

Big Data Analytics

Fall 2025

Lecture 1

Dr. Tariq Mahmood



Defining Big Data



- Extremely large data sets that have grown beyond the ability to manage and analyze with traditional data processing tools
- Conventional IT can no longer effectively handle either the size of the data set or the scale and growth of the data set
- **Primary difficulties:** acquisition, storage, searching, sharing, analytics and visualization



- *John R. Mashey popularized the term*
- *Big data is high-volume, high-velocity and high-variety information assets that demand cost-effective, innovative forms of information processing for enhanced insight and decision making.* – Gartner
- <https://cloud.google.com/learn/what-is-big-data>
- https://en.wikipedia.org/wiki/Big_data



List of Global Big Data Businesses

- Energy
- Retail (Online/Brick and Mortar)
- Insurance
- Finance
- Actuarial
- Telco
- Healthcare
- Crime Prevention
- Governance
- Traffic
- Agricultural
- Media (Print/Electronic)
- Advertising
- Software houses
- Academic Institutions
- E-Commerce
- Consultancy Companies (Projects – Training – Workshops)
- Trading/Brokerage
- Oil and Gas
- Automotive



Big Data and Cloud Solution Companies

- Amazon AWS
- MS Azure
- Google
- HP
- Oracle
- Dell
- SAP
- SAS
- Apache (ASF)
- Cloudera
- Databricks
- Informatica
- Snowflake
- Confluent
- Redis Labs
- Matillion
- VMWare
- MongoDB
- REDIS
- Neo4J
- Dell
- Intel
- Nvidia
- Cisco
- IBM

<https://builtin.com/big-data/big-data-companies-roundup>

<https://themanifest.com/big-data/companies>

<https://db-engines.com/en/ranking>



5 BIG DATA USE CASES

Customer
Sentiment
Analysis

Behavioral
Analytics

Predictive
Support

Fraud
Detection

Customer
Segmentation

BIG DATA USE CASES:



1. OPTIMIZE FUNNEL CONVERSION



5. MARKET BASKET ANALYSIS AND PRICING OPTIMIZATION



2. BEHAVIORAL ANALYTICS



6. PREDICT SECURITY THREATS



3. CUSTOMER SEGMENTATION



7. FRAUD DETECTION



4. PREDICTIVE SUPPORT



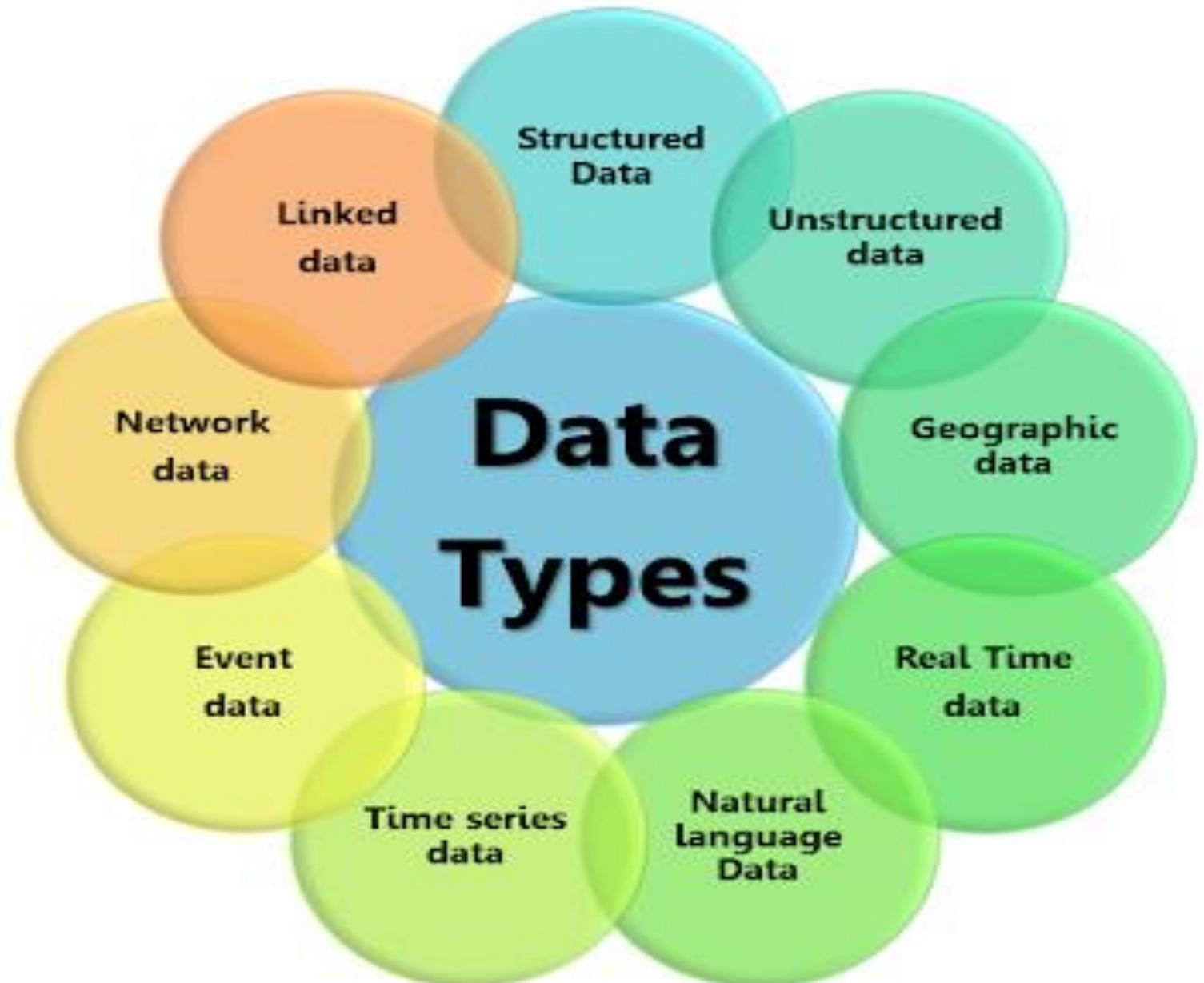
8. INDUSTRY SPECIFIC

Big Data opportunities across industries and use cases

Innovative analytic use cases are cutting across structured, unstructured and semi structured data

Finance	Government	Telecom	Manufacturing	Energy	Healthcare
<ul style="list-style-type: none">• Customer knowledge• Event marketing• Risk management	<ul style="list-style-type: none">• Law enforcement• Video surveillance and security• Traffic flow optimization	<ul style="list-style-type: none">• Analysis and customer retention• Analysis of network usage• Monitoring hearings and ad optimization	<ul style="list-style-type: none">• Supply chain optimization• Defect tracking• RFID Correlation• Warranty management	<ul style="list-style-type: none">• Weather forecasting• Natural resource exploration	<ul style="list-style-type: none">• Drug development• Scientific research• Evidence based medicine• Healthcare outcomes
Horizontal use cases					
<ul style="list-style-type: none">• Sentiment analysis• Social CRM / network analysis• Churn mitigation• Brand monitoring• Cross and Up sell• Loyalty & promotion analysis• Web application optimization	<ul style="list-style-type: none">• Marketing campaign optimization• Brand management• Social media analytics• Pricing optimization• Internal risk assessment• Customer behavior analysis• Revenue assurance	<ul style="list-style-type: none">• Logistics optimization• Clickstream analysis• Influencer analysis• IT infrastructure analysis• Legal discovery• Equipment monitoring• Enterprise search			

Sources: IDC: 2012 "Worldwide Big Data Technology and Services Forecast: 2011-2015, Gartner: 2012 "Big Data Drives Rapid Changes in Infrastructure and \$232 Billion in IT Spending Through 2016



Unstructured Data Types for Big Data Analysis



This slide is 100% editable. Adapt it to your needs and capture your audience's attention.



This slide is 100% editable. Adapt it to your needs and capture your audience's attention.



This slide is 100% editable. Adapt it to your needs and capture your audience's attention.



This slide is 100% editable. Adapt it to your needs and capture your audience's attention.



This slide is 100% editable. Adapt it to your needs and capture your audience's attention.



This slide is 100% editable. Adapt it to your needs and capture your audience's attention.



This slide is 100% editable. Adapt it to your needs and capture your audience's attention.



This slide is 100% editable. Adapt it to your needs and capture your audience's attention.



Why BDA can fail?

- Asking the wrong questions and solving the wrong problem
- Lack of right data
- Lack of right talent
- Not deploying products – we think deployment is the last step (DevOps)
- Not involving the business
- Selection of the wrong technology stack – lack of resource optimization
- Indifference to technological requirements



<https://www.kdnuggets.com/2021/09/sparkbeyond-avoid-data-science-projects-fail.html>

https://medium.com/@daniel_3607/addressing-the-85-data-project-failure-rate-is-your-companys-greatest-chance-to-succeed-da37967372d6

<https://www.globalbigdataconference.com/news/141927/85-of-big-data-projects-fail-but-your-developers-can-help-yours-succeed.html>

<https://hbr.org/2020/02/use-this-framework-to-predict-the-success-of-your-big-data-project>

<https://www.datascience-pm.com/project-failures/>

Disruption

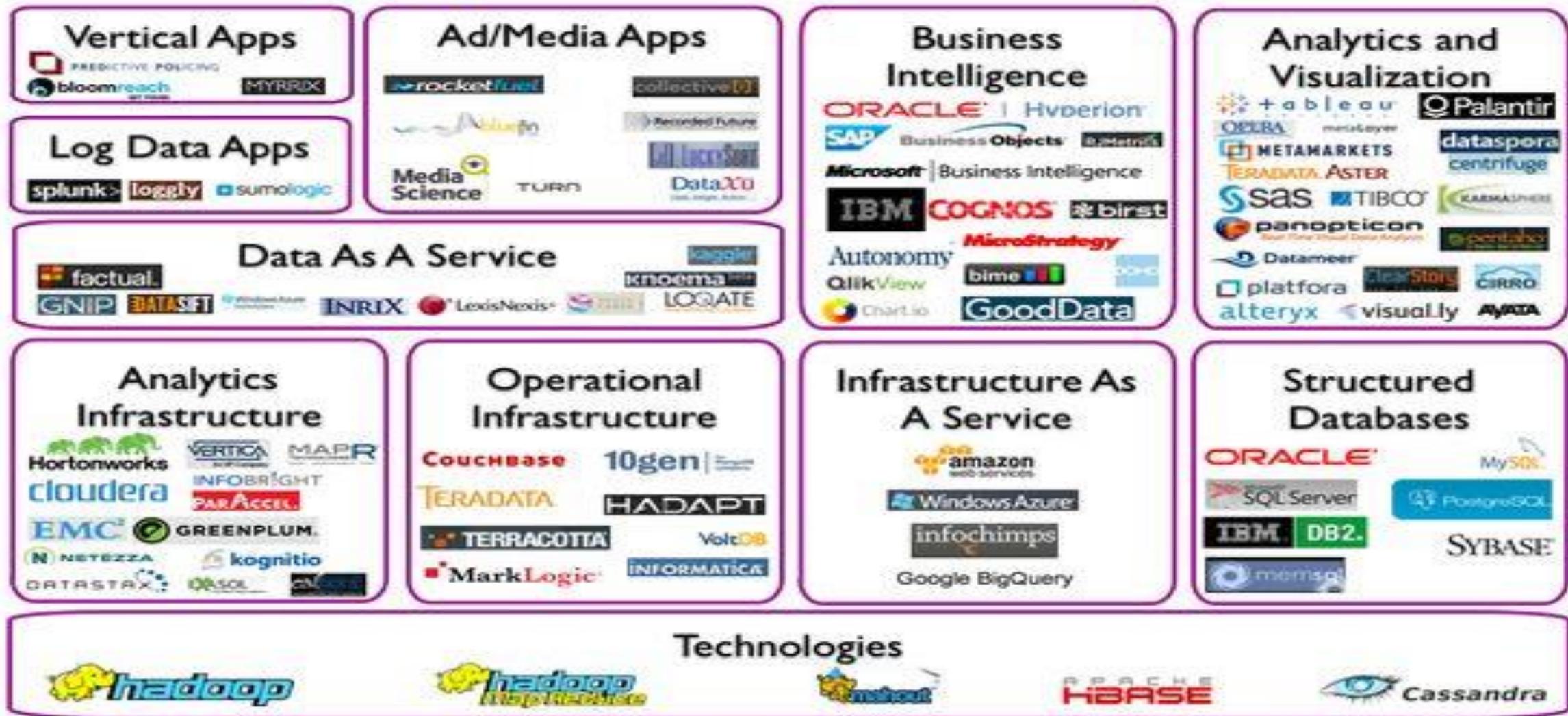


- Not Disruptive Now – Not Scary Anymore
 - Due to cloud solutions
 - If you fail, it is not anymore due to hardware or software limitations
- BDA Project Generic Process:
 - Acquire understanding of business operations
 - Understand the big data business problem
 - Do a benefit analysis – Will BDA be aligned with goals? Will there be actual insights? Will they be actionable? What are short and long term benefits?
 - Measure risk: What's the loss if you didn't do it? - Resolution of Pain Points
 - Determine previous successful use cases
 - Review BDA Implementation Options – Technology stack is cloud now
 - Implement

Big Data in a state of Flux: Evolution of Tools, Technologies and Procedures

However, Orgs Cannot Wait-and-See!

