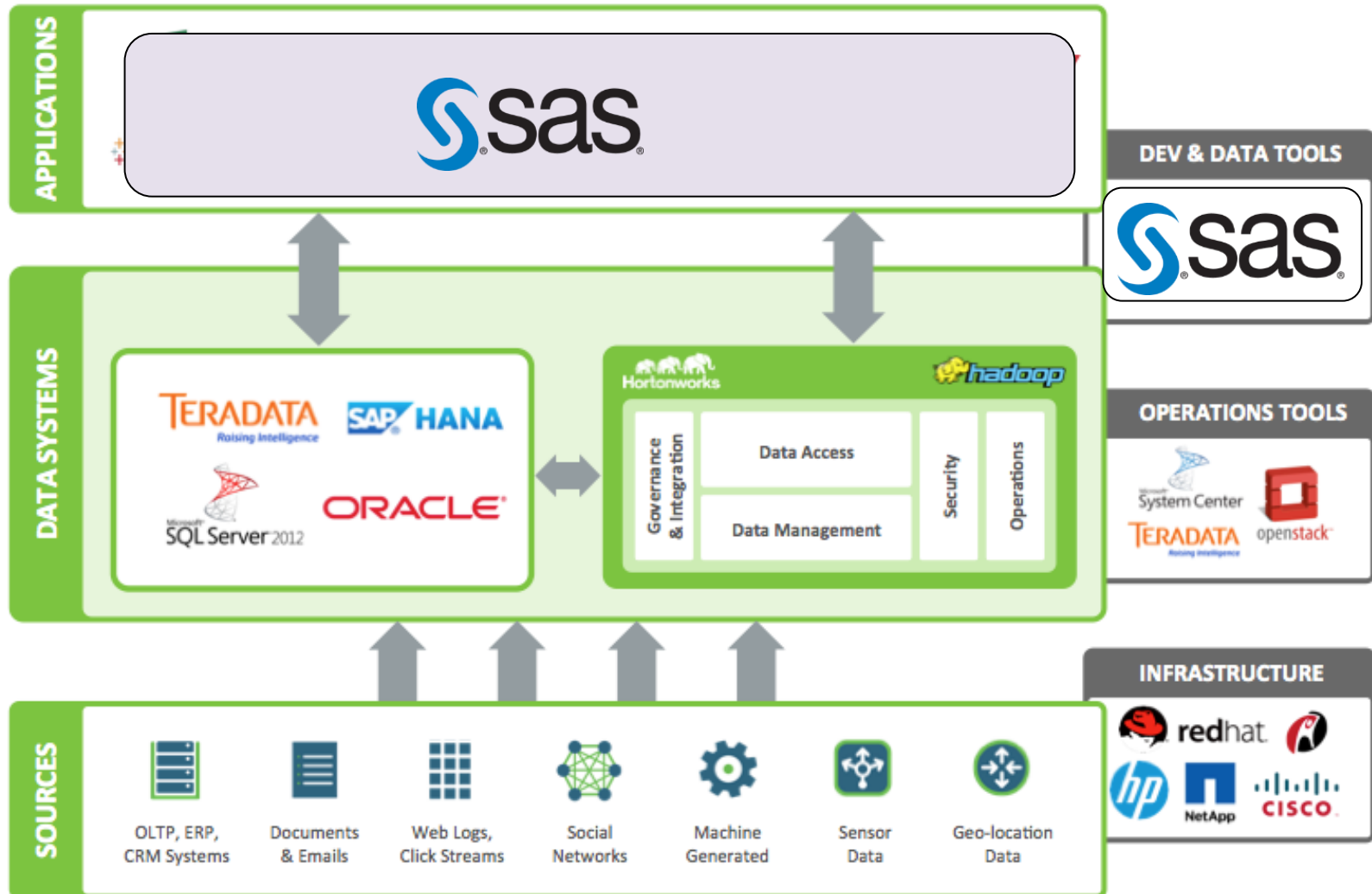
SAS Big data Strategy

– SAS areas

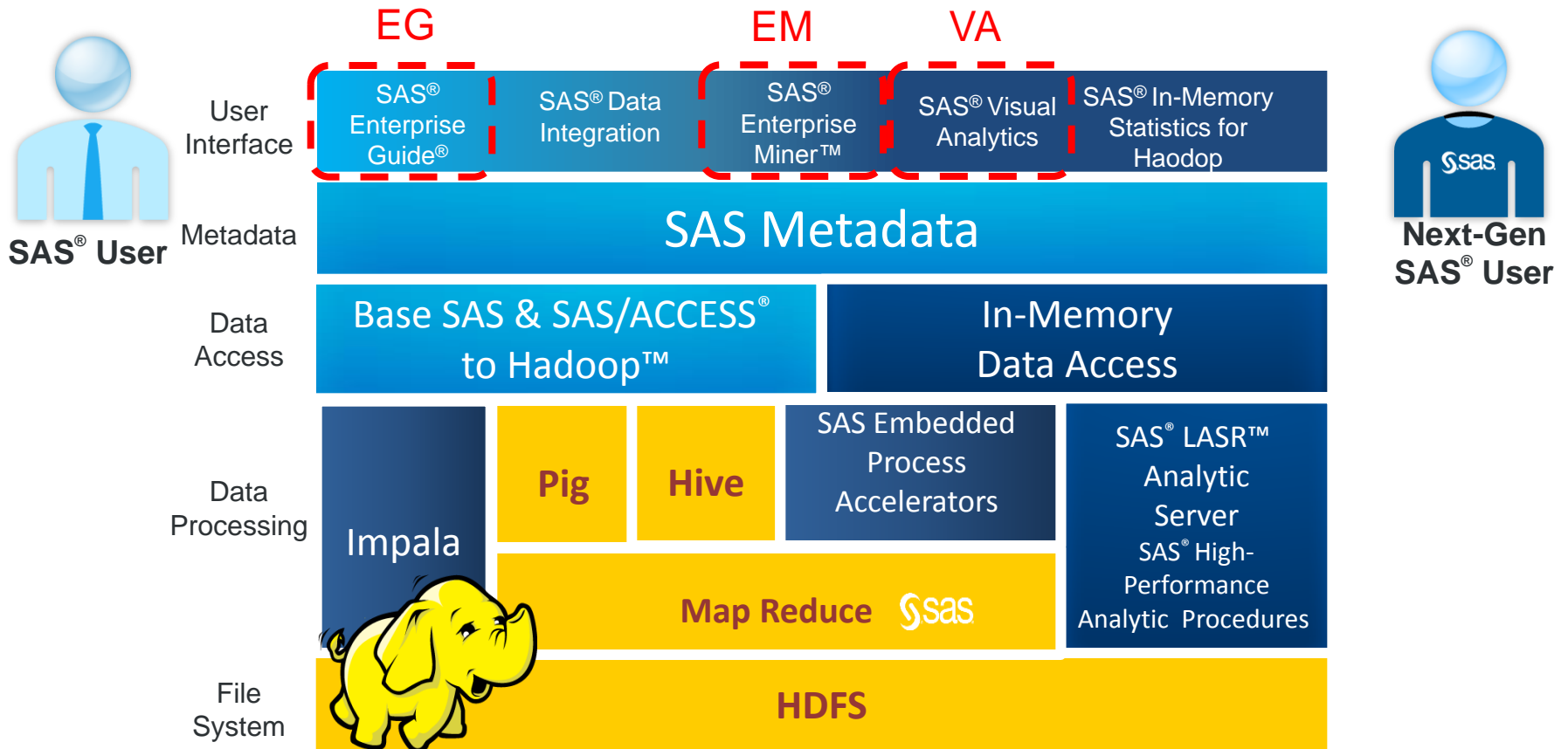


SAS Big data Strategy

– SAS areas

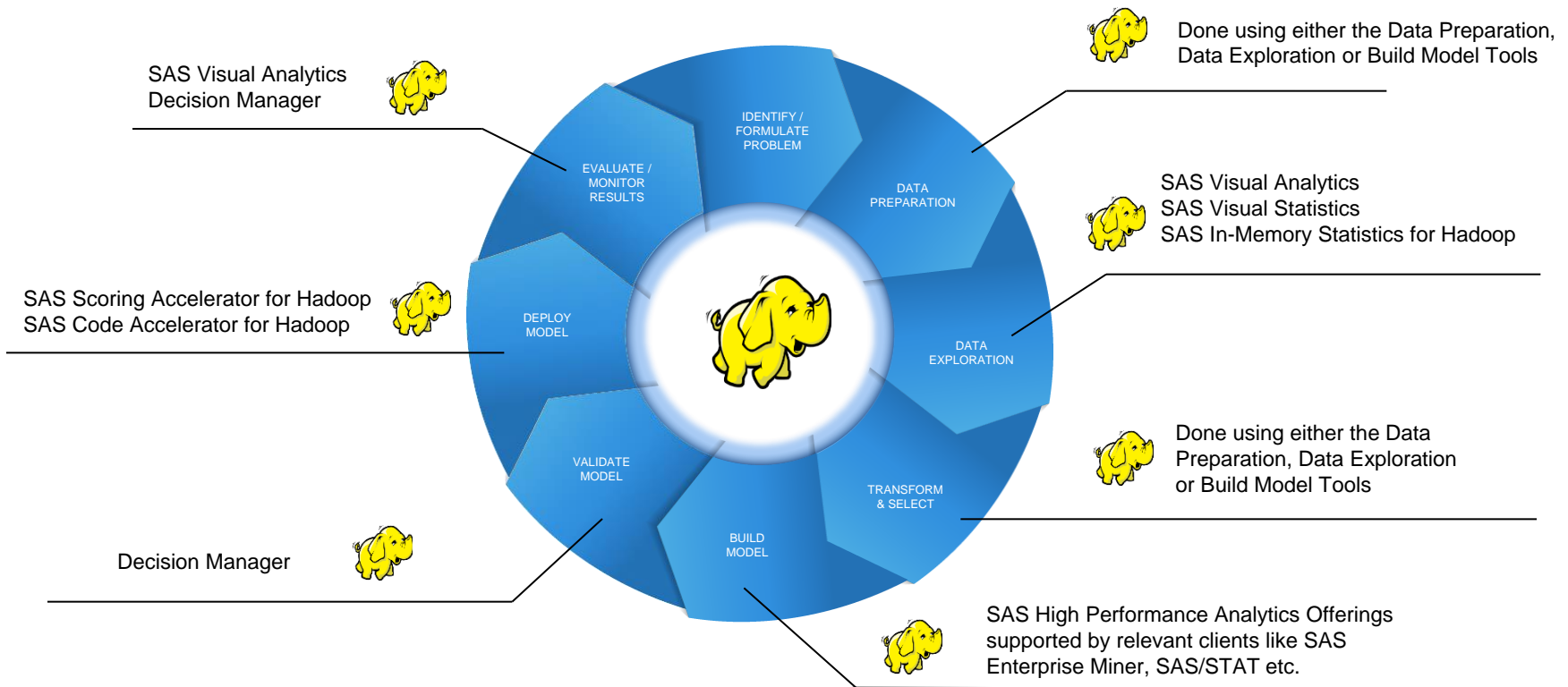


