# Outline

- Big data platform and Technologies

  - IBM Big data platform

- Digging into Big data technology

  - Big data technology stack

  - Big data analytics platforms and software

# Big data platform

Comprehensive,

enterprise-ready,

integrated

# Main tasks in Big data

- Tasks within the domain of Big Data often involve data mining as a prevalent method, yet at a larger scale.
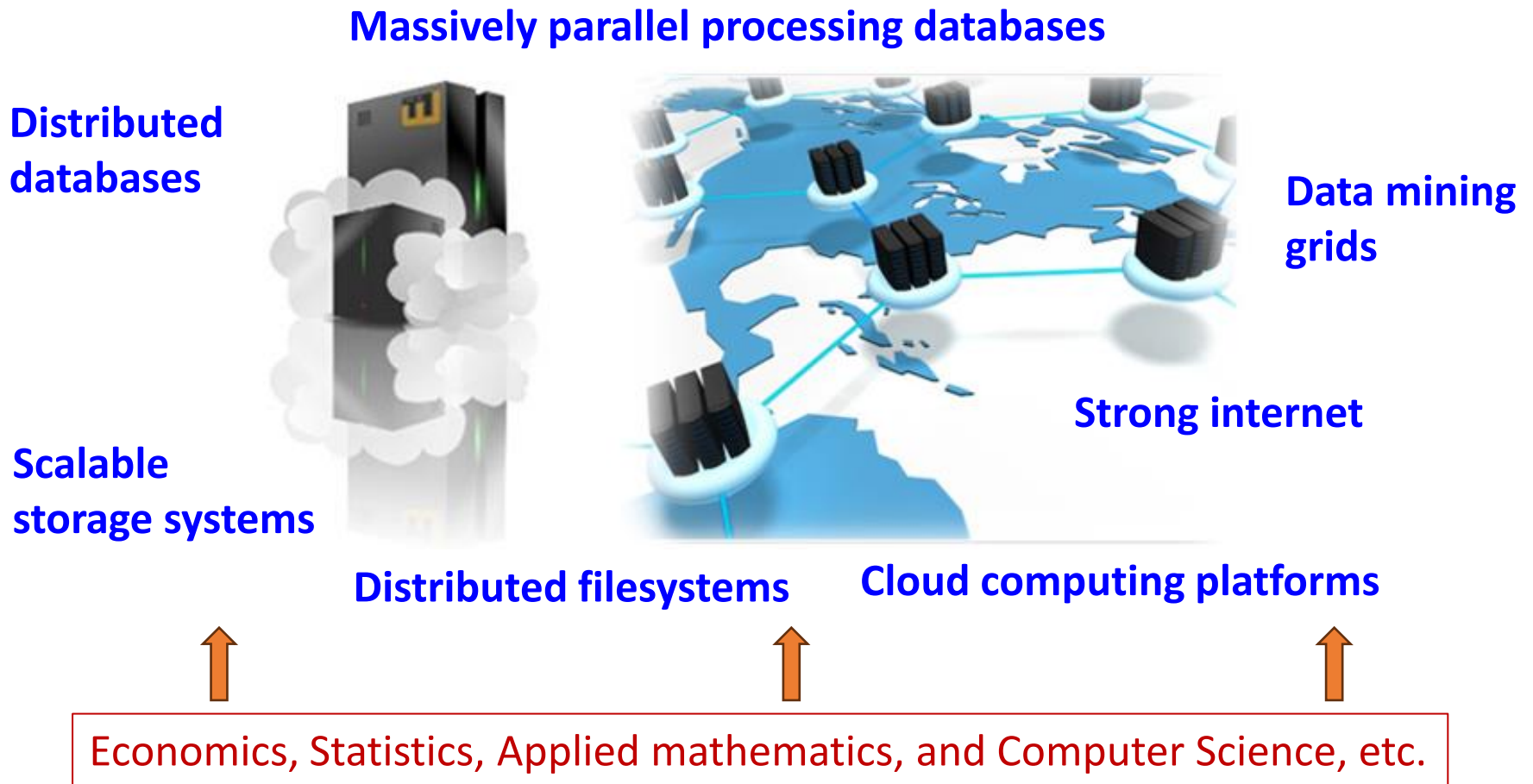


Data aggregation



Data analysis



Data manipulation



Data visualization

# Big data are multidisciplinary

- Technologies in Big data involves multidisciplinary studies.

**Massively parallel processing databases**

**Distributed databases**

**Data mining grids**

**Strong internet**

**Scalable storage systems**

**Distributed filesystems**

**Cloud computing platforms**

Economics, Statistics, Applied mathematics, and Computer Science, etc.

# A Big data platform should offer

**Comprehensive**

Every dimension of Big data challenge is addressed.

**Enterprise-ready**

Features of performance, security, usability and reliability included.

**Integrated**

Big data technologies to enterprise should be simplified and accelerated

Integration with information supply chain, including databases, data warehouses, and BI applications.