Big Data Landscape



Big Data Landscape (Version 2.0)

Infrastructure



SAP sas IBM Google ORACLE Microsoft VMware amazon iotdata

Framework

Hadoop
Apache
MapReduce
HDFS

Query / Data Flow

Hbase

Batchprocess

Cassandra

mongrelDB

SciDB

Sqoop

Open Source Projects

Coordination / Workflow

ZooKeeper

talend

oozie

Real - Time

Storm

Statistical Tools

SciPy

Machine Learning

TensorFlow

Cloud Deployment

CloudBees

Analytics



METAMARKETS

TRADINGTA

Autonomy

luminate

Applications



Data Sources



Withings

Personal Data

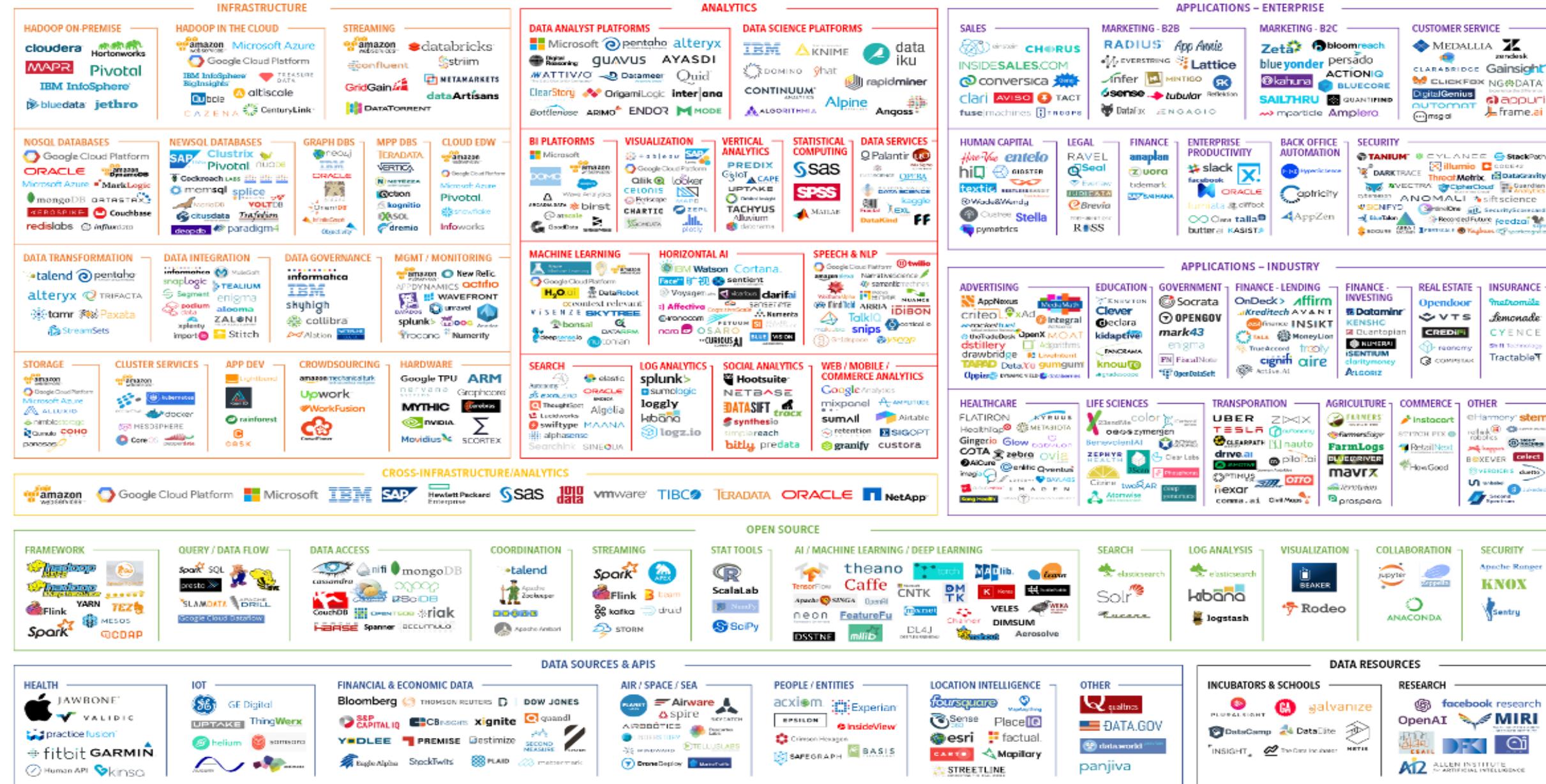
BASIS

JAWBONE

RunKeeper

Nike+ Fitbit

BIG DATA LANDSCAPE 2017

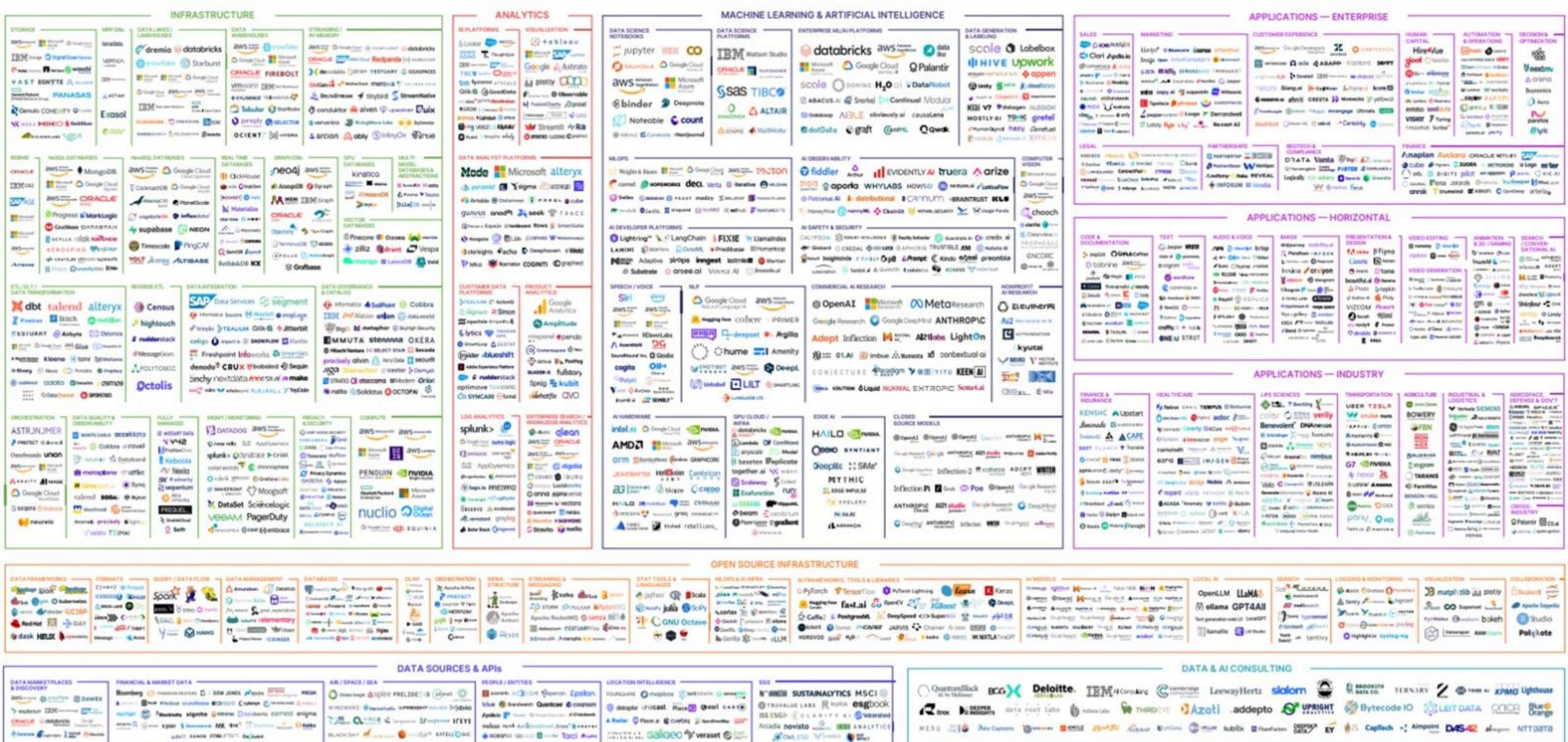


Last updated 4/5/2017

© Matt Turck (@mattturck), Jim Hao (@jimrhao), & FirstMark (@firstmarkcap)

mattturck.com/bigdata2017

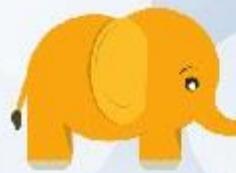
FIRSTMARK
EARLY STAGE VENTURE CAPITAL



Hadoop



- Hadoop has been around some time (Google MapReduce and Big Table, Doug Cutting, Hadoop, Yahoo): – works well for tasks which are deep and computationally extensive and no time restriction (clustering)
- Removes management overhead
- Fault Tolerant
- Commodity Hardware
- Different Configs for Different Hardware – a blend!
- Network is a bottleneck in BDA so Hadoop works where the data resides, hence minimizing communication



Hadoop Ecosystem


oozie
(Work flow)
HCatalog
Table & schema Management 
Pig
(Scripting) **Hive**
(Sql Query) 
Mahout
(Machine Learning) **Drill**
(Interactive Analysis)
AVRO
(JSON) **Thrift**
(Cross Language Service)
APACHE HBASE
HBASE
(Columnar Store)
Sqoop
(Data Collection)
Zookeeper
(Coordination)
Apache Ambari
(Management & Monitoring)
Mapreduce
(Data Processing)
Yarn
(Cluster Resource Management)
HDFS
(Hadoop Distributed File system)



Hadoop Success Stories

- <https://www.forbes.com/sites/bernardmarr/2015/11/16/big-data-success-stories-beyond-hadoop/>
- <https://www.singlestore.com/blog/stampeding-away-from-hadoop-three-customer-success-stories/>
- <https://ellicium.com/success-stories/managed-services-for-hadoop/>
- <https://www.bmc.com/blogs/hadoop-examples/>
- https://community.microstrategy.com/s/article/Customer-Success-Stories-with-the-Hadoop-Gateway?language=en_US

Note!

