

เอกสารประกอบการทำกิจกรรม COP

Basic



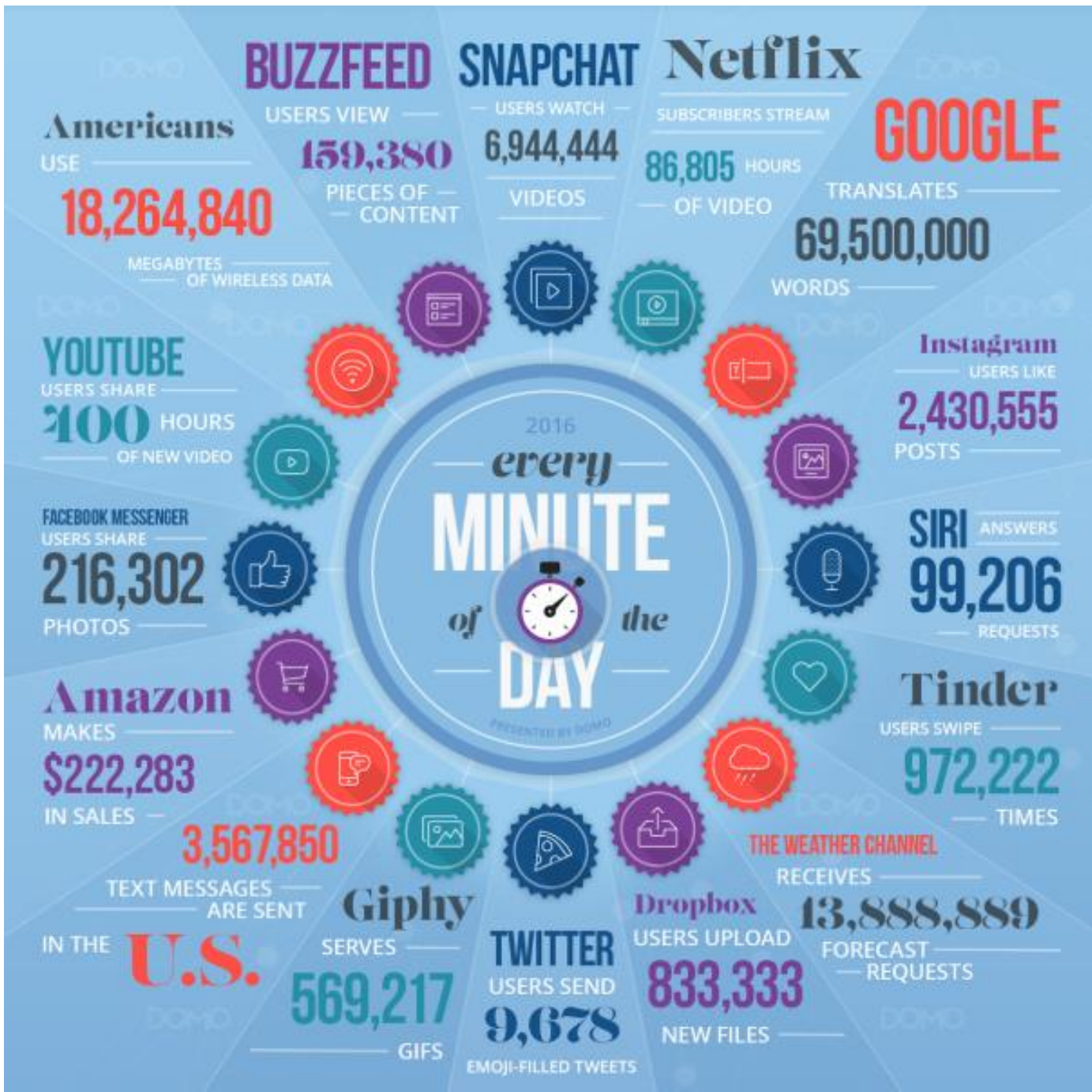
ศส.กรอ.



Narongrich Chanprasert

BI Consultant

Depth First Co.,Ltd.