SAS® Within the HADOOP ECOSYSTEM



SAS enables the entire lifecycle around HADOOP

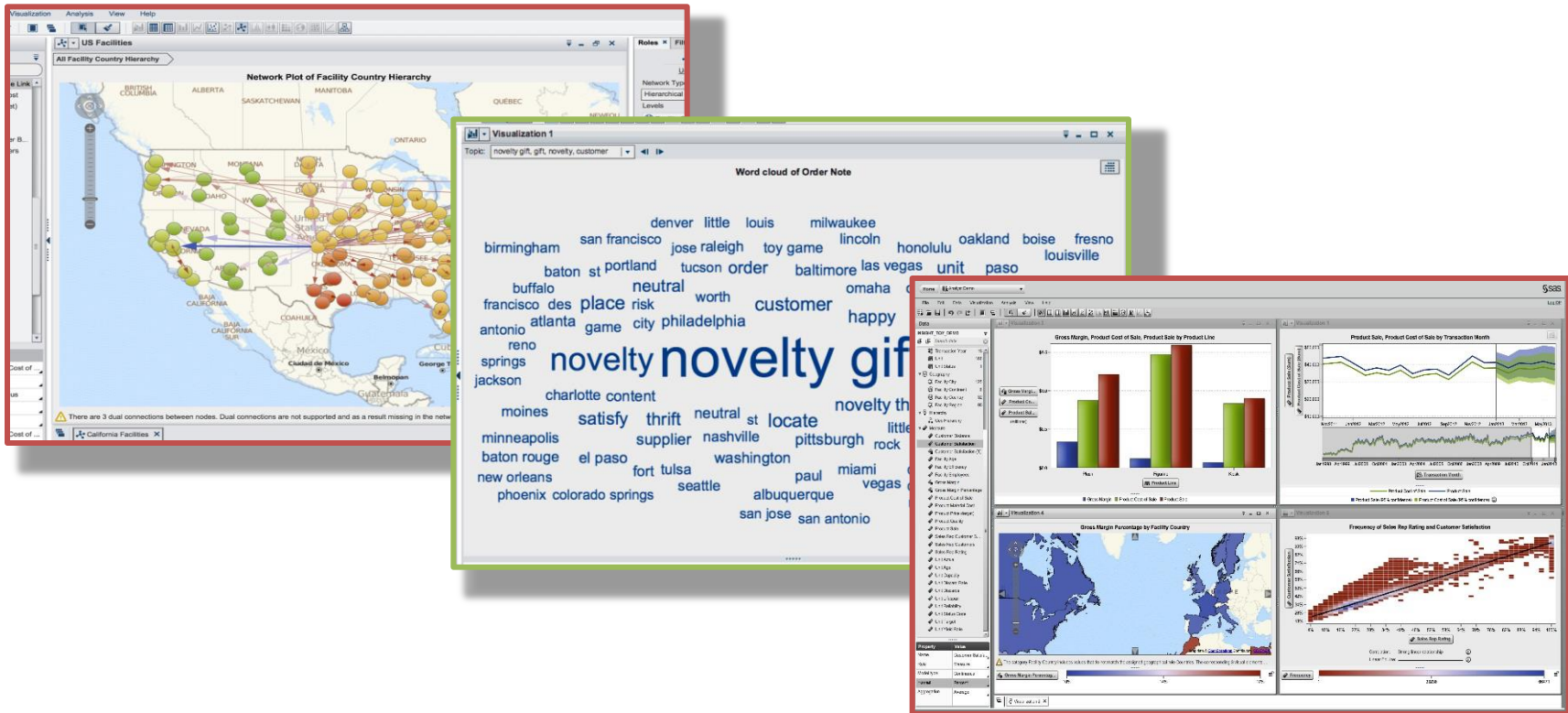
SAS enableS the entire lifecycle around HADOOP



SAS[®] VISUAL ANALYTICS

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

Visualization



References

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