- Not a plug and play experience
- Skill set, experience, technical requirements, business decisions
- HDFS and HBase for Data management
- MapReduce / Spark / Oozie as a Processing framework
- Hive as Development framework
- Tableau / PowerBI / SiSense / Microstrategy for BI



Skillset Needed

- **Data Science:** Ability to translate the significance of data in a way that can be easily understood by others
- **Machine Learning/Statistics/Linear Algebra:** Important to understand the basic mathematics behind each algorithm
- **Data Cleaning:** Wrangling, Preprocessing, Dredging
- **Data Analysis experience:** SQL-based, DWH-based, BI-based
- **Data Visualization (BI tools and Python-based dashboards)**
- **Data Discovery:** Experience in Extraction of Insights
- Skill in Linux OS
- Skill in Python development using open source technologies



Big Data Techniques

- Parallel Processing, Distributed Computing and High-Performance Computing
- Motherboard, DRAM, Cache, SSD/HDD, GPU)
- In-memory processing, columnar SQL access
- NoSQL Databases
- Containerization (Docker)
- Infrastructures and architectures (lambda, kappa, htap)
- Cloud Computing
- High-Speed Networks
- AI and Data Science



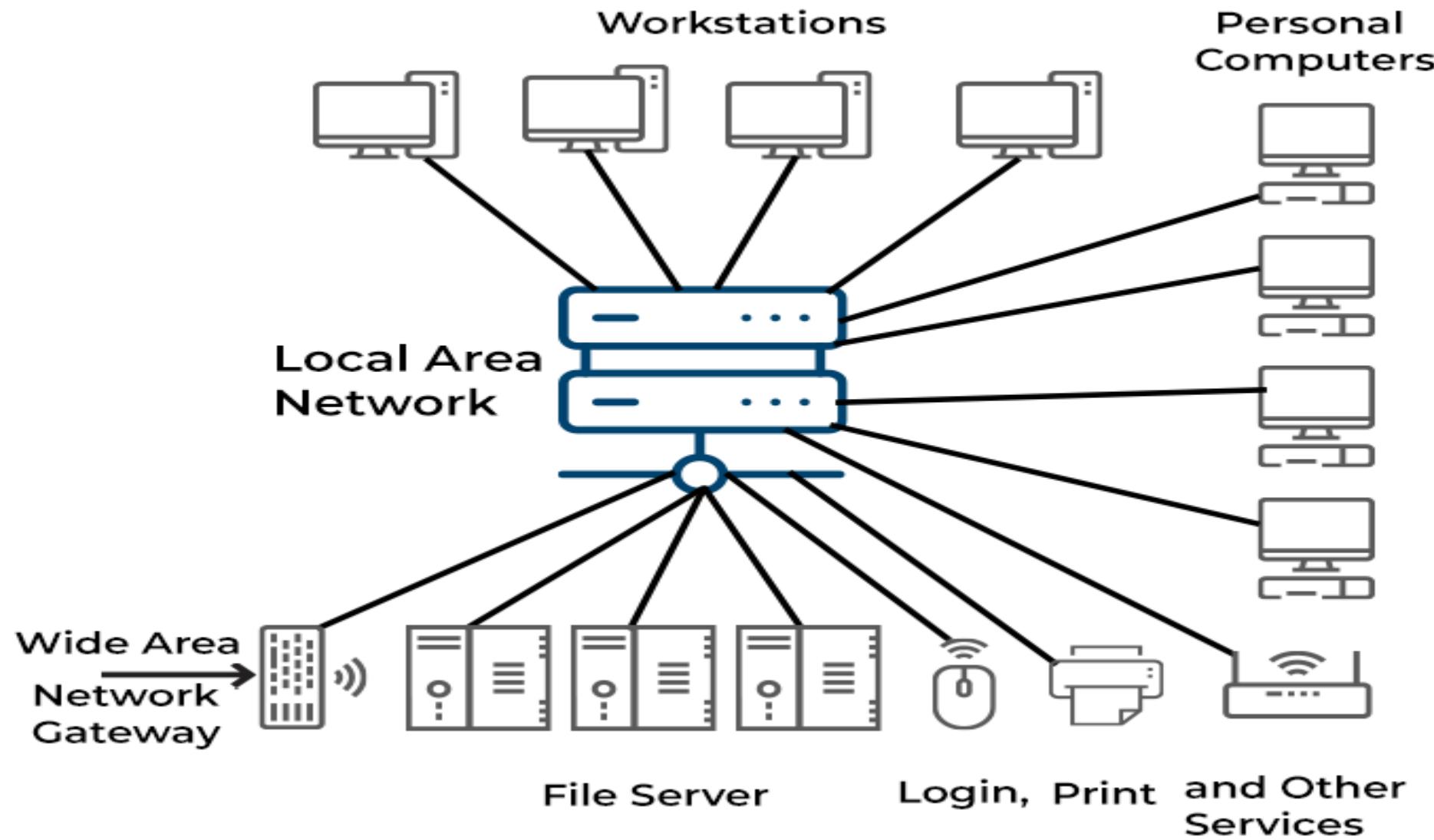
Yes, parallel processing is used in normal Windows operations on a Core i7 processor. Modern operating systems like Windows are designed to take advantage of multi-core processors, such as the Intel Core i7, to improve performance and responsiveness. Here's how parallel processing is utilized:

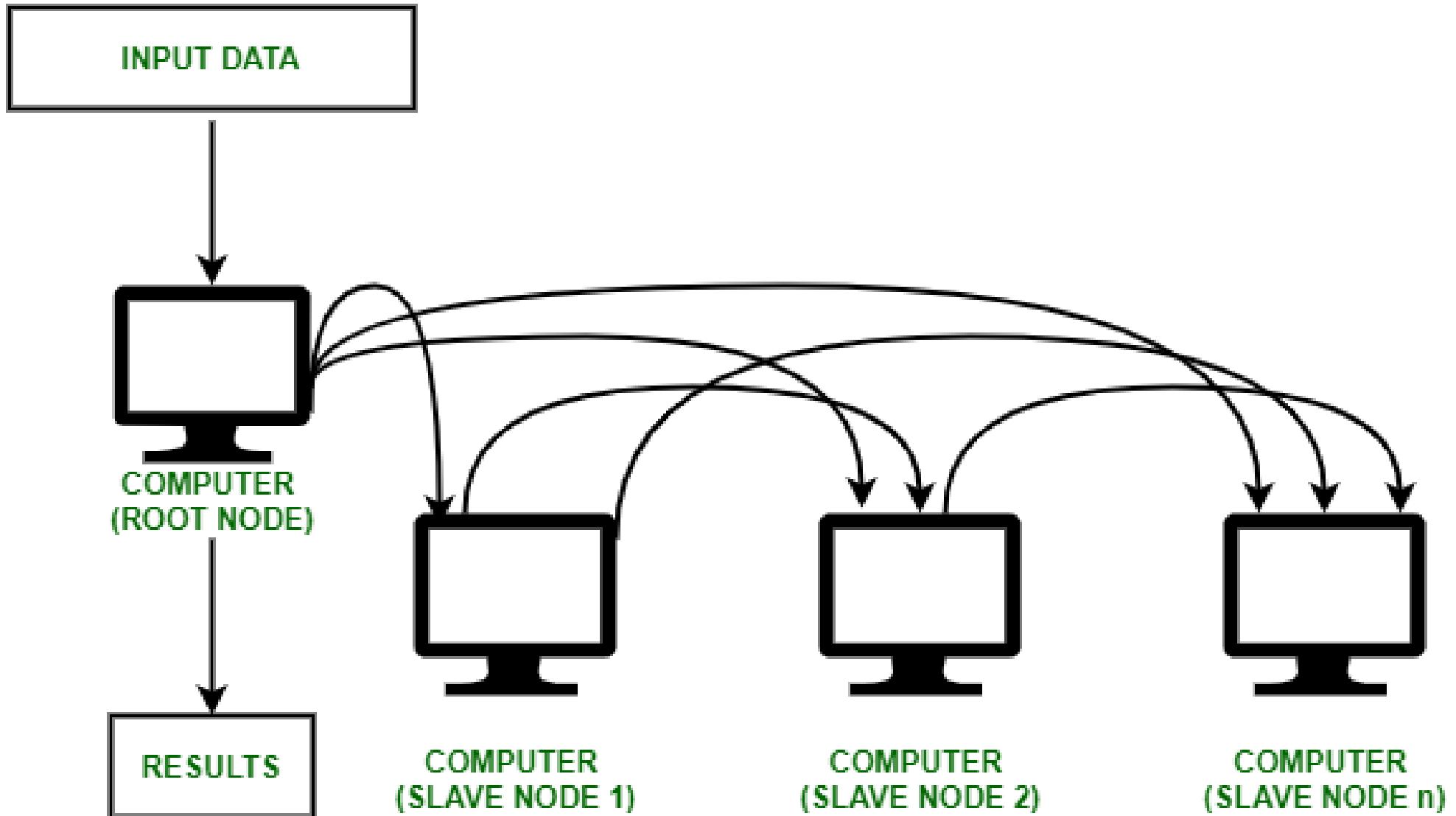
1. **Multithreading:** Many applications, including the Windows operating system itself, are designed to run multiple threads simultaneously. Each thread can be processed by a different core, allowing tasks to be completed more quickly.
2. **Task Scheduling:** The Windows operating system includes a task scheduler that efficiently distributes processes across the available cores. This helps in balancing the load and ensuring that no single core is overwhelmed, leading to smoother performance.
3. **Background Processes:** Windows often runs several background processes concurrently, such as system maintenance tasks, updates, and security scans. Parallel processing allows these tasks to run without significantly impacting the performance of the main tasks or applications you are using.
4. **Application Performance:** Many modern applications, especially those that are resource-intensive like video editing software, games, and scientific computing applications, are optimized for multi-core processors. They can split their workload across multiple cores to run more efficiently.