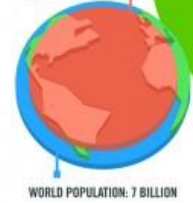
The 4V of Big Data

40 ZETTABYTES

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



6 BILLION PEOPLE have cell phones



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

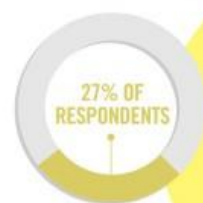
Velocity
ANALYSIS OF STREAMING DATA

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



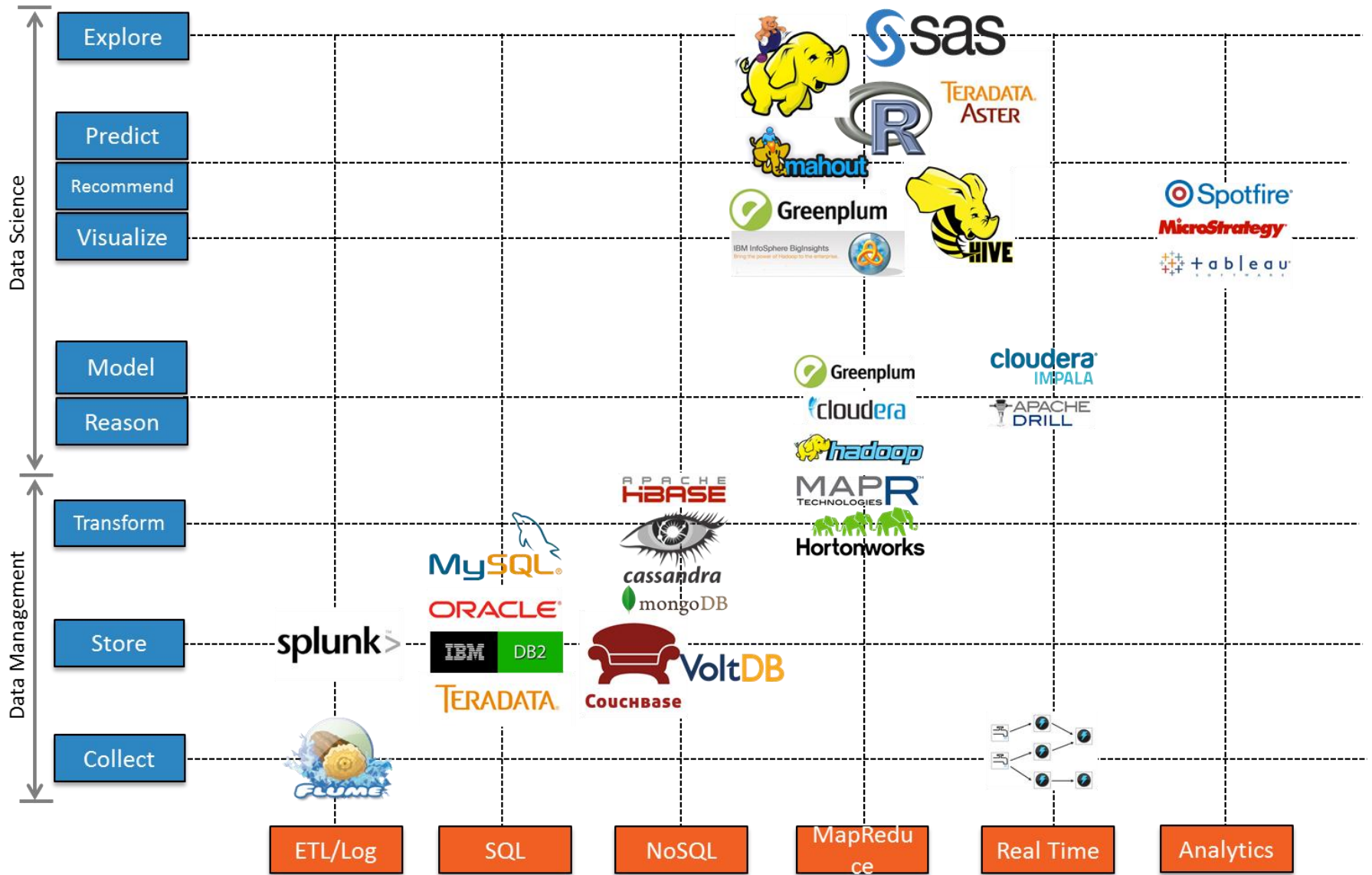
in one survey were unsure of how much of their data was inaccurate