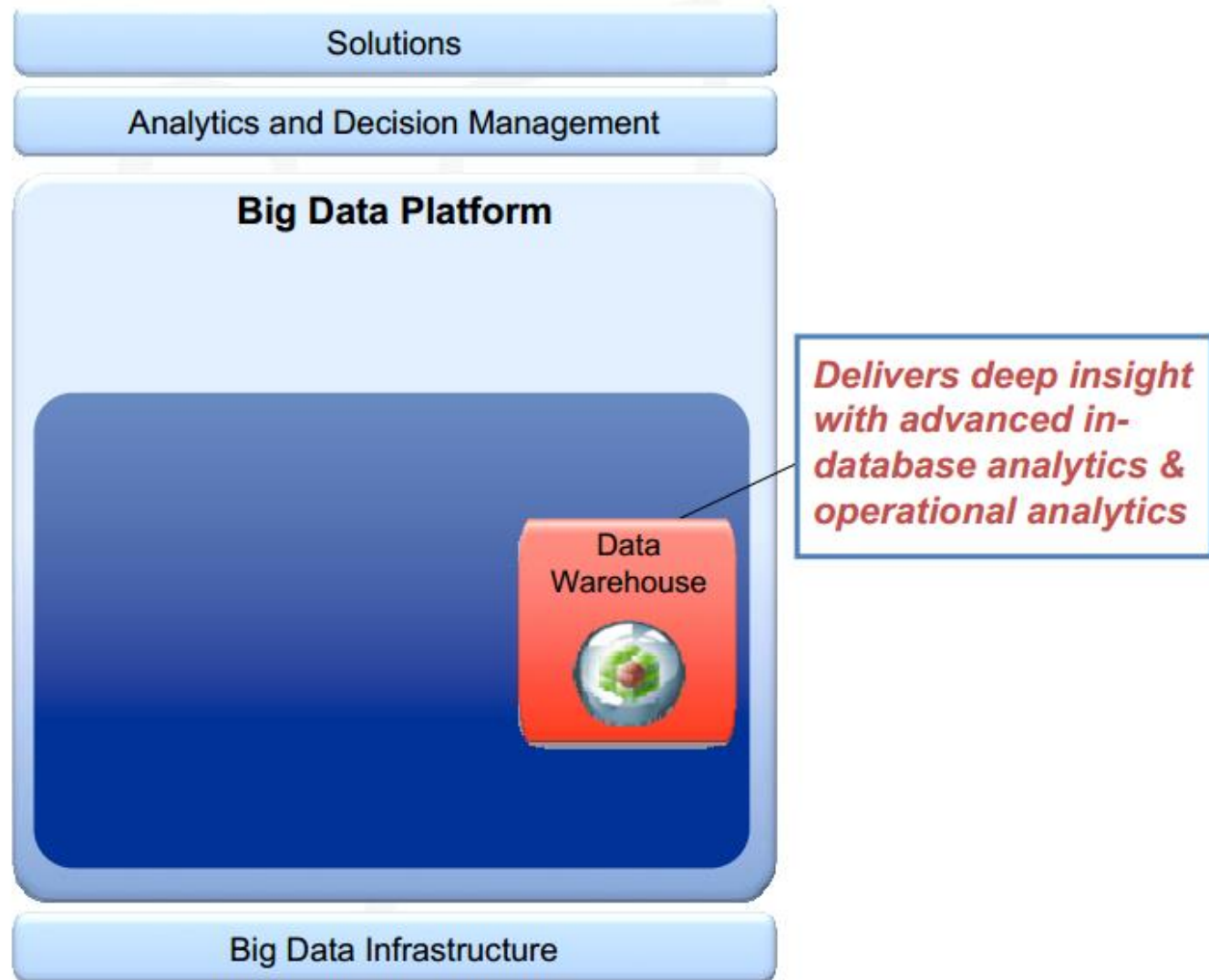- **Moreover, a Big Data platform should also offer**
  - Open-source based, low latency reads/updates, ad-hoc queries, scalability, extensible, robust fault-tolerant, minimal maintenance.
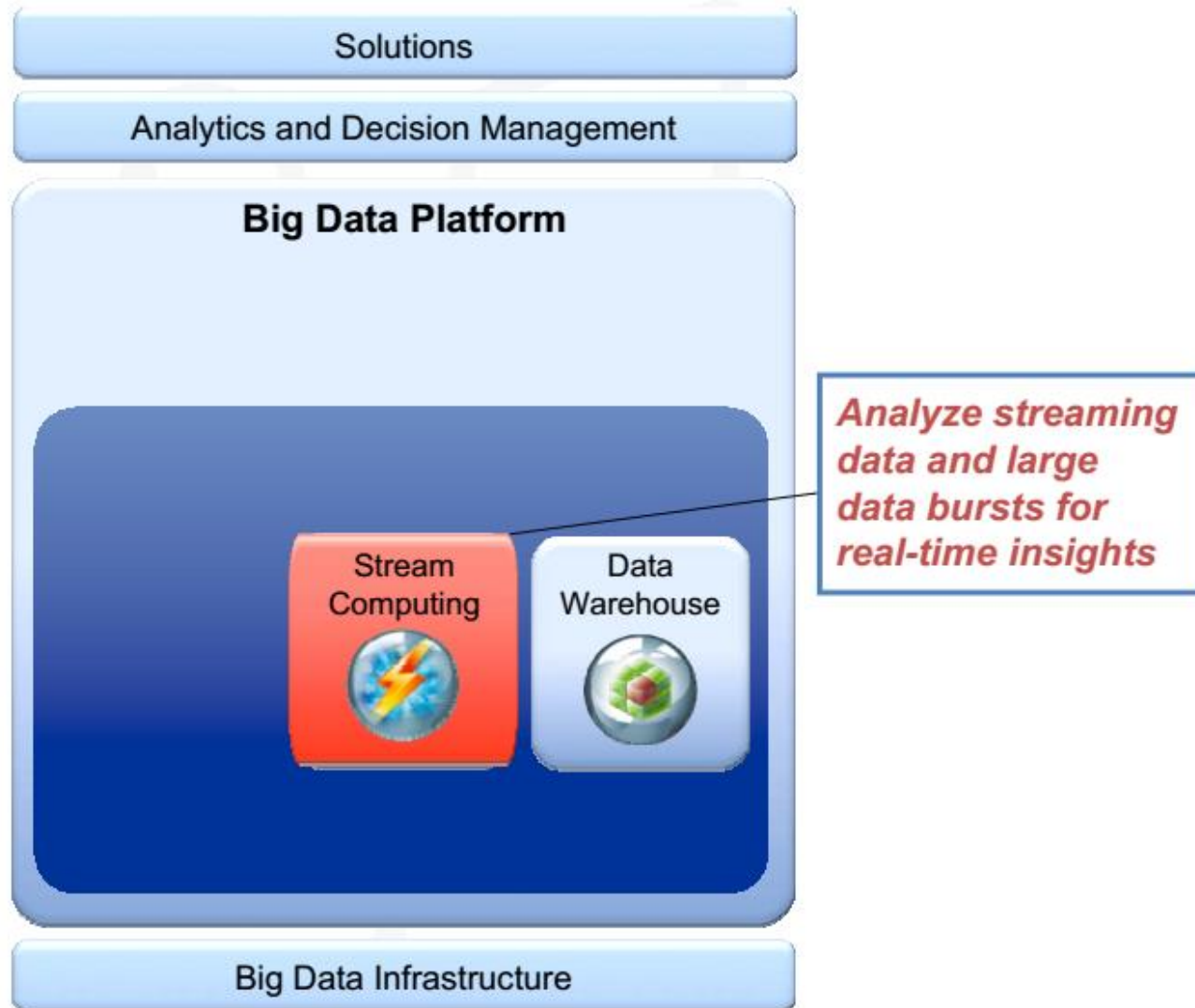
# IBM Big data platform