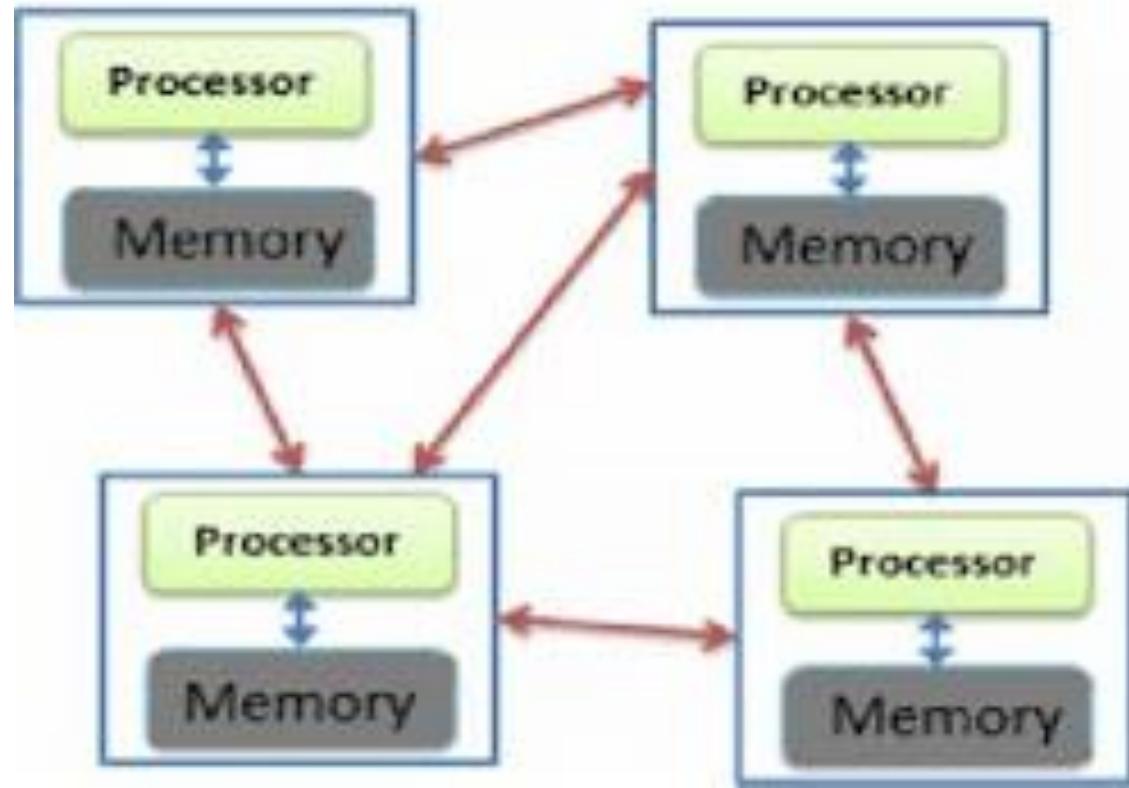
A DISTRIBUTED SYSTEM



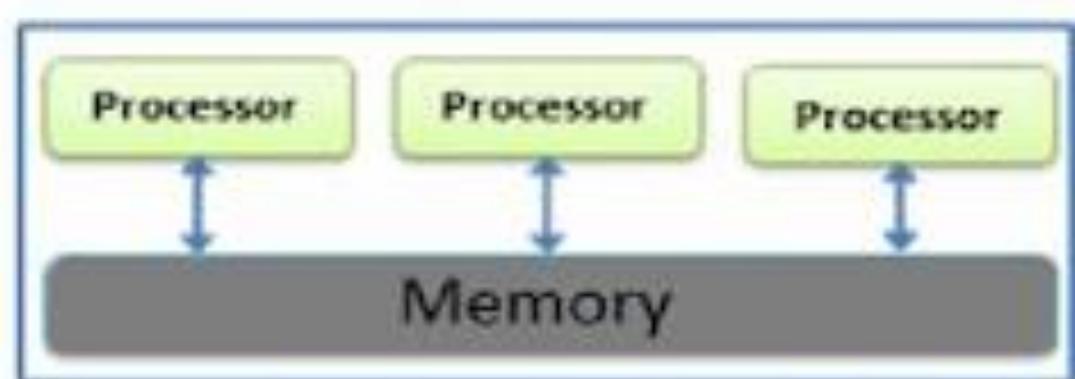




Distributed Computing



Parallel Computing



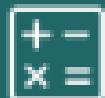
Data Scientist

also known as Data Managers, statisticians.



A data scientist will be able to take data science projects from end to end. They can help store large amounts of data, create predictive modelling processes and present the findings.

Skills: Mathematics, Programming, Communication



Will use programmes such as:

SQL, Python, R

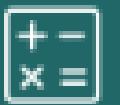
Data Engineers

also known as database administrators and data architects.



They are versatile generalists who use computer science to help process large datasets. They typically focus on coding, cleaning up data sets, and implementing requests that come from data scientists.

Skills: Programming, Mathematics, Big data



Will use programmes such as:

Hadoop, NoSQL, and Python

Data Analysts

also known as business Analysts.



They typically help people from across the company understand specific queries with charts.

Skills: Statistics, Communication, Business knowledge



Will use programmes such as:

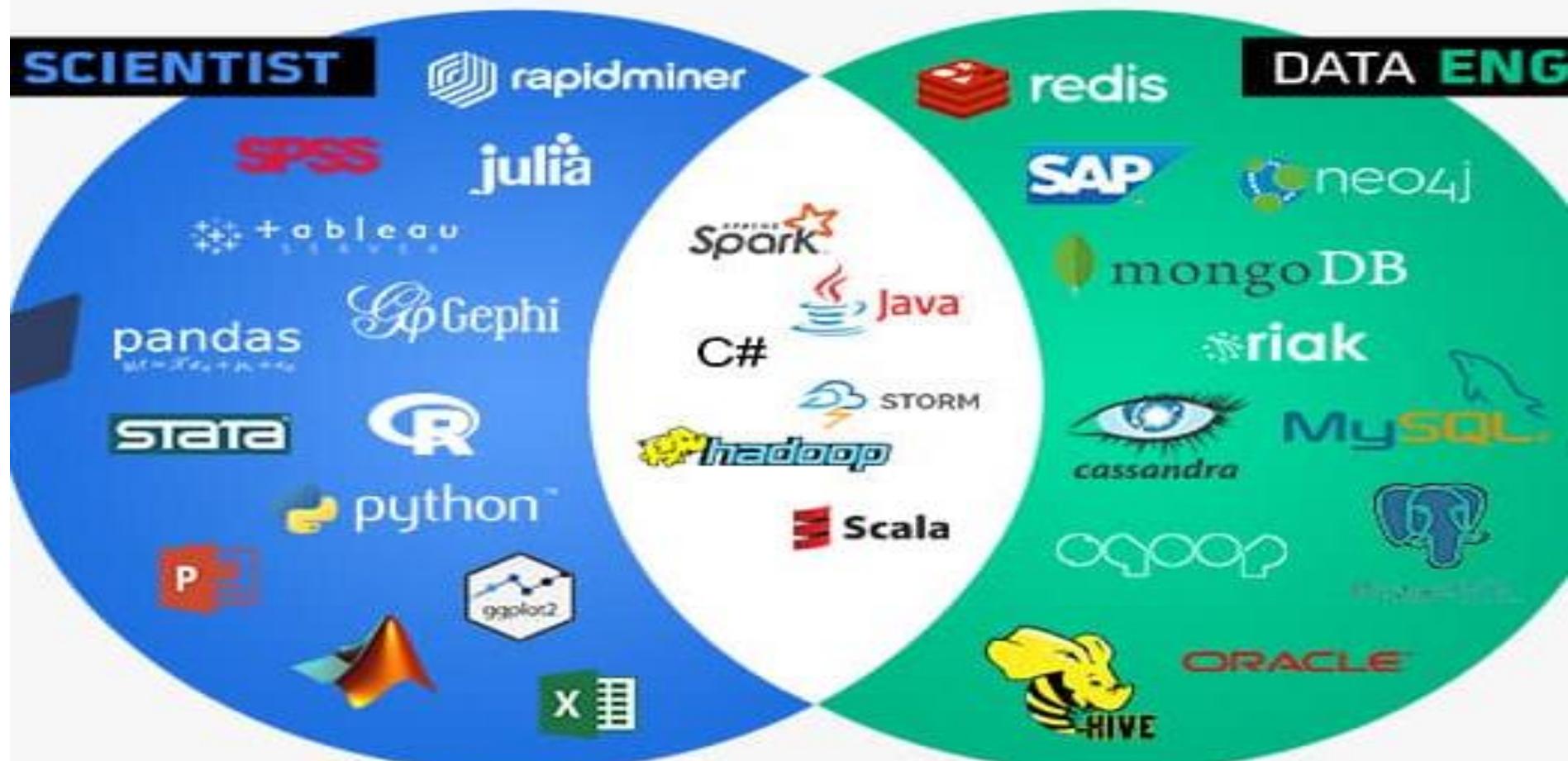
Excel, Tableau, SQL

Data Engineer vs Data Scientist

Data Engineer	Data Scientist
Programming Languages like R, Python	Programming languages like R, Python
Database management skills including SQL	Math, Statistics
Knowledge of big data architecture such as Hadoop, Kafka, Hive, Impala	Knowledge of database queries to access data
Knowledge of security protocols to ensure data is secure.	Machine learning skills.
Cloud computing skills, e.g., Amazon web services.	Communication skills
Data Visualization skills.	Data Visualization skills.
Basics of machine learning added benefit.	NLP, AI added benefit.



LANGUAGES, TOOLS AND SOFTWARE

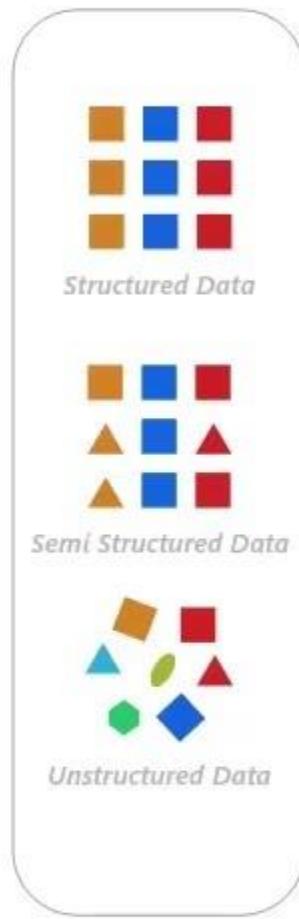


Big Data Analytics (BDA)

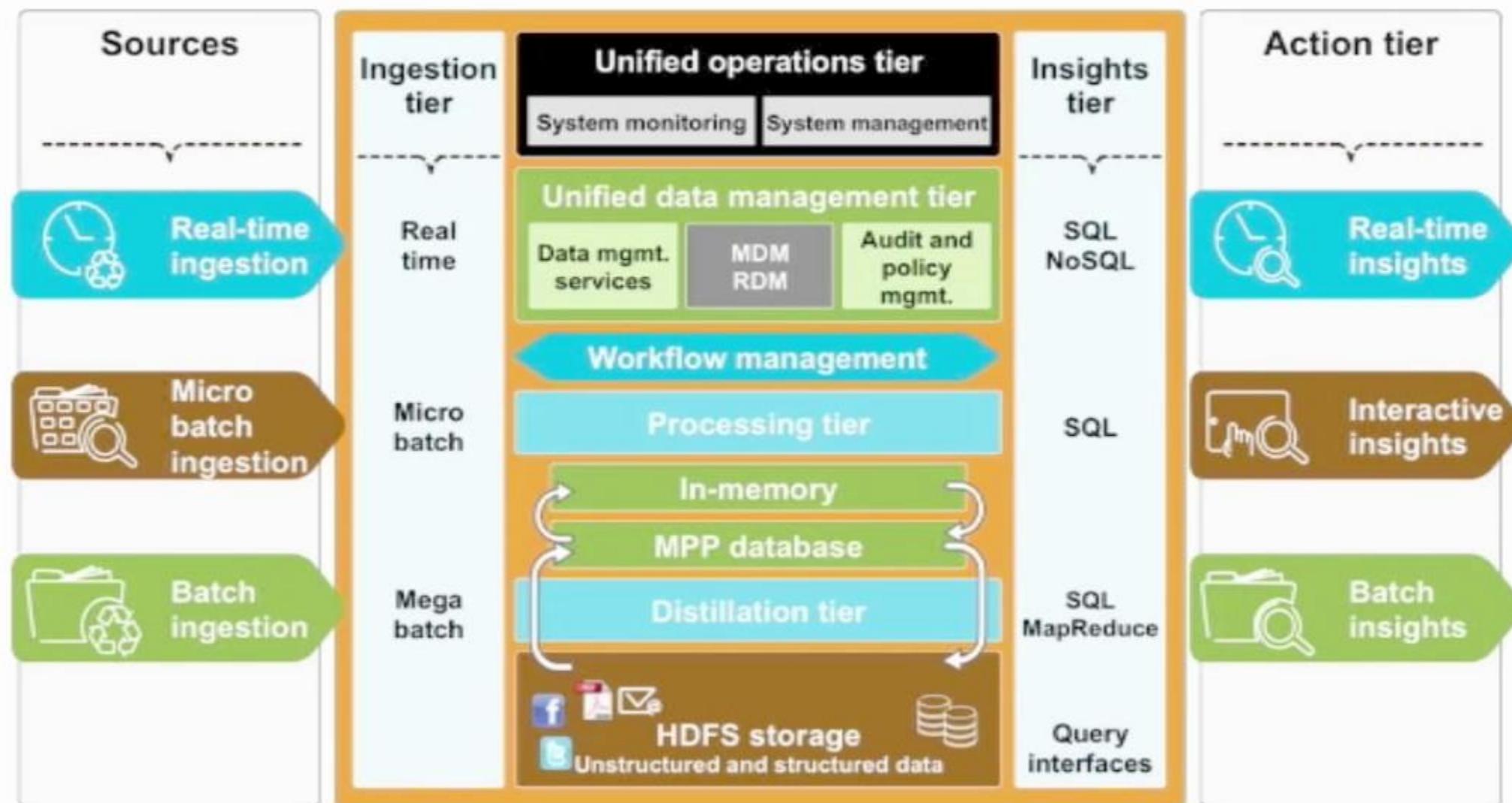
- Big Data Analytics \equiv Derive Value \equiv Strike Gold!
- Learns models which give you information that is hidden, non-trivial, critical, previously unknown and potentially useful (valuable)
- Benefits: Drive business, formulate business strategies, and maintain competitive edge.