Veracity
UNCERTAINTY OF DATA

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



Big Data Technology



Big Data – Technology, Platforms & Products

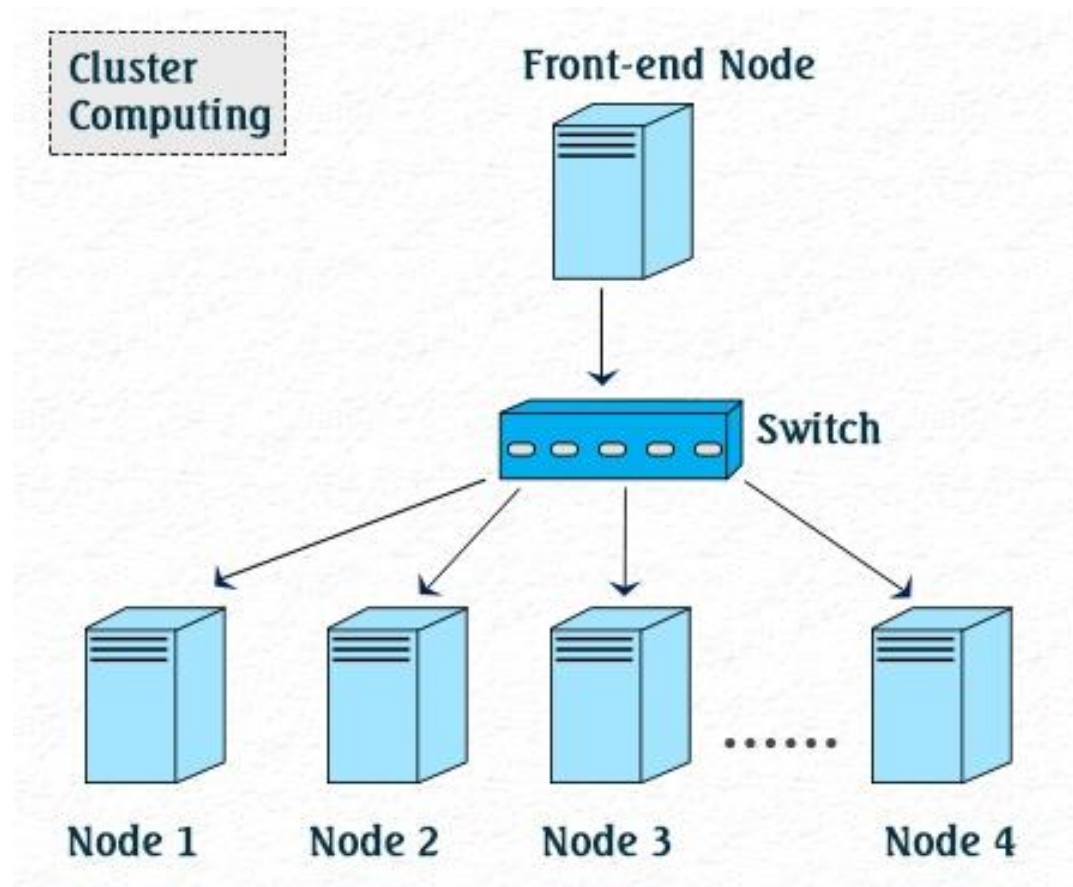
NoSQL



MongoDB เป็น open-source document database โดยเป็นฐานข้อมูลแบบ NoSQL คือ ไม่มี relation (ความสัมพันธ์) ของตารางแบบ SQL ทั่วไป แต่จะเก็บข้อมูลเป็นแบบ JSON (JavaScript Object Notation) แทน การบันทึกข้อมูลทุกๆ record ใน MongoDB เราจะเรียกมันว่า Document ซึ่งจะเก็บค่าเป็น key และ value จะเห็นว่ามันก็คือ JSON

“Big Data ไม่ใช่ Hadoop
แต่ Hadoop เป็นส่วนหนึ่ง
ของ Big Data”

Hadoop Concept



Hadoop Ecosystem

- Apache Hadoop (HDFS + YARN)
- Map Reduce
- Apache Hive
- Spark
- Cloudera Impala
- Sqoop
- Hbase
- Apache Ambari