- Give a solution which is designed specifically with the needs of the enterprise in the mind.
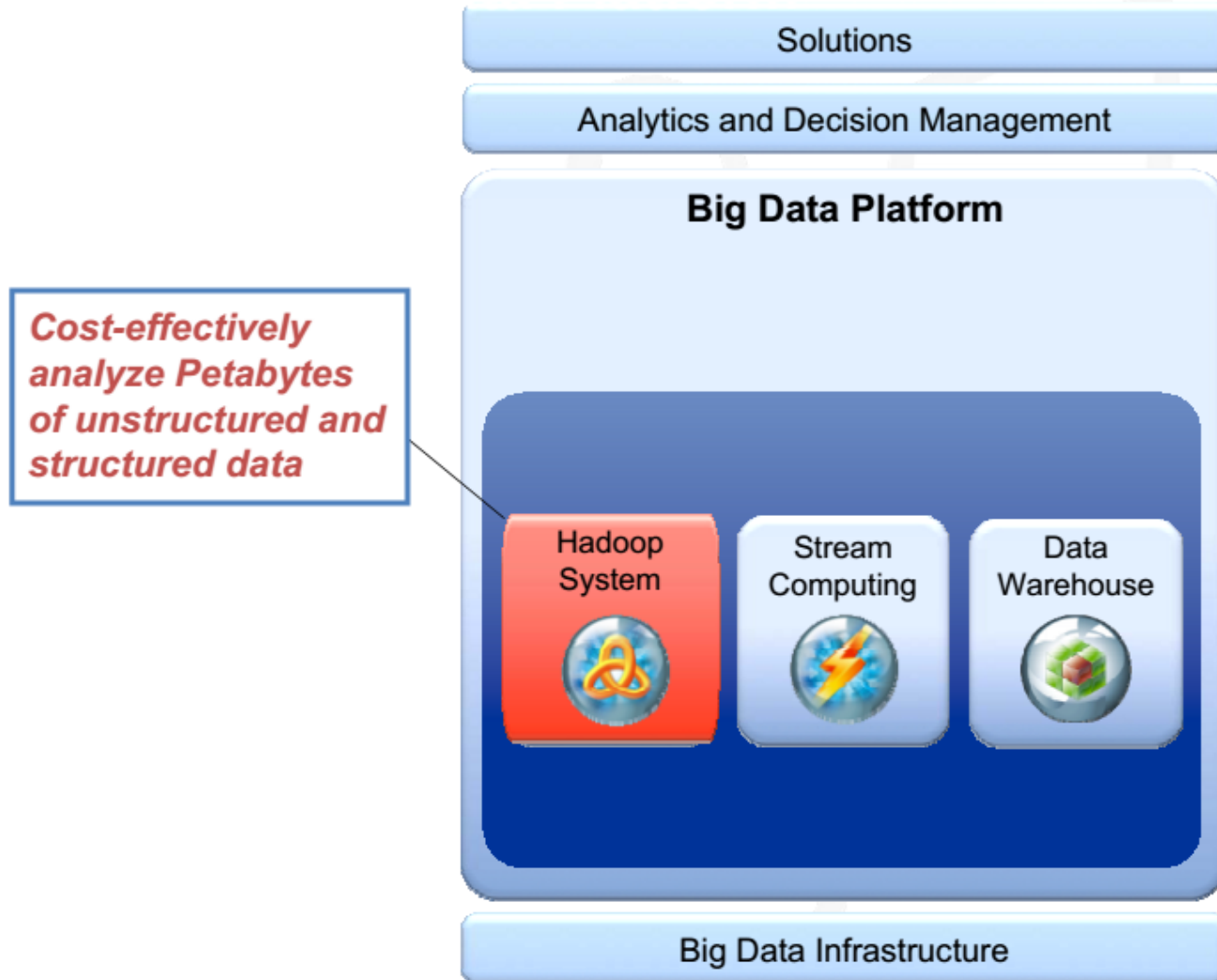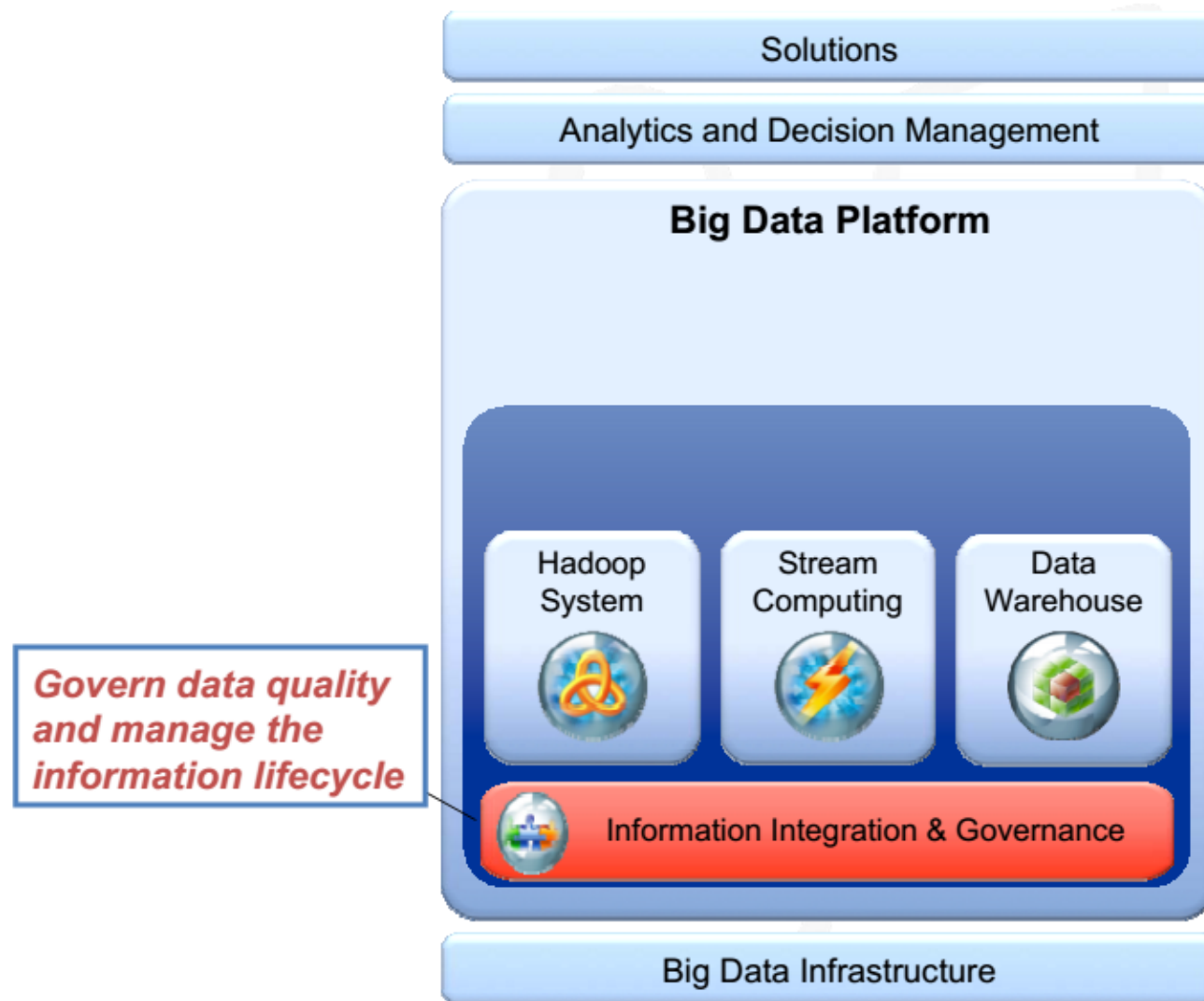
| Solutions |
|---|
| Analytics and Decision Management |