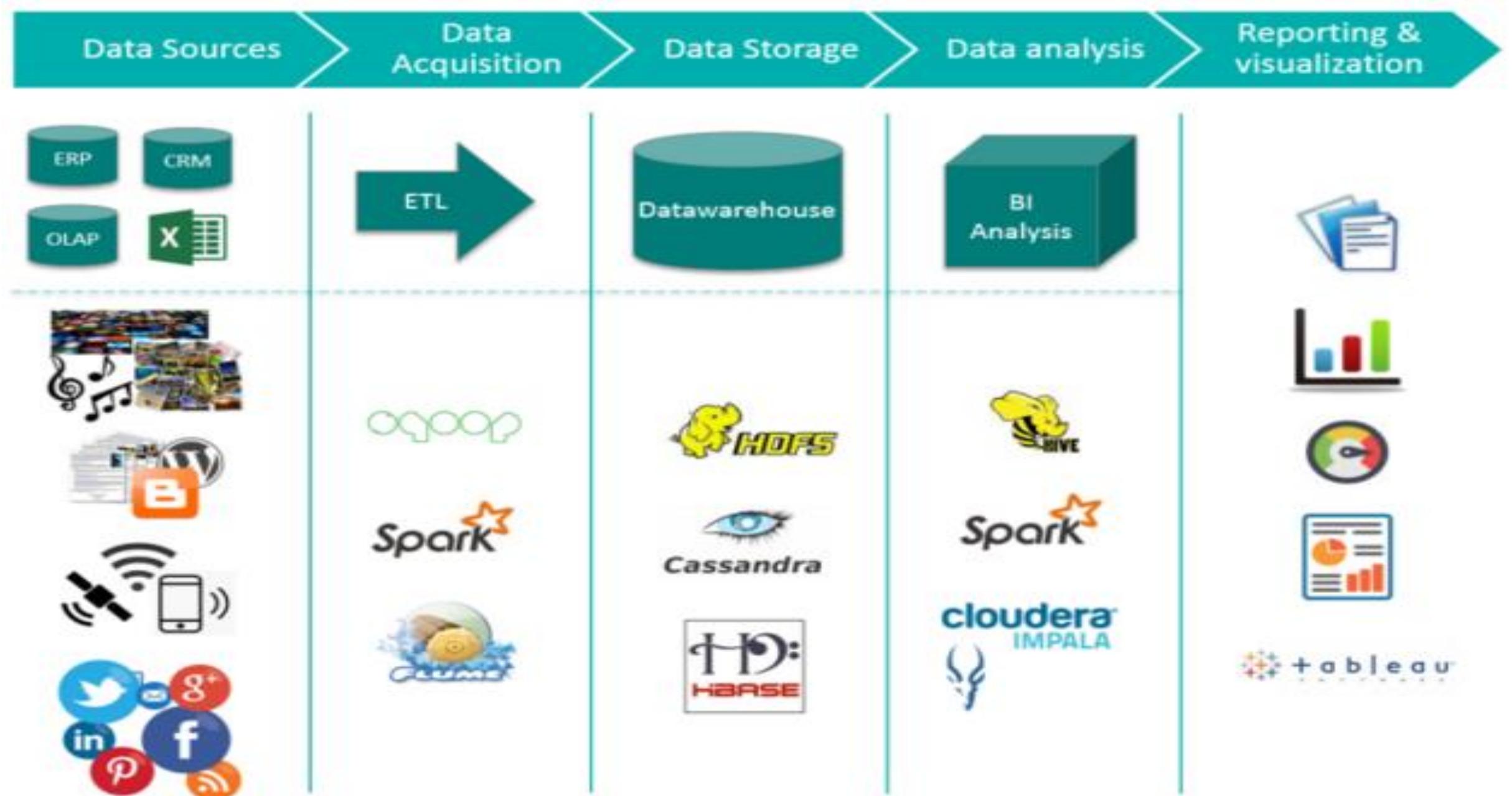


Big Data Workflow



Data Lake





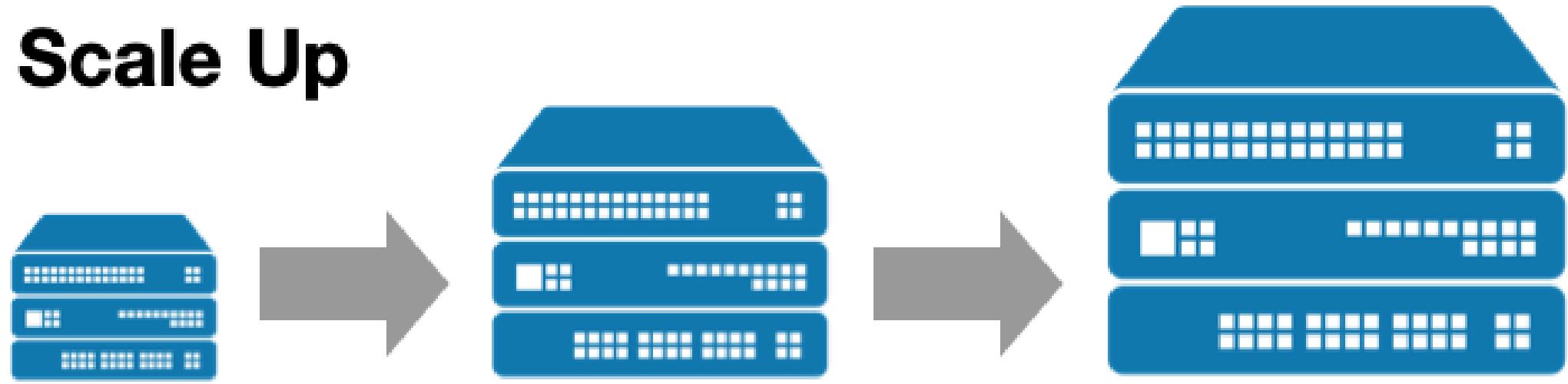
BDA Situation

Era	Tools	Data Focus	Output	Role of AI/LLMs
Classic BDA (2010s)	Hadoop, Hive, Spark	Structured/semi-structured	Batch reports, dashboards	Minimal
Modern BDA (2020s)	Cloud DW (Snowflake, BigQuery), Spark Streaming	Structured + some unstructured	Real-time dashboards, ML predictions	ML models integrated
BDA + LLMs (2023–2025)	LLMs, Vector DBs, RAG, AI copilots	Structured + unstructured (text, image, audio)	Conversational insights, predictive + prescriptive analytics	LLMs generate queries, automate wrangling, summarize, recommend

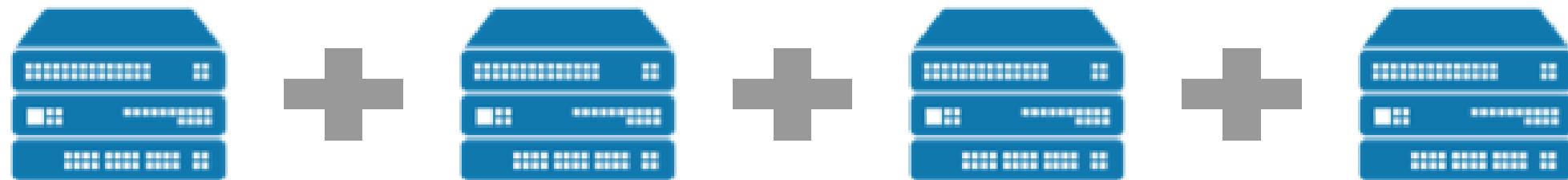
BDA Situation

Era	Architecture	Tools & Tech	Data Focus	Output	Role of AI/LLMs
Classic BDA (2010s)	Data Lakes (HDFS-based) → centralized storage of raw data	Hadoop, Hive, Spark, Pig	Mostly structured/semi-structured (logs, tables, JSON)	Batch analytics, ETL to warehouse, static dashboards	Minimal → manual SQL/HQL
Modern BDA (2020s)	Cloud Data Lakes + Early Data Fabric (integration across sources)	Spark Streaming, Snowflake, BigQuery, Lakehouse (Databricks Delta Lake)	Structured + some unstructured (images, IoT streams)	Real-time dashboards, ML predictions	ML integrated into pipelines, some automation
BDA + LLMs (2023–2025)	Intelligent Data Fabric over Data Lakes (AI-driven metadata mgmt, knowledge graphs, vector DBs)	LLMs, RAG, Vector DBs (Pinecone, FAISS, Milvus), Databricks MosaicML, Snowflake Cortex	All data types: structured + unstructured (text, docs, video, audio, sensor data)	Conversational insights, predictive + prescriptive analytics, autonomous actions	LLMs power natural language queries, automate wrangling, augment decision-making, bridge data silos