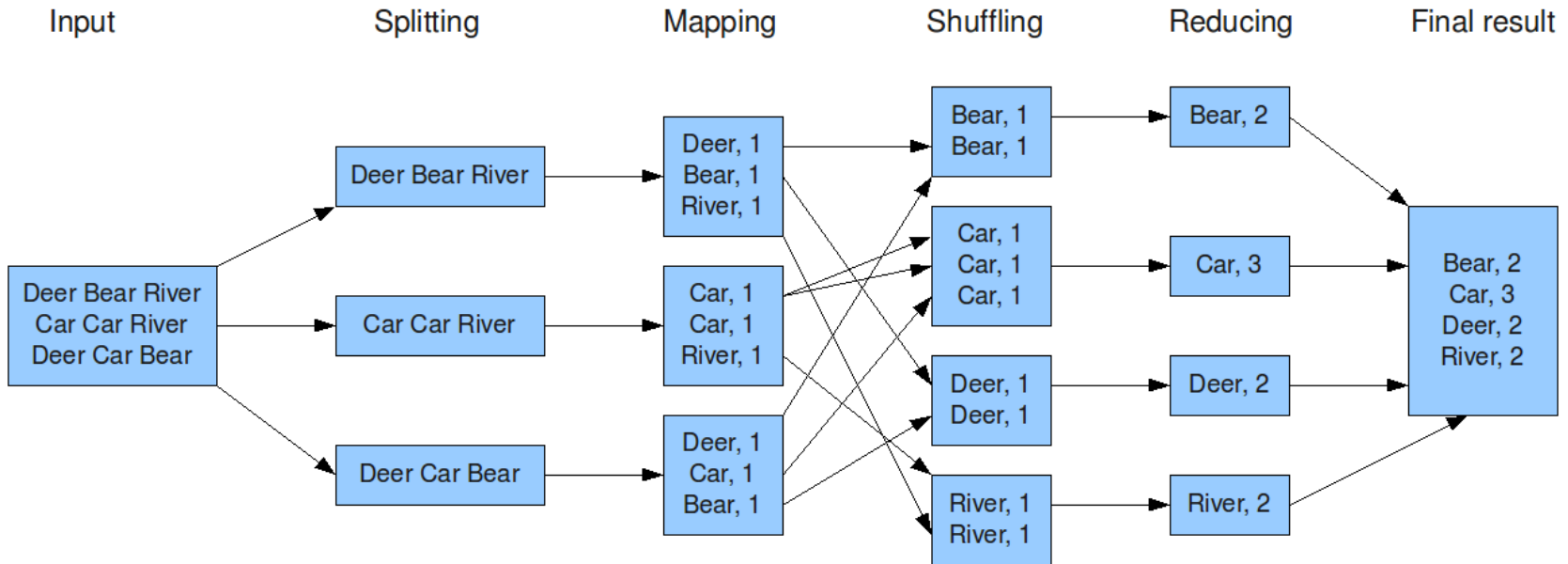
Apache Hadoop (HDFS + YARN)



- เป็นพื้นฐาน ของ Big Data Tools โดยเฉพาะ HDFS ที่เป็น File System ที่อยู่บน Distributed system

Map Reduce

The overall MapReduce word count process



- เป็น programming model ที่ช่วยในการใช้ทรัพยากรให้คุ้มค่า โดยระบบจะทำการกระจาย Task ไปรันแบบ Parallel บนเครื่องหลายๆ เครื่อง (ลดขั้นตอนการดำเนินการ)

Apache Hive



- เป็นเครื่องมือสำหรับผู้ต้องการสืบค้น (Query) ข้อมูลที่เก็บใน HDFS ด้วยภาษาลักษณะ SQL โดย Hive จะทำหน้าที่ในการแปล SQL like ให้มาเป็น Map/Reduce แล้วก็ทำการรันแบบ Batch

Spark



- เป็น Data Processing Framework ตัวหนึ่งที่นิยมกันแพร่หลาย Spark ทำงานได้รวดเร็วกว่าตัว Hadoop เพราะ Hadoop ทำงานบน Disk แต่ Spark ทำงานบน Memory

Cloudera Impala



- เป็นเครื่องมือสำหรับผู้ต้องการสืบค้น (Query) มีการทำงานคล้ายๆ Hive แต่ที่แตกต่างคือ Impala จะทำงานกับข้อมูลที่อยู่บน Memory ซึ่งแน่นอนว่าต้องเร็วกว่าแต่มันก็ยังมีข้อจำกัดในเรื่อง Memory