**Big Data Platform**

**Big Data Infrastructure**

# IBM Big data platform

# IBM Big data platform



Solutions

Analytics and Decision Management

**Big Data Platform**

Stream Computing

Data Warehouse

*Analyze streaming data and large data bursts for real-time insights*

Big Data Infrastructure

# IBM Big data platform



Solutions

Analytics and Decision Management

**Big Data Platform**

*Cost-effectively analyze Petabytes of unstructured and structured data*

Hadoop System

Stream Computing

Data Warehouse

Big Data Infrastructure

# IBM Big data platform

# IBM Big data platform

Solutions

Analytics and Decision Management

**Big Data Platform**

*Speed time to value with analytic and application accelerators*

Accelerators

| Hadoop System | Stream Computing | Data Warehouse |

Information Integration & Governance

Big Data Infrastructure

# IBM Big data platform

*Discover, understand, search, and navigate federated sources of big data*

**Solutions**

**Analytics and Decision Management**

**Big Data Platform**

Visualization & Discovery

Application Development

Systems Management

Accelerators

Hadoop System

Stream Computing

Data Warehouse

Information Integration & Governance

**Big Data Infrastructure**

# Components in a Big data platform

# Digging into Big data technology
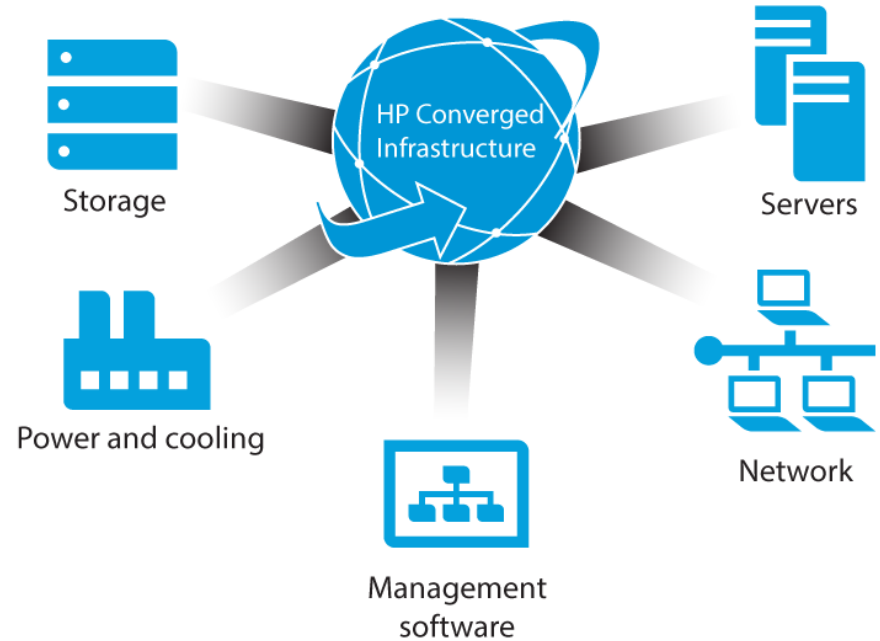
Digging deeper, better insights

# Big data technology stack

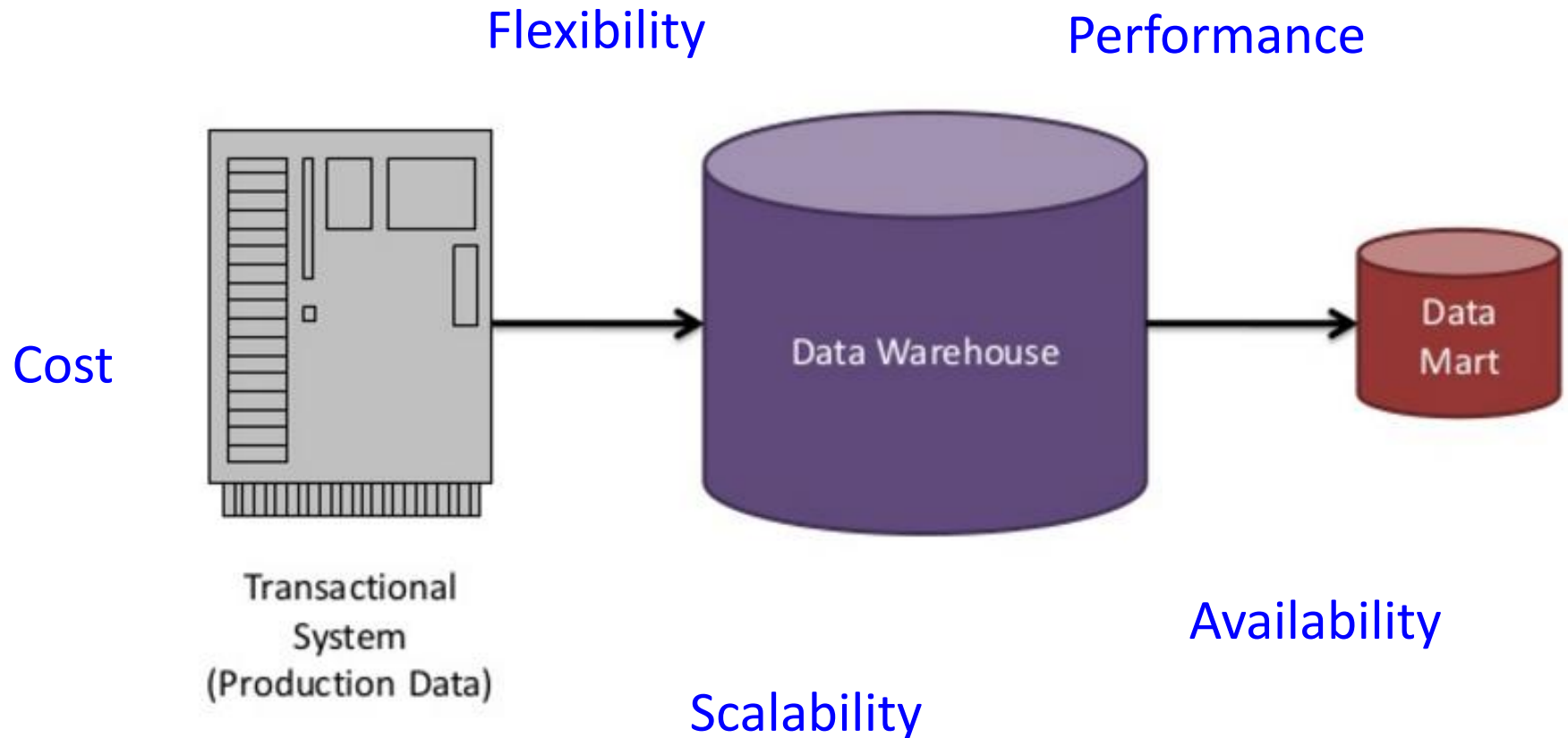# Layer 0: Redundant physical infrastructure