Scale Up



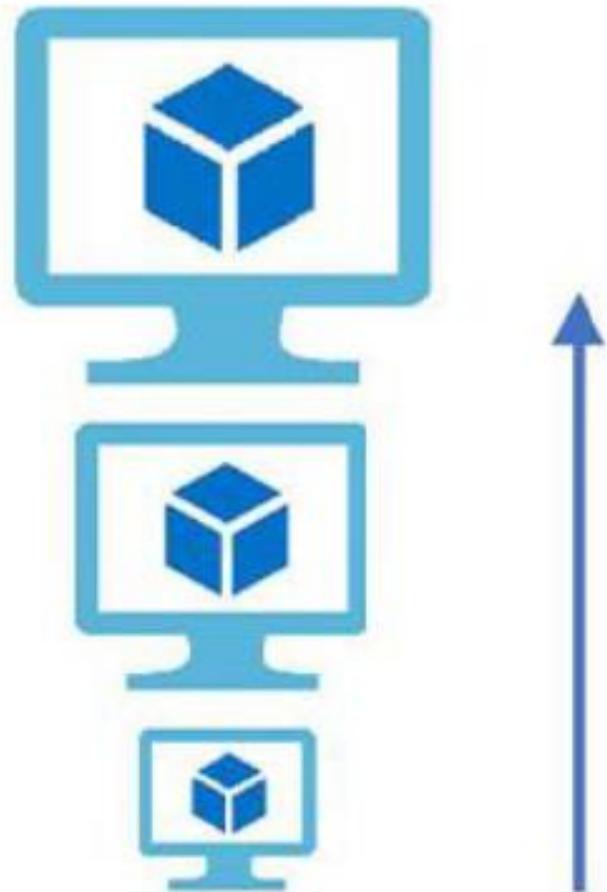
Scale Out





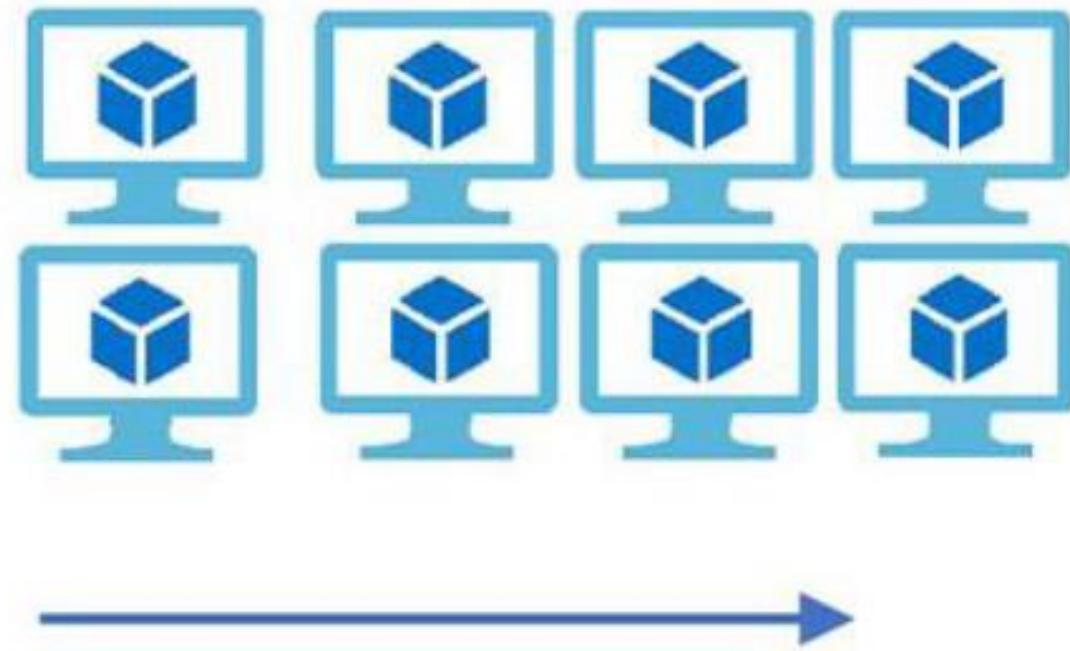
Vertical Scaling

(Increase size of instance (RAM ,
CPU etc.))



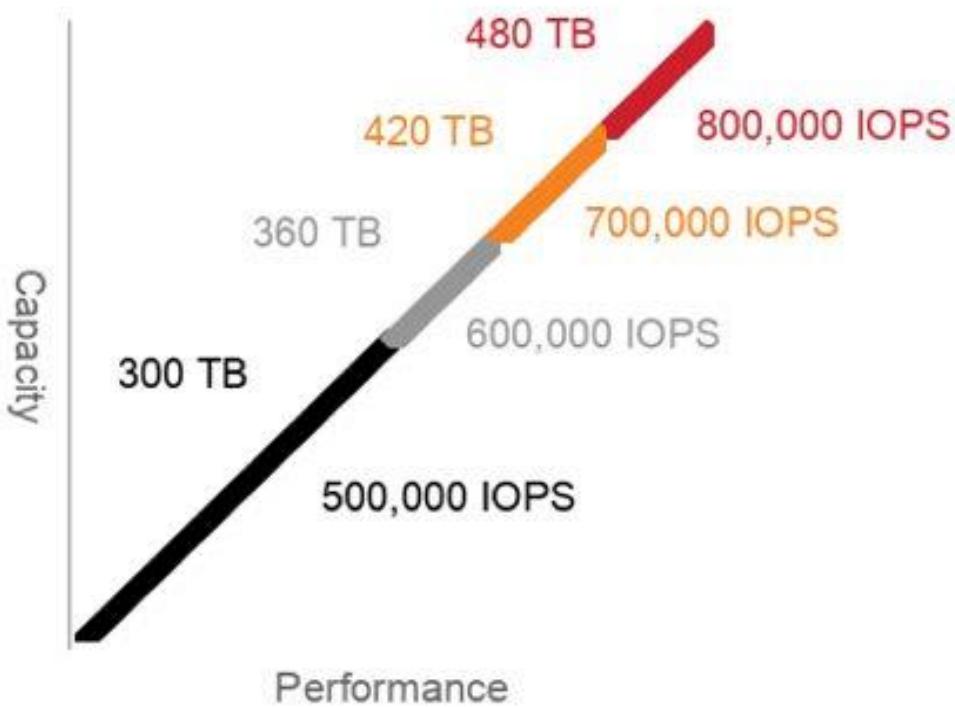
Horizontal Scaling

(Add more instances)

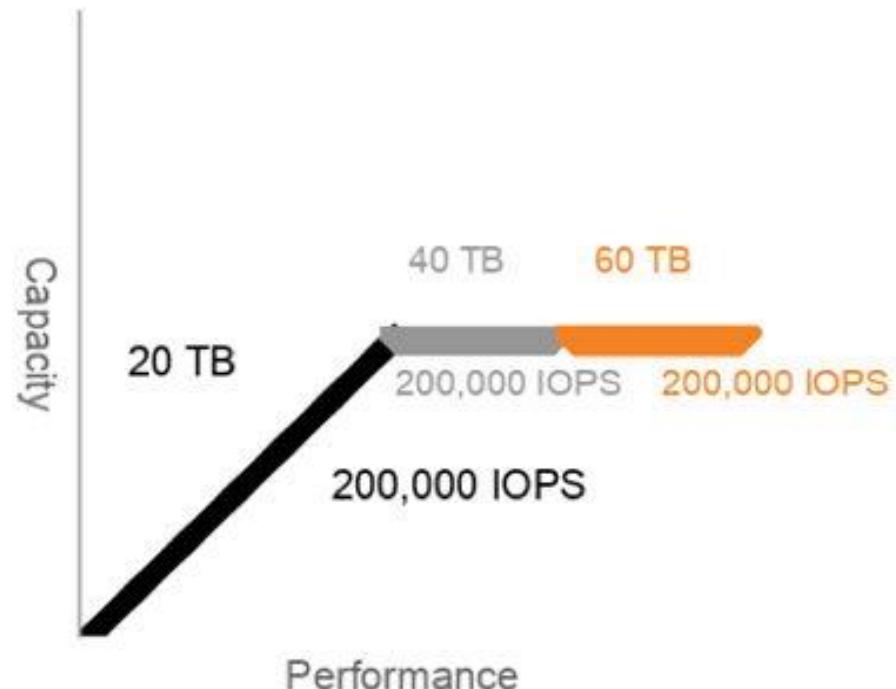


Scalability – Scale Up and Scale Out

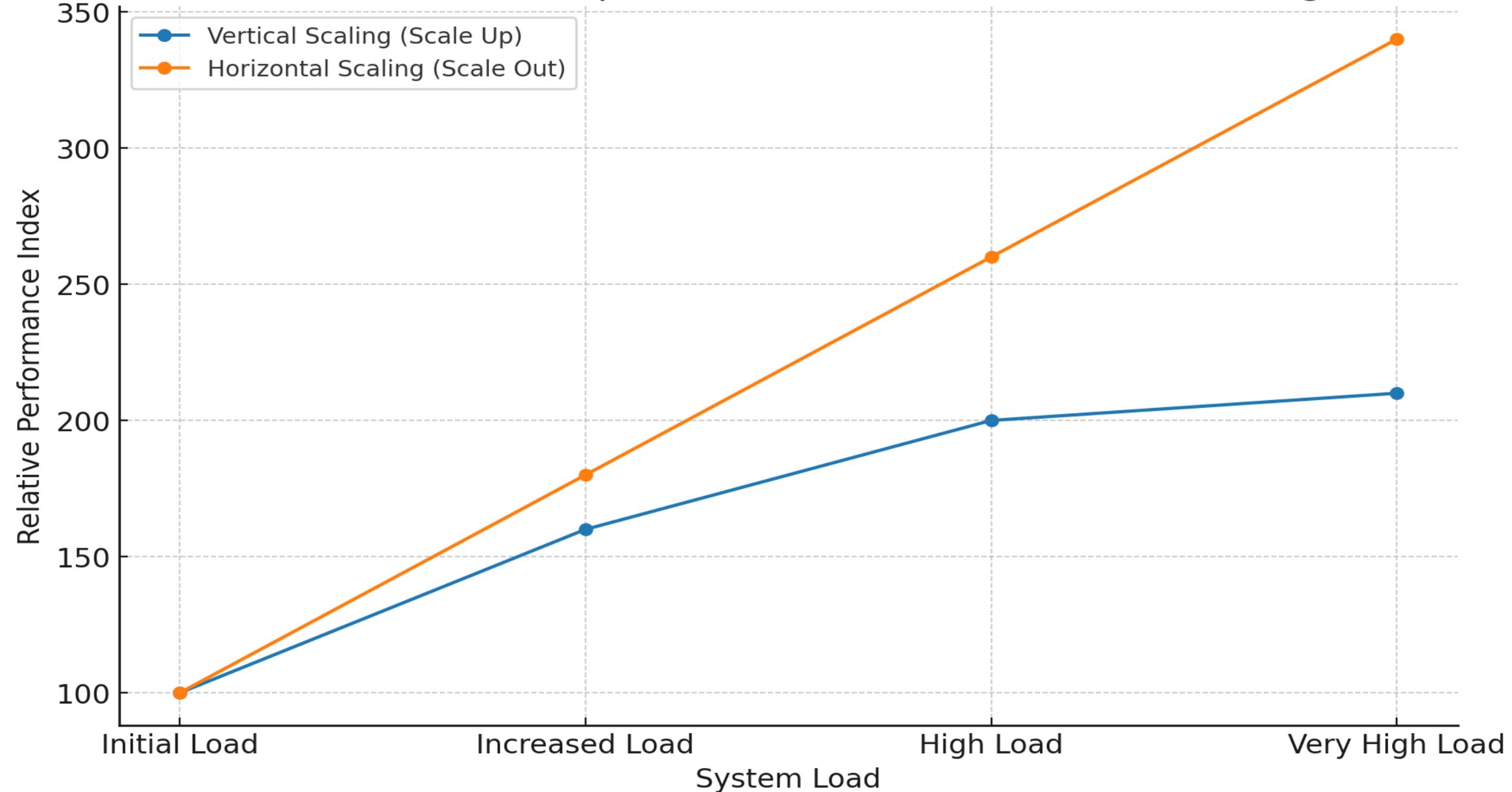
Scale-Out



Scale-Up



Performance Comparison: Vertical vs Horizontal Scaling



In-Memory Spark

Speed

Run programs up to 100x faster than Hadoop MapReduce in memory, or 10x faster on disk.

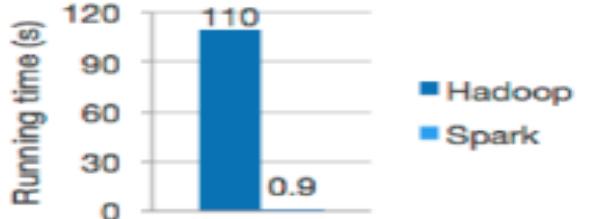
Spark has an advanced DAG execution engine that supports cyclic data flow and in-memory computing.



Ease of Use

Write applications quickly in Java, Scala or Python.