Sqoop



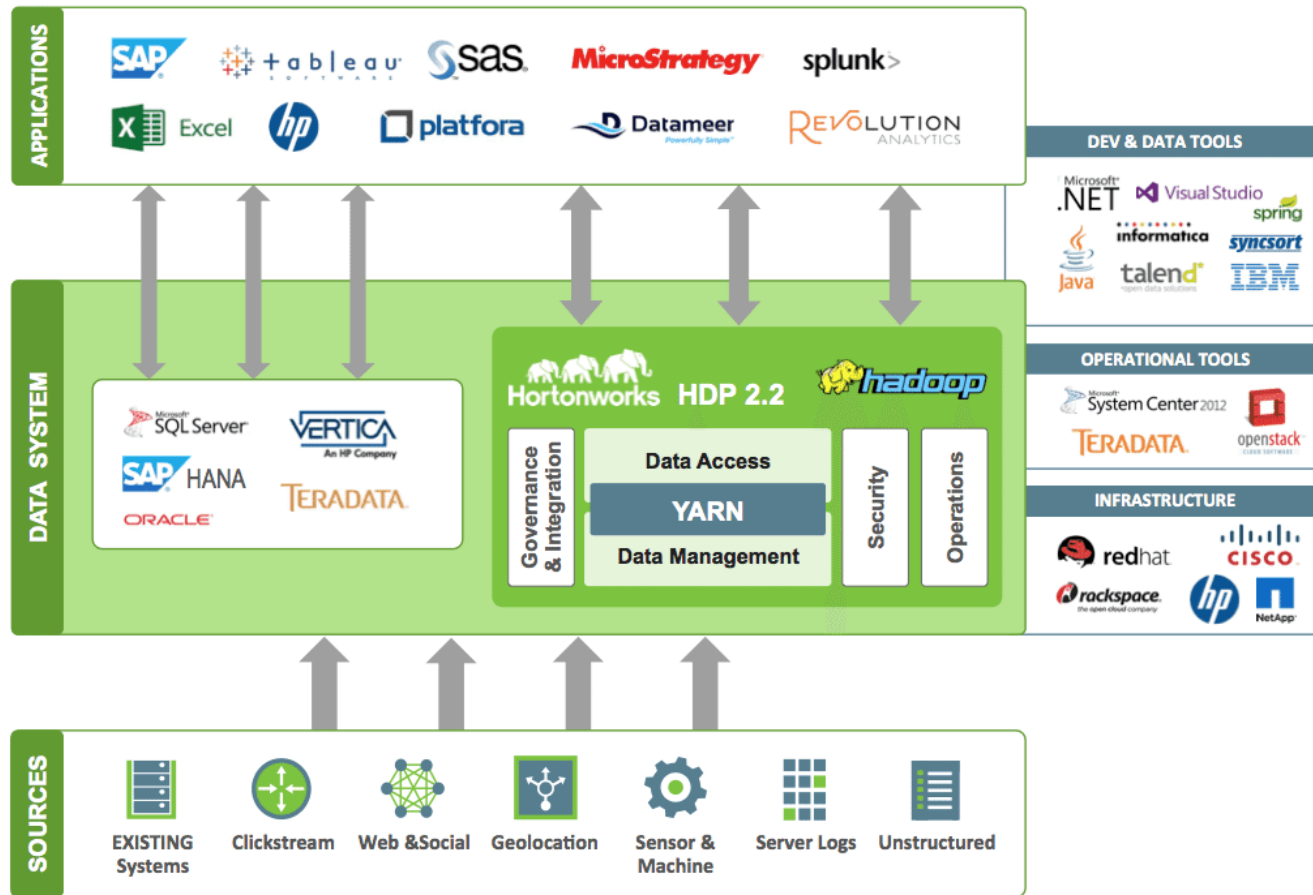
- เป็น Framework ที่จัดการการถ่ายโอนข้อมูลระหว่างข้อมูลรูปแบบ Table บน RDBMS กับข้อมูลรูปแบบ HDFS บน Hadoop