- The **physical infrastructure** is the lowest level.
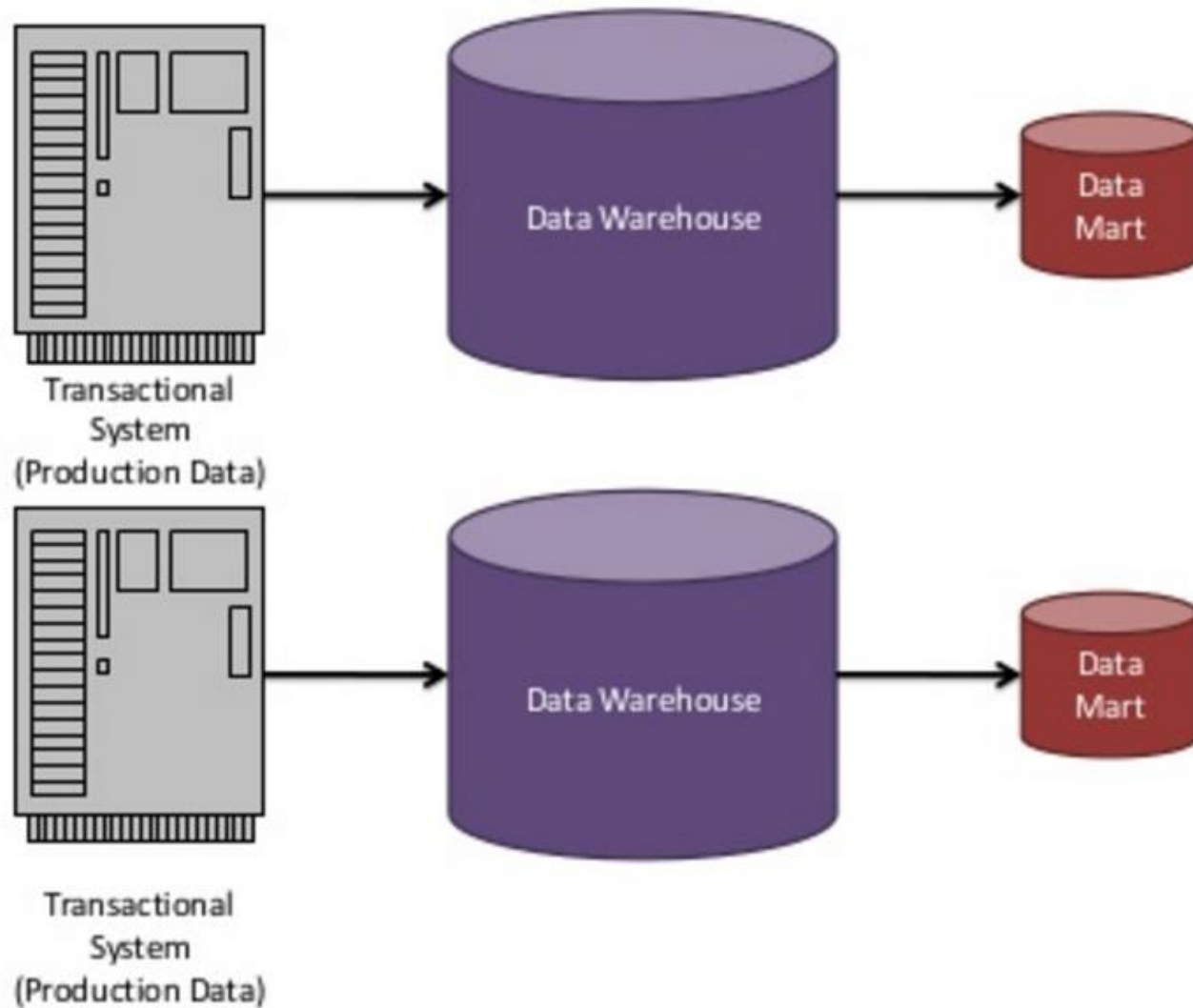    - Hardware, network, etc.



- Your company might already have a data center or made investments in physical infrastructures.