Spark offers over 80 high-level operators that make it easy to build parallel apps. And you can use it *interactively* from the Scala and Python shells.



Logistic regression in Hadoop and Spark

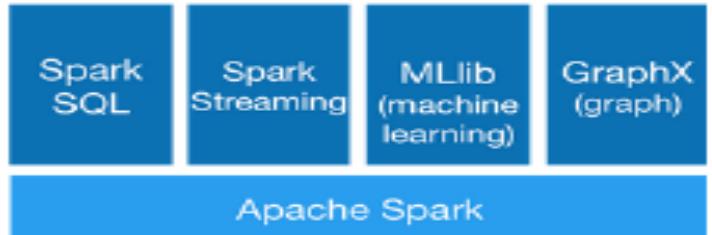
```
file = spark.textFile("hdfs://...")  
file.flatMap(lambda line: line.split())  
.map(lambda word: (word, 1))  
.reduceByKey(lambda a, b: a+b)
```

Word count in Spark's Python API

Generality

Combine SQL, streaming, and complex analytics.

Spark powers a stack of high-level tools including [Spark SQL](#), [MLlib](#) for machine learning, [GraphX](#), and [Spark Streaming](#). You can combine these frameworks seamlessly in the same application.

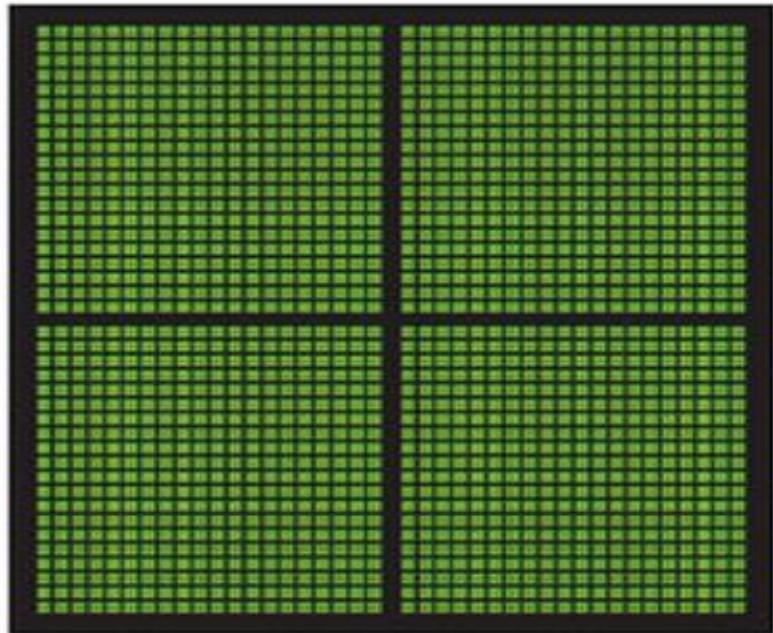
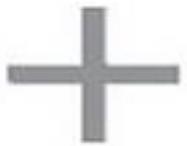


GPUs vs CPUs

GPUs have thousands of cores to process parallel workloads efficiently



CPU
MULTIPLE CORES



GPU
THOUSANDS OF CORES



GPUs vs CPUs

CPU:

+-----+

| Core 1 |

+-----+

| Core 2 |

+-----+

| Core 3 |

+-----+

| Core 4 |

+-----+

(few, powerful)

GPU:

+-----+-----+

| Core 1 | Core 2 | Core 3 | ...

+-----+-----+

| Core 4 | Core 5 | Core 6 | ...

+-----+-----+

| ... | ... | ... |

+-----+-----+

| ... thousands of simple cores ...

+-----+

(many, parallel)

GPUs vs CPUs

Feature	CPU (Central Processing Unit)	GPU (Graphics Processing Unit)
Core Count	Few (4–32 cores in most systems)	Hundreds to thousands of cores
Core Type	Powerful, complex, general-purpose	Smaller, simpler, specialized for parallel
Best At	Sequential tasks, logic-heavy computation	Parallel tasks, matrix/vector operations
Throughput	Lower total throughput, high per-core power	Very high throughput, low per-core power
Latency	Low latency per task	Higher latency per task
Workload Examples	OS tasks, databases, sequential apps	AI/ML training, graphics, simulations
Memory	Large cache, fast access	High bandwidth VRAM (GDDR/HBM)
Flexibility	Runs a wide variety of programs	Optimized for parallel workloads



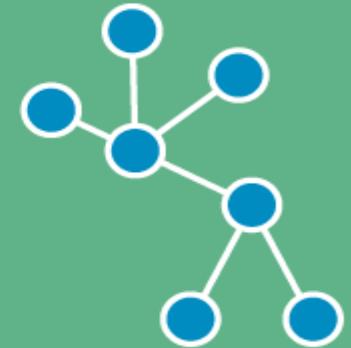
NoSQL Databases

- Stands for **Not Only SQL**
- Class of non-relational data storage systems
- Usually do not require a fixed table schema nor do they use the concept of joins
- All NoSQL offerings relax one or more of the ACID properties
- CAP theorem
- [Document](#), [Key-Value Pairs](#), [Columnar Stores](#), [Graphs](#)
- [NewSQL databases \(NoSQL+ACID\)](#)

Key-Value



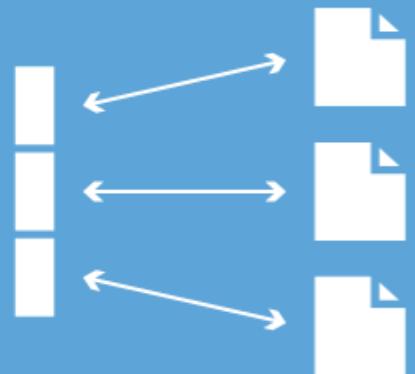
Graph DB

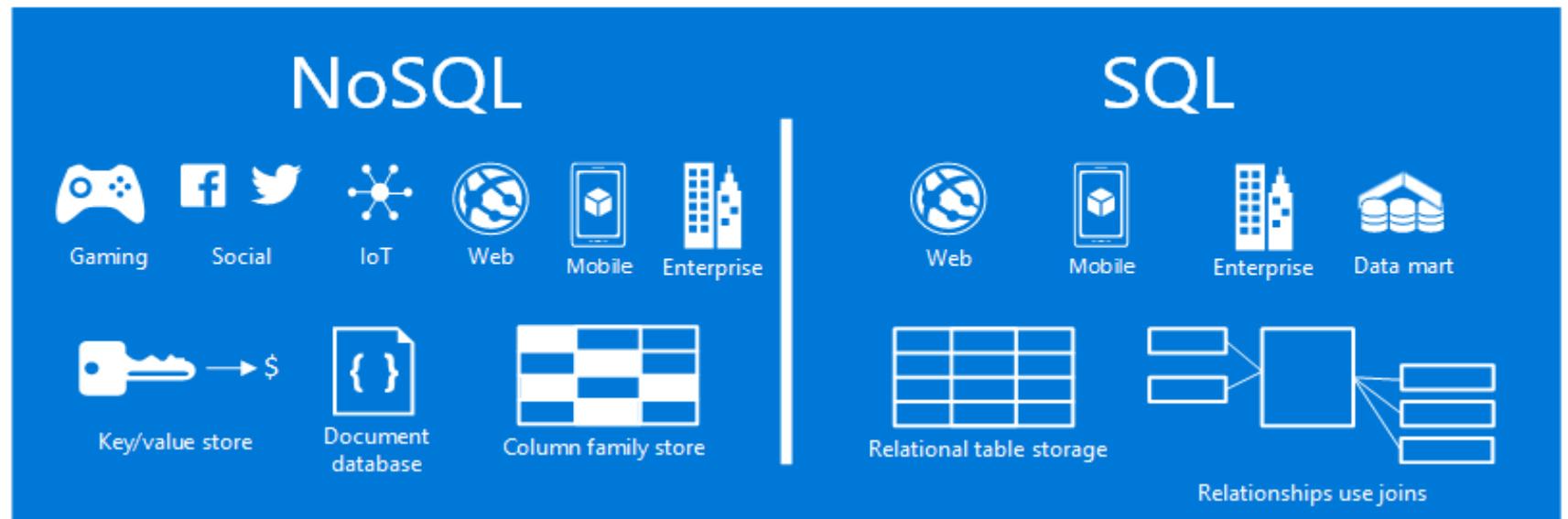


Column Family

1				1
	1		1	1
				1
1				1
1				1
	1			1
				1

Document







Big Data Problems

- Where to store Big Data?
- When to store Big Data?
- How to store Big Data?
- How to process Big Data?
- How to analyze Big Data?
- Organization, Searching and Meta-Data
- How to manage access?
- How to copy, move and backup?
- Provenance – meta data



BIG DATA USAGE BY TECHNOLOGY



Big Data as a Service	Traditional Big Data	Traditional Database
Scalability on demand through a combination of cloud computing and distributed architecture	Scalability in processing and storage achieved through distributed architecture	Lack of resources such as computational power and storage capacity.
Virtualized data storage on a distributed platform.	Data storage on HDFS or distributed platform	Integrated hard data storage such as NAS, SAN, and traditional disks
Structured and unstructured data on cloud environment	Structured and unstructured data	Structured data
Advanced analytics functions with on-demand computing power	Advanced analytics functions	Reporting using tools such as OLAP
Ubiquitous accessibility	Limited accessibility	Limited accessibility
Analytical capability derived from out-of-box domain-specific algorithms along with custom coding	Analytical capability derived through custom coding	Analytical capability derived through custom coding

DIFFERENCE BETWEEN BIG DATA AND INTERNET OF THINGS

IoT Analytics 3-Tier Architecture

using the example of a wind turbine farm

TIER 1

Performs individual wind turbine real-time performance analysis and optimization



100 99 99 99
100 99 100 99
99 100 81 82
100 99 99 98
100 99 99 99

Purpose-built T1 edge gateways leverage real-time data compression techniques to only send a subset of the critical data

TIER 2

Optimizes performance and predicts maintenance needs across the wind turbines in the same wind farm



TIER 3

The data-lake-enabled core analytics platform that aggregates the critical data across all wind farms and individual turbines, and combines the sensor data with external data sources



To drive meaningful business impact, you will need to begin with the business and not the technology:

Engage the business stakeholders on day one

Align the business and IT teams

Understand the organization's key business and operational initiatives

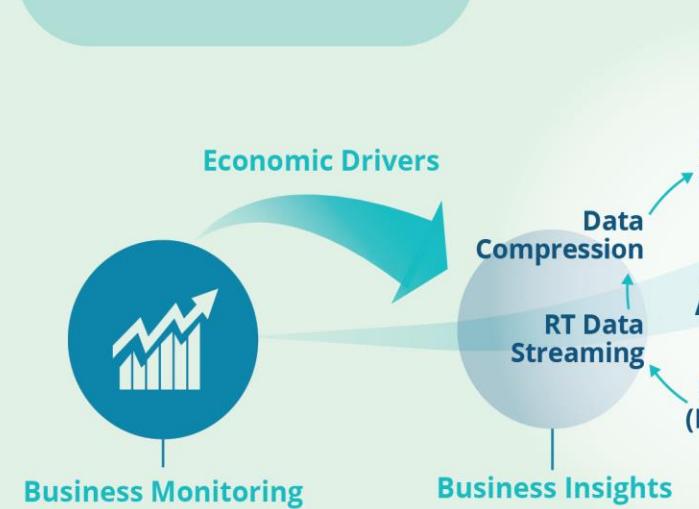
Identify and prioritize the use cases (decisions/goals) that support those business initiatives

Economic Drivers



Business Monitoring

BIG DATA BUSINESS MODEL MATURITY INDEX



“
IoT is more than just another data source. IoT represents the ability to take actions at the point of data capture; to apply machine learning at the source of the data in order to optimize operational decisions. That's the power of the edge.