Apache HBase



- เป็น DBMS ที่อยู่บน Hadoop เหมาะสำหรับการ Read/Write ข้อมูลที่เป็น Realtime

Hortonworks



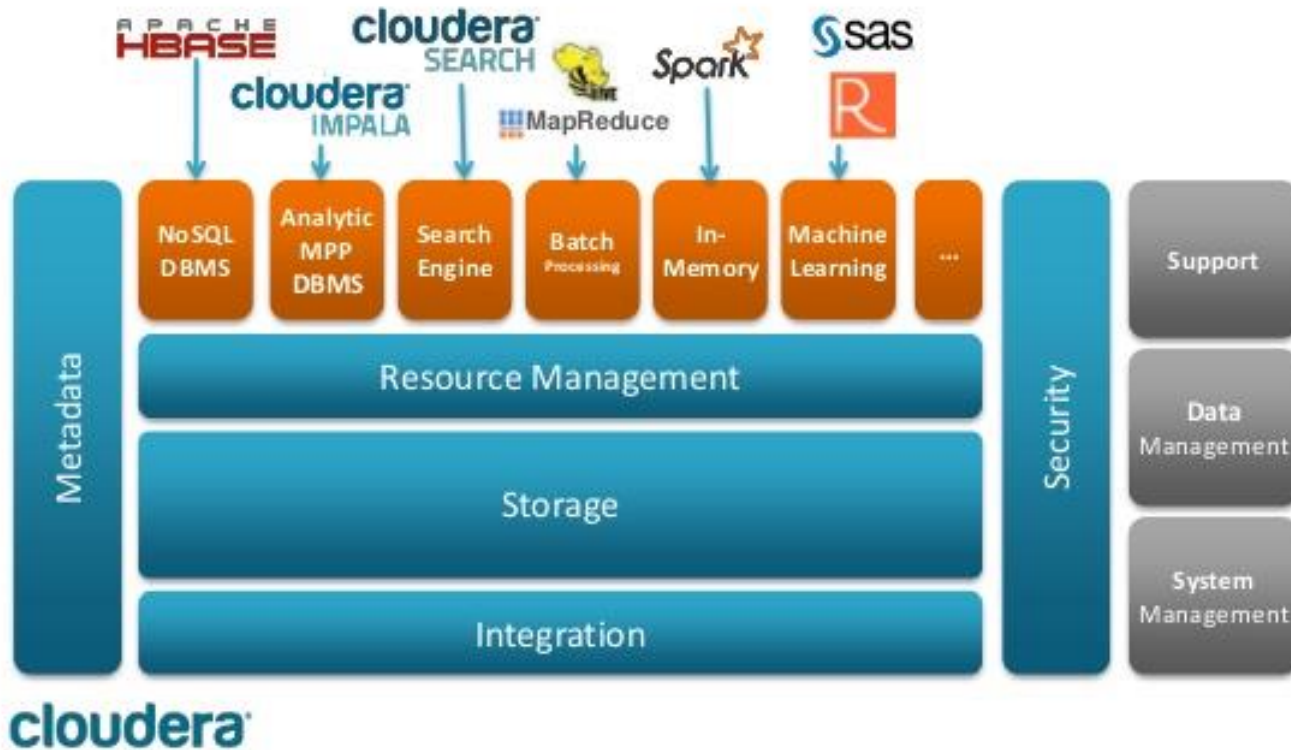
Apache Ambari



- โปรแกรมบริหาร Cluster ที่เป็น Opensource ช่วยทำให้เพิ่มประสิทธิภาพในการบริหาร Server

Cloudera

CDH: the App Store for Hadoop



Lab



Basic + a b | e a u

Narongrich Chanprasert

BI Consultant