- Hence, you may want to find a way to utilize existing assets.
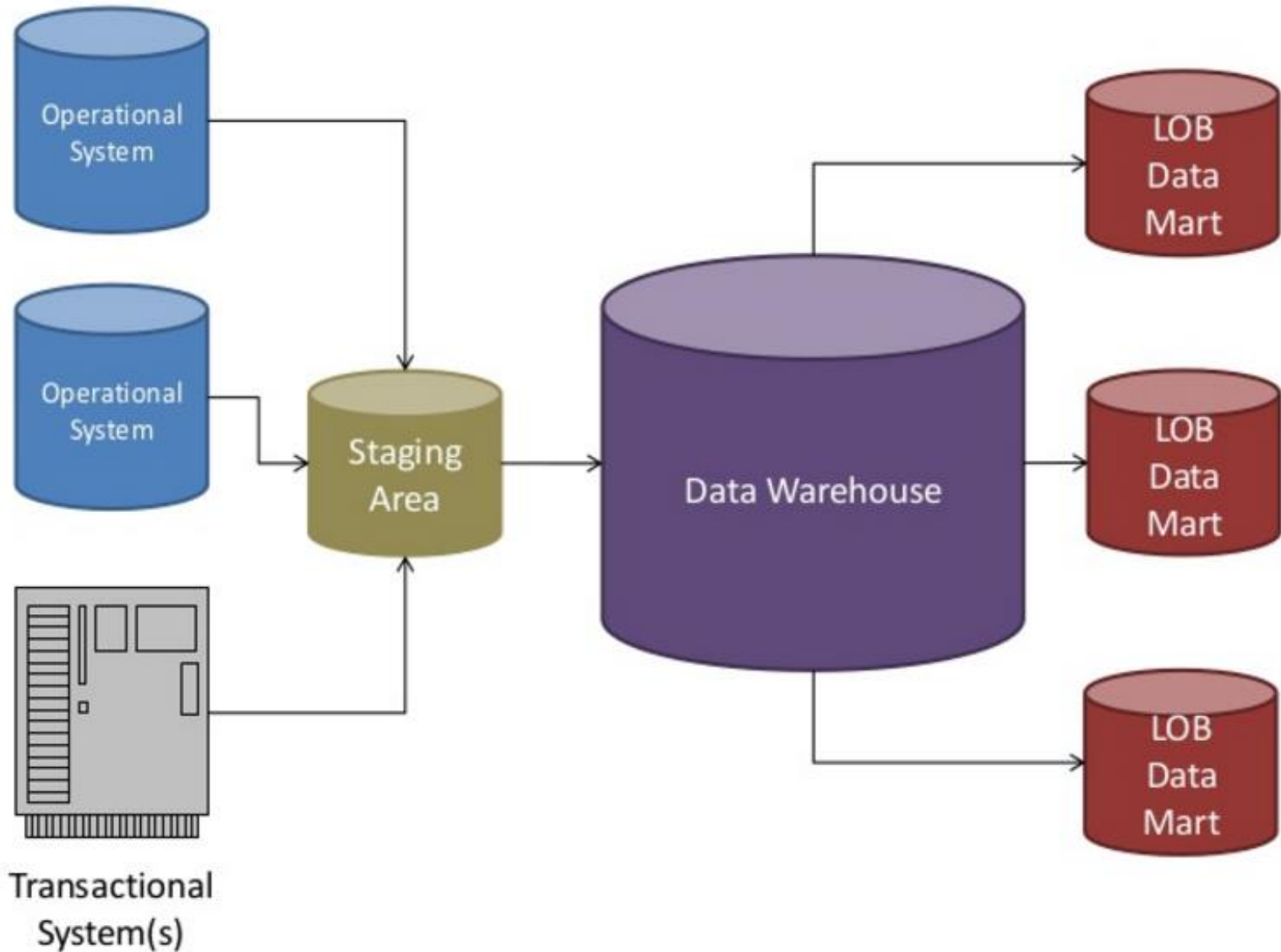
# Where most of this began?

- A prioritized list of these principles should include statements about the following



Flexibility

Performance

Cost

Data Warehouse

Data Mart

Transactional System (Production Data)

Availability

Scalability

# It grows bigger..



Transactional System (Production Data) → Data Warehouse → Data Mart

Transactional System (Production Data) → Data Warehouse → Data Mart

# ….then very big

# Why redundant?

- Most big data implementations need to be highly available.

- That is, networks, servers, and physical storage must be both resilient and redundant.

- A system is resilient to failure or changes when sufficient redundant resources are in place, ready to jump into action.

# Layer 1: Security infrastructure

- Security and privacy requirements for big data are similar to those for conventional data environments.

- They must be closely aligned to specific business needs.

**Data access**

The data should be available only to those who have a legitimate business need for examining or interacting with it.

Protection from unauthorized usage or access are offered by most APIs.

**Application access**

**Data encryption**

Most challenging, extremely stress the systems' resources
Encrypt only data elements that require this level of security

The inclusion of mobile devices and social networks exponentially increases both the amount of data and the opportunities for security threats.

**Threat detection**

# Layer 2: Operational databases

- The core of any Big data environment is <span style="color:red">database engines</span> holding collections of data elements relevant to a business.

**Atomicity**

If any part of the transaction or the underlying system fails, the entire transaction fails.

Only transactions with valid data will be performed.
**Consistency**

**Isolation**

Multiple simultaneous transactions do not interfere with each other. All valid transactions will execute until completed and in the order, they were submitted for processing.

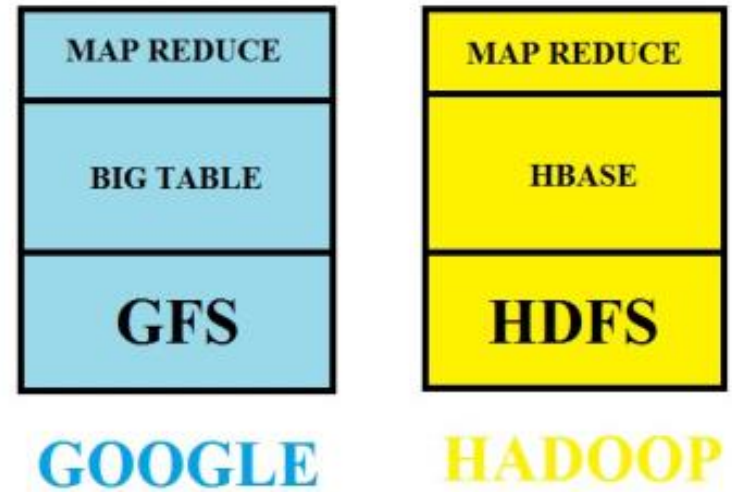After the data from the transaction is written to the database, it stays there "forever."
**Durability**

# Layer 3: Organizing Data Services and Tools

- Prepare an ecosystem of tools and technologies to gather and assemble data in preparation for further processing

- Technologies in this layer include the following:

  - A distributed file system

  - Serialization services

  - Coordination services

  - Extract, transform, and load (ETL) tools

  - Workflow services

# Hadoop, MapReduce and Big Table

- New technologies to store, access, and analyze huge amounts of data



| MAP REDUCE | MAP REDUCE |
| BIG TABLE | HBASE |
| GFS | HDFS |
| GOOGLE | HADOOP |

- Proved to be the sparks that led to a new generation of data management.

- Addressing one of the most fundamental problems: the capability of processing massive amounts of data efficiently, cost effectively, and in a timely fashion.

# Layer 4: Traditional and advanced analytics

- **What** does your business now **do with all the data** in all its forms to try to make sense of it **for the business**?

  - Managing big data holistically requires different analysis approaches, *depending on the problem being solved*, to help the business to successfully plan.

  - Some analyses will use a traditional data warehouse, while the others will take advantage of advanced predictive analytics.

- **Key techniques:** Analytical data warehouses and data marts, Big data analytics, Reporting and visualization, etc.