Producers

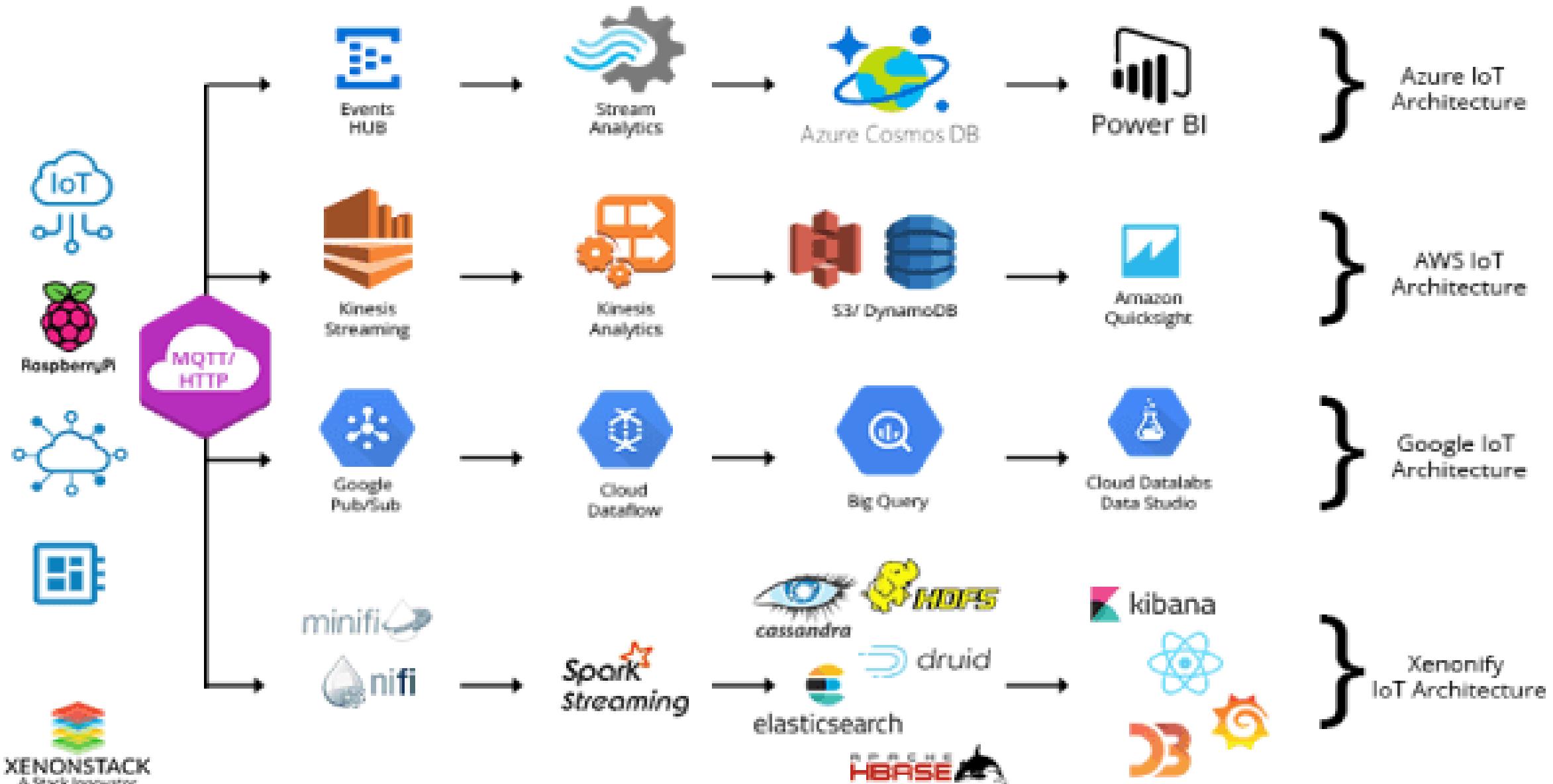
Ingestion

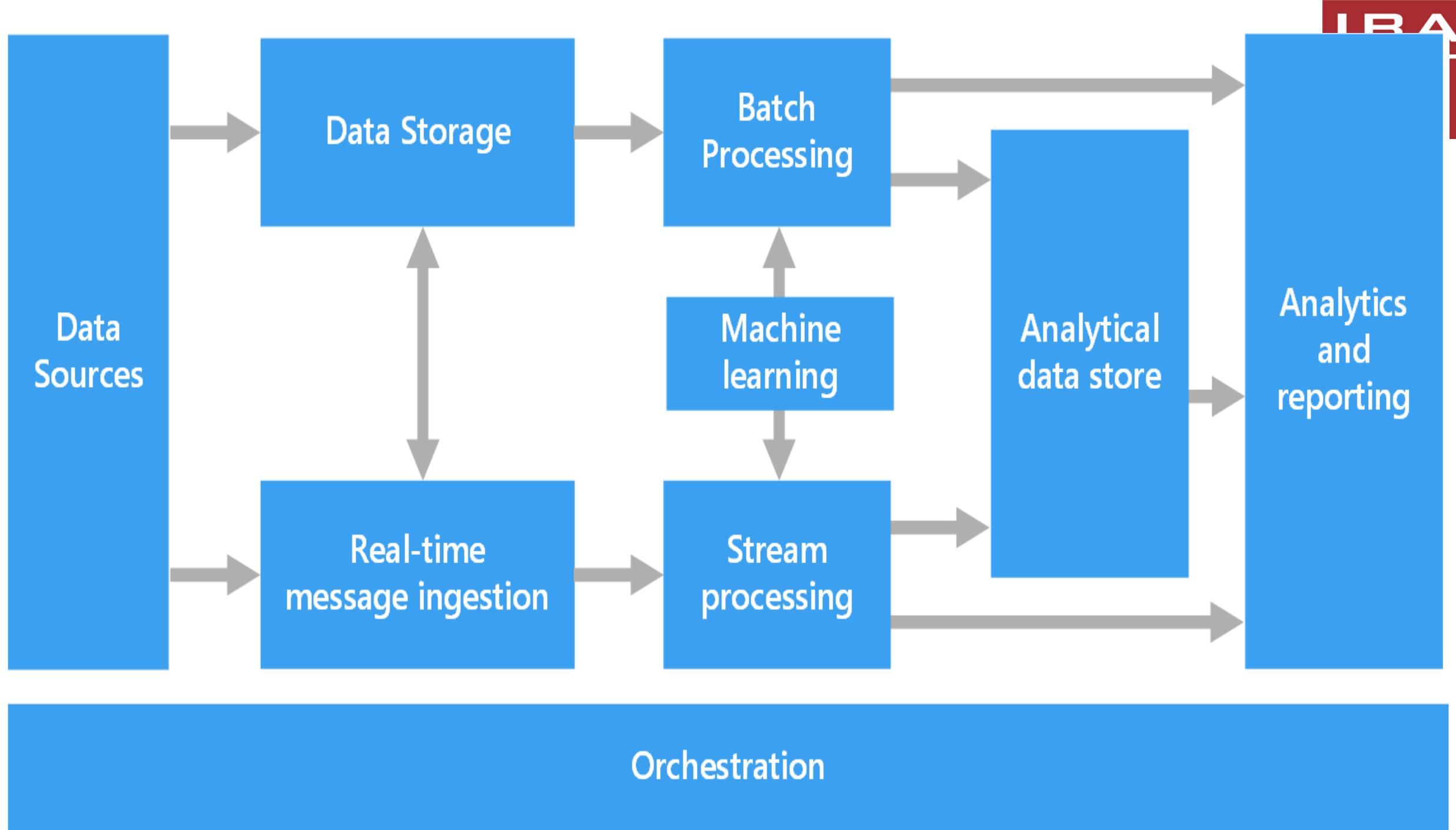
Stream Processing

Storage

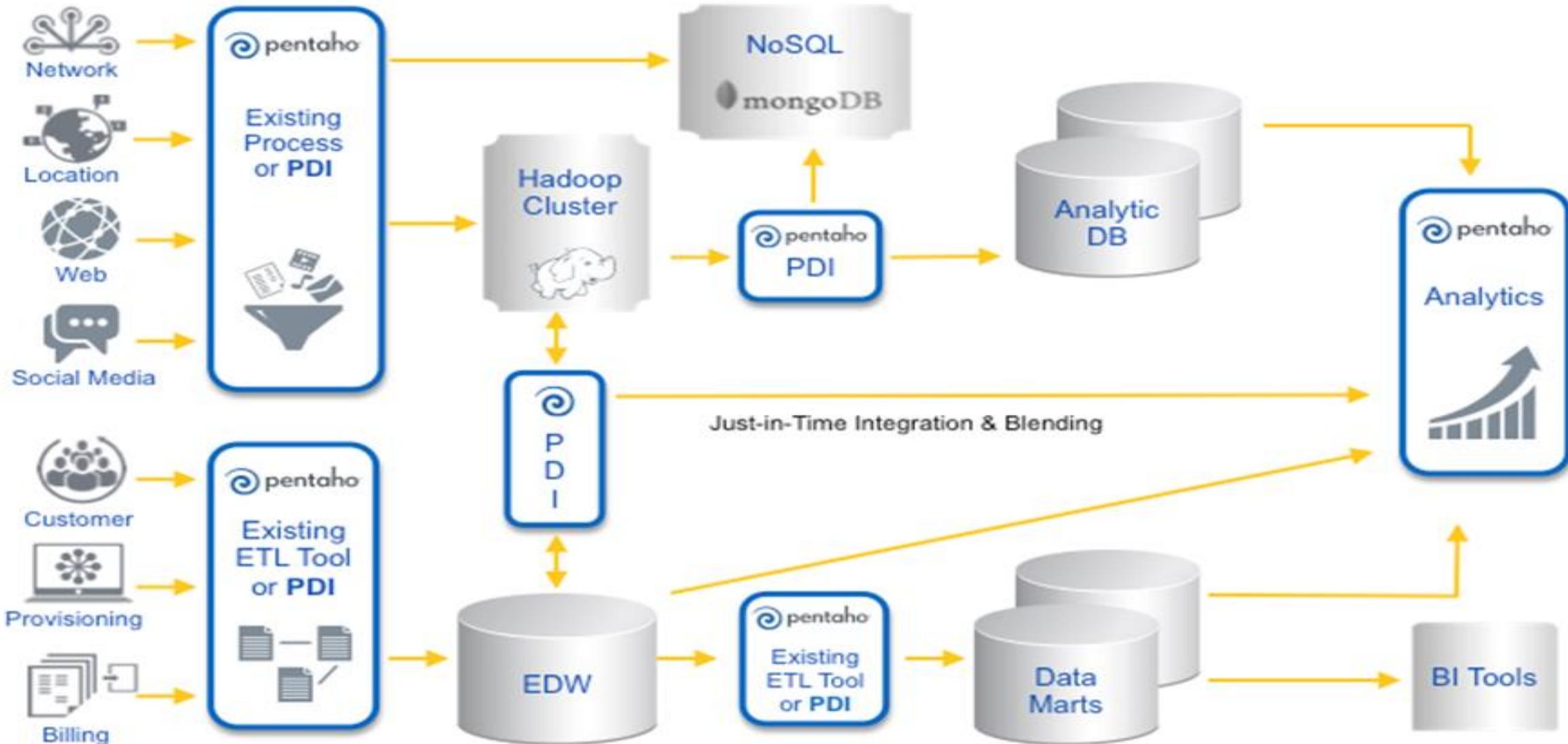
Presentation