Depth First Co.,Ltd.

Data usage



รายงานระดับปฏิบัติการ และรายงานสำหรับผู้บริหาร

- รายงานสำหรับระบบงาน
ปรกติ เพื่อให้ผู้ใช้งานระดับ
ปฏิบัติการ

- รายงานนำเสนอในรูปแบบ
ตาราง หรือออกเป็น

Spreadsheet

- รายงานนำเสนอข้อมูลที่เป็น
Data เป็นข้อมูลที่แสดง
รายละเอียด

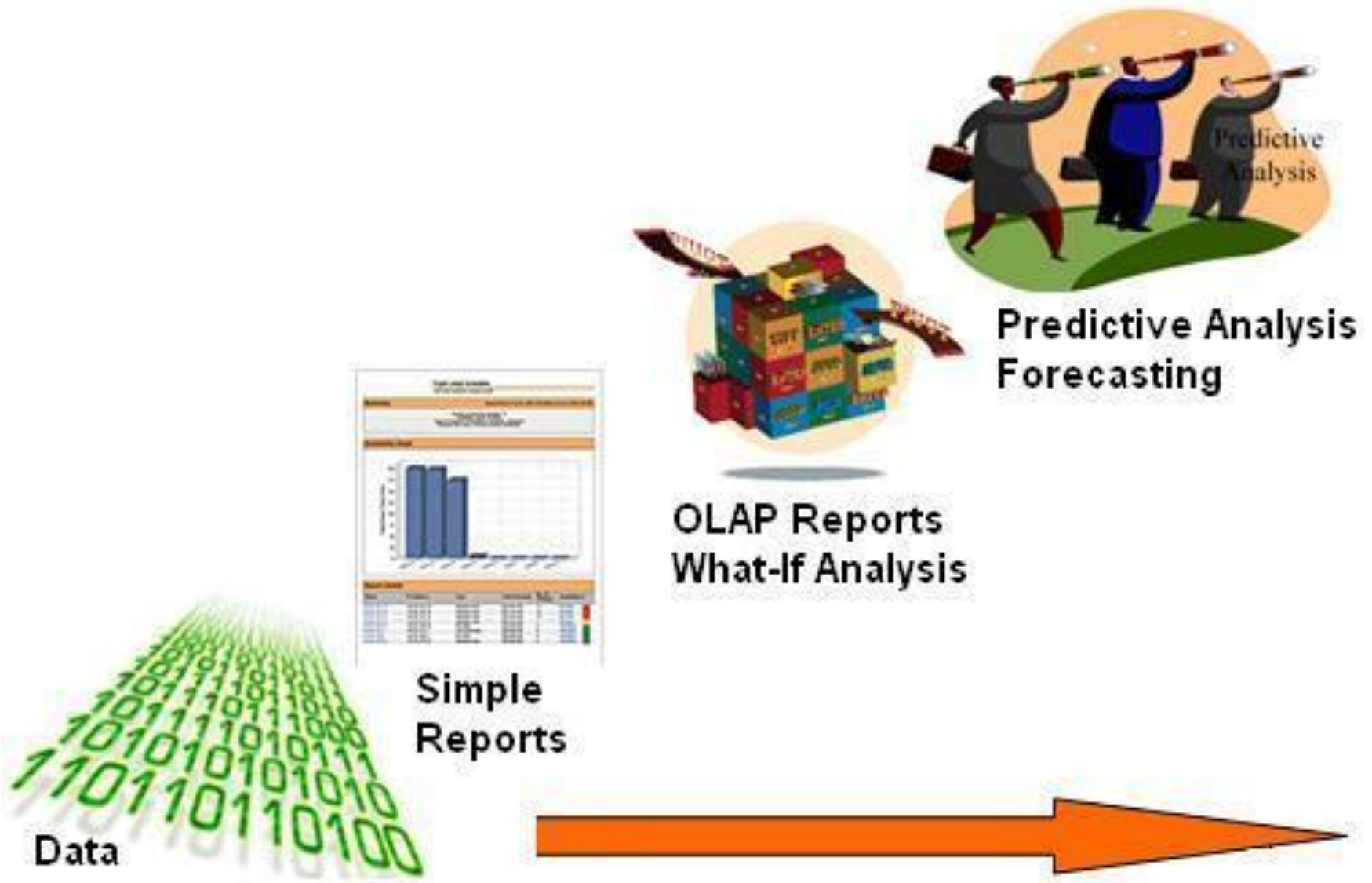
- รายงานสำหรับระบบงาน
ปรกติ เพื่อให้ผู้ใช้งานระดับ
ผู้บริหาร