# Big data platform and analytics software

- Features of Big data platform and analytics software

Data ingestion, Data management, ETL and Warehouse, Hadoop system and Stream Computing

Analytics/Machine learning, Content management, Data integration and governance

Provide efficiency in workplace
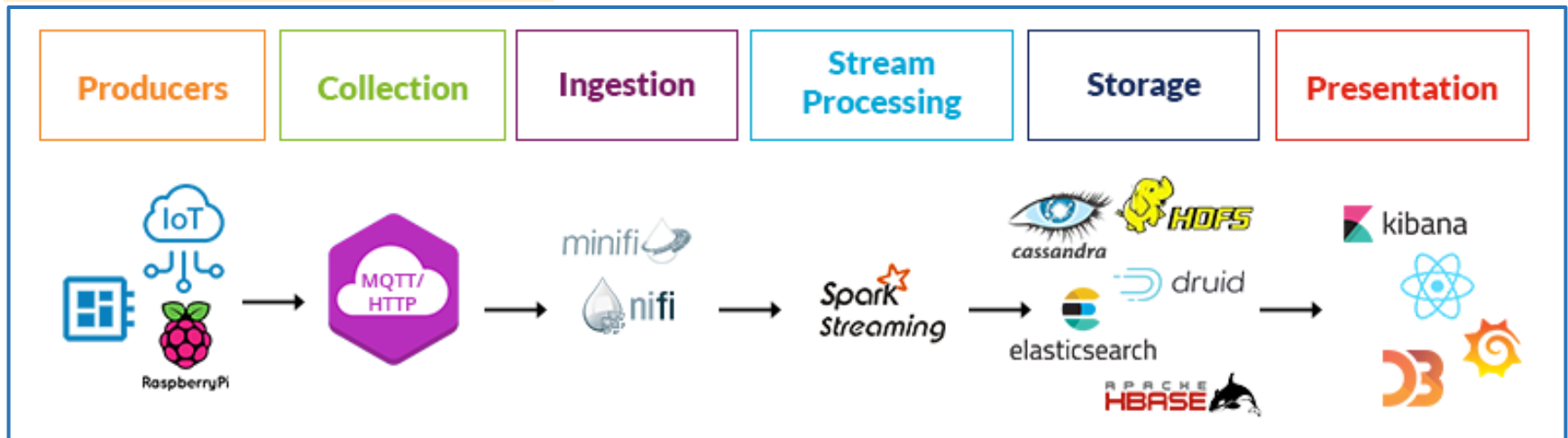Provide accurate data
Give answer to complex questions
It is secure

# Big data analytic platform tools

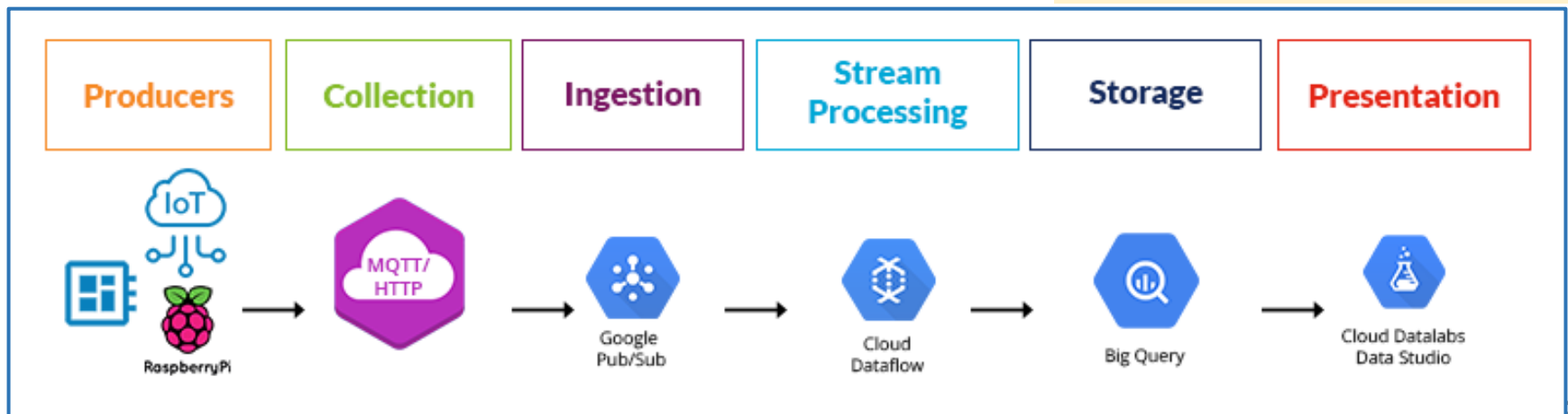- There are some key Big data analytic platform tools available for enterprise use

# IoT Analytics Platform
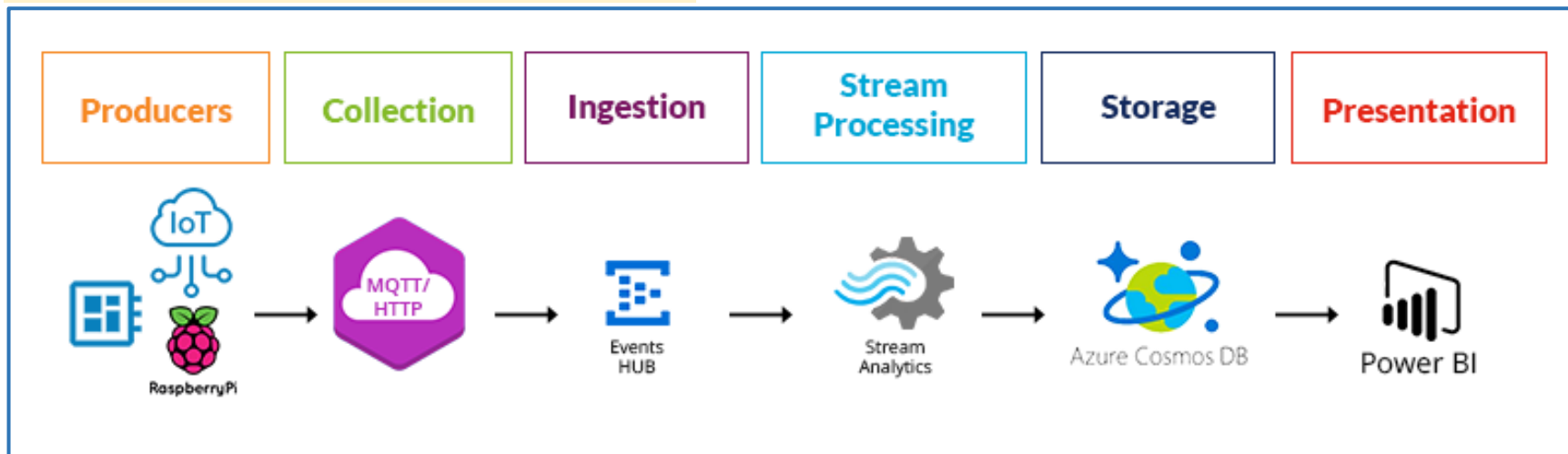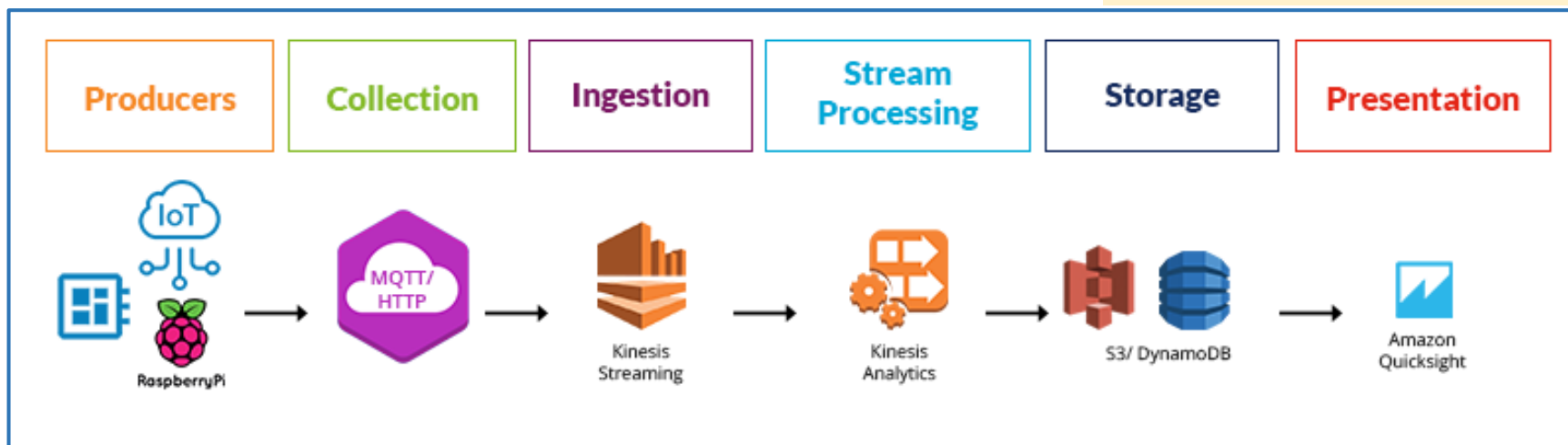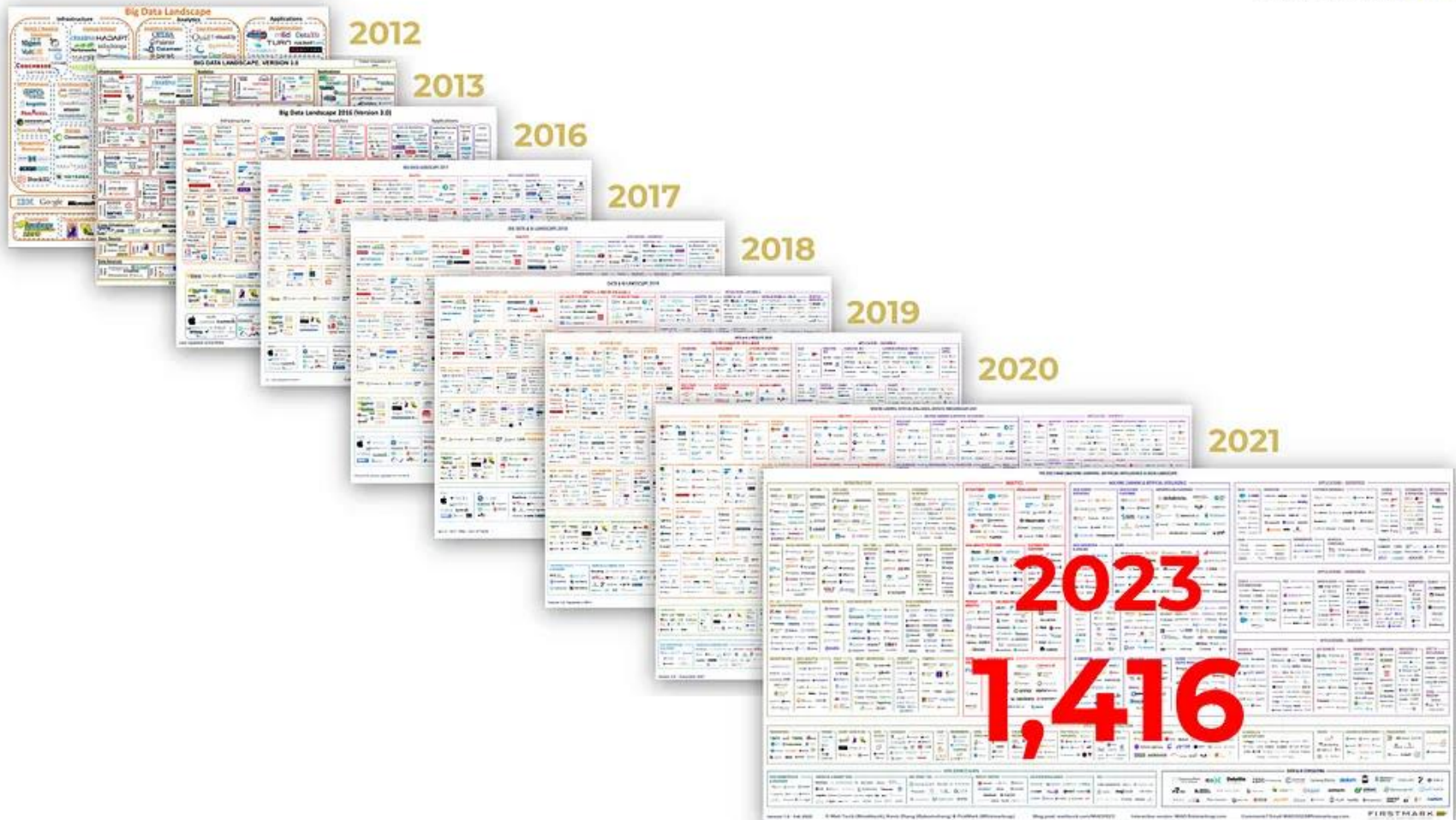# for Real-Time Data Ingestion

## Microsoft Azure IoT Architecture



## AWS IoT Architecture

The 2023 MAD (ML/AI/Data) Landscape

# Big data vs. Data science

- In Data science, the data can be of all sizes, which is related to a business or scientific case.

- Big data offers techniques to handle large-scale data at different steps.



| OBTAIN | SCRUB | EXPLORE | MODEL | INTERPRET |
| --- | --- | --- | --- | --- |
| **O** | **S** | **E** | **M** | **N** |
| Gather data from relevant sources | Clean data to formats that machine understands | Find significant patterns and trends using statistical methods | Construct models to predict and forecast | Put the results into good use |