- รายงานนำเสนอในรูปแบบ
กราฟ แผนภาพ หรือตาราง

Pivot

- รายงานนำเสนอข้อมูลที่เป็น
Information เป็นข้อมูล
เชิงสรุป

Roadmap to embrace Business Intelligence



Common Business Intelligence Process

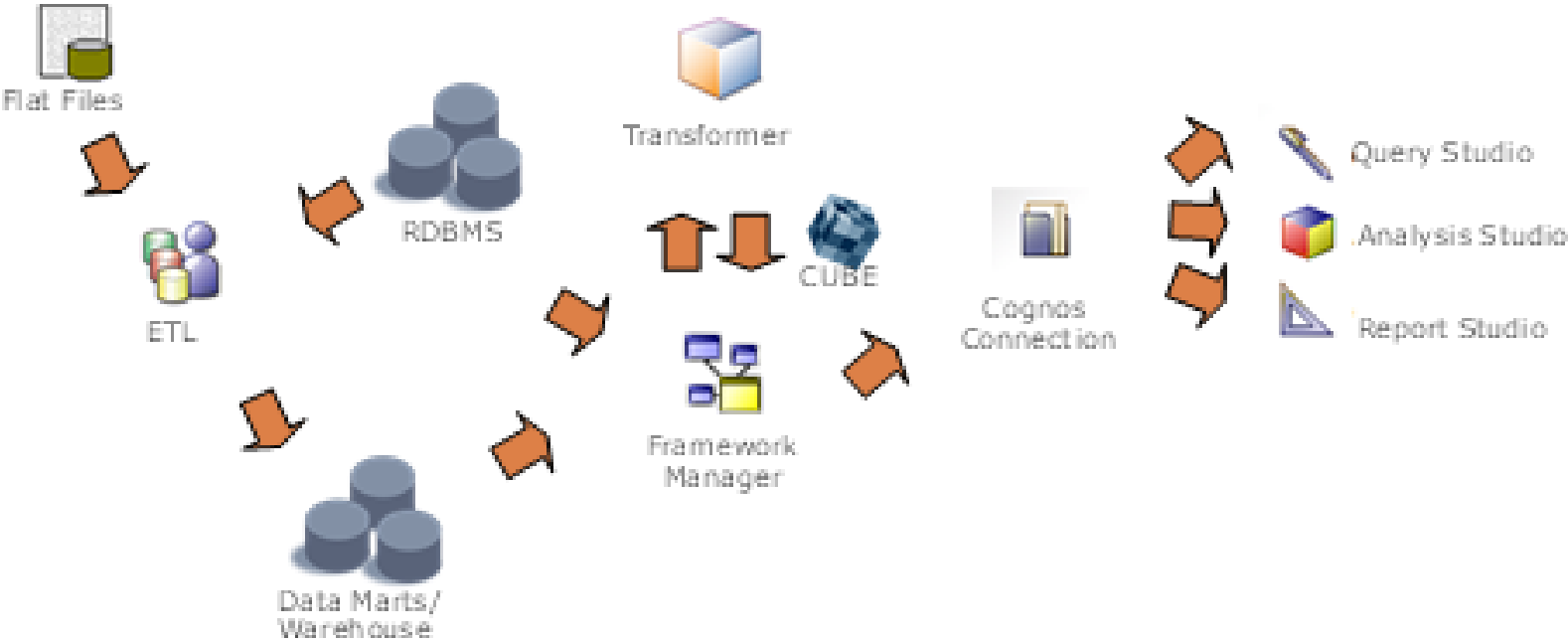
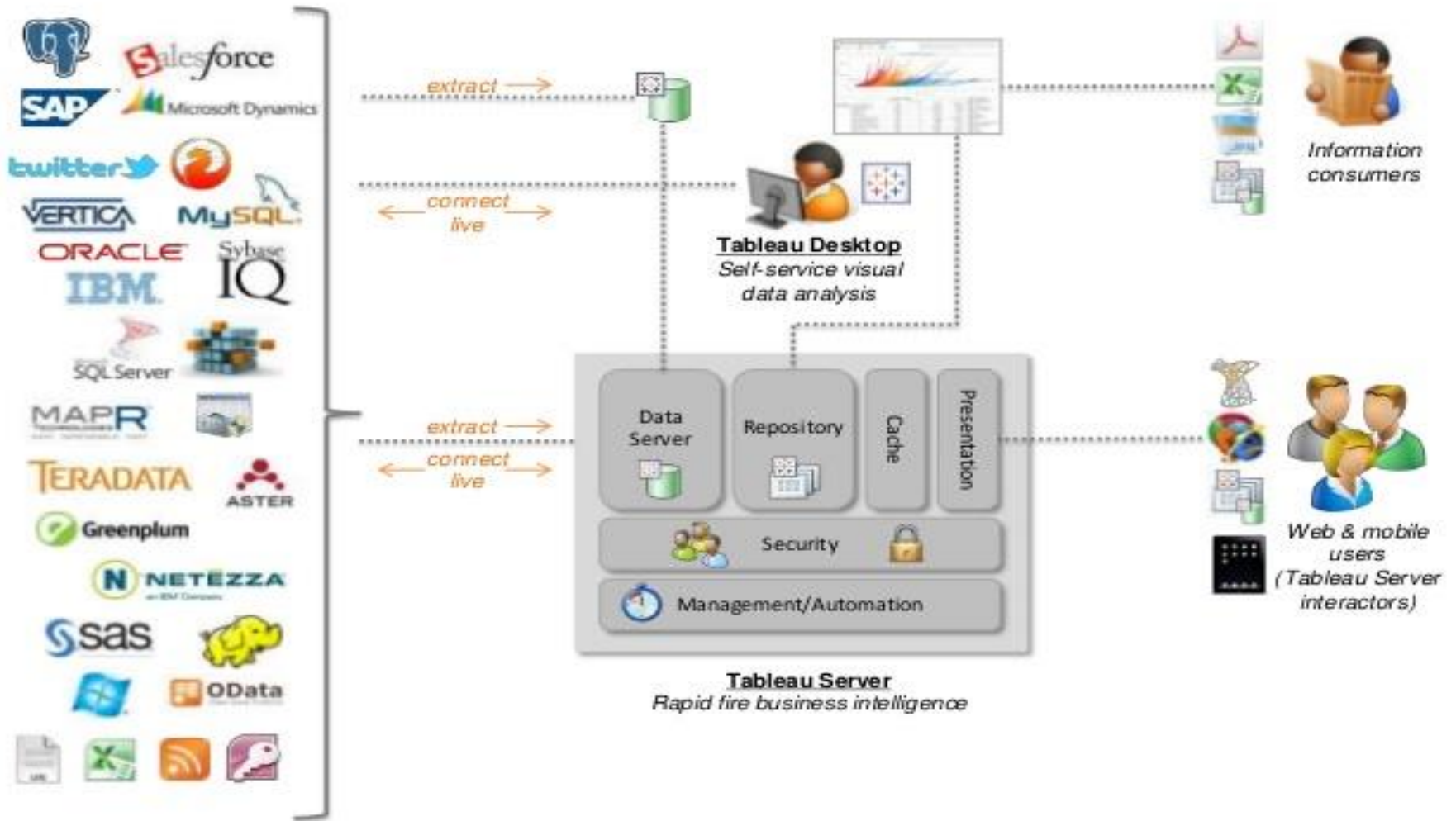


Tableau Architecture



Basic Data Warehouse Concept

•Dimension

- คือ มุมมองในการวิเคราะห์ข้อมูล ส่วนใหญ่เป็นข้อมูลเชิงคุณภาพ เช่น ข้อมูลชื่อ หรือว่ารายละเอียด ต่างๆ

•Measure

- คือ ข้อมูลเชิงปริมาณที่มีค่าวัดมากน้อยได้ สามารถนำมาประมวลผลเชิงผลรวม **Aggregation** ได้ เช่น ผลรวม, ค่าเฉลี่ย, **min, max** เป็นต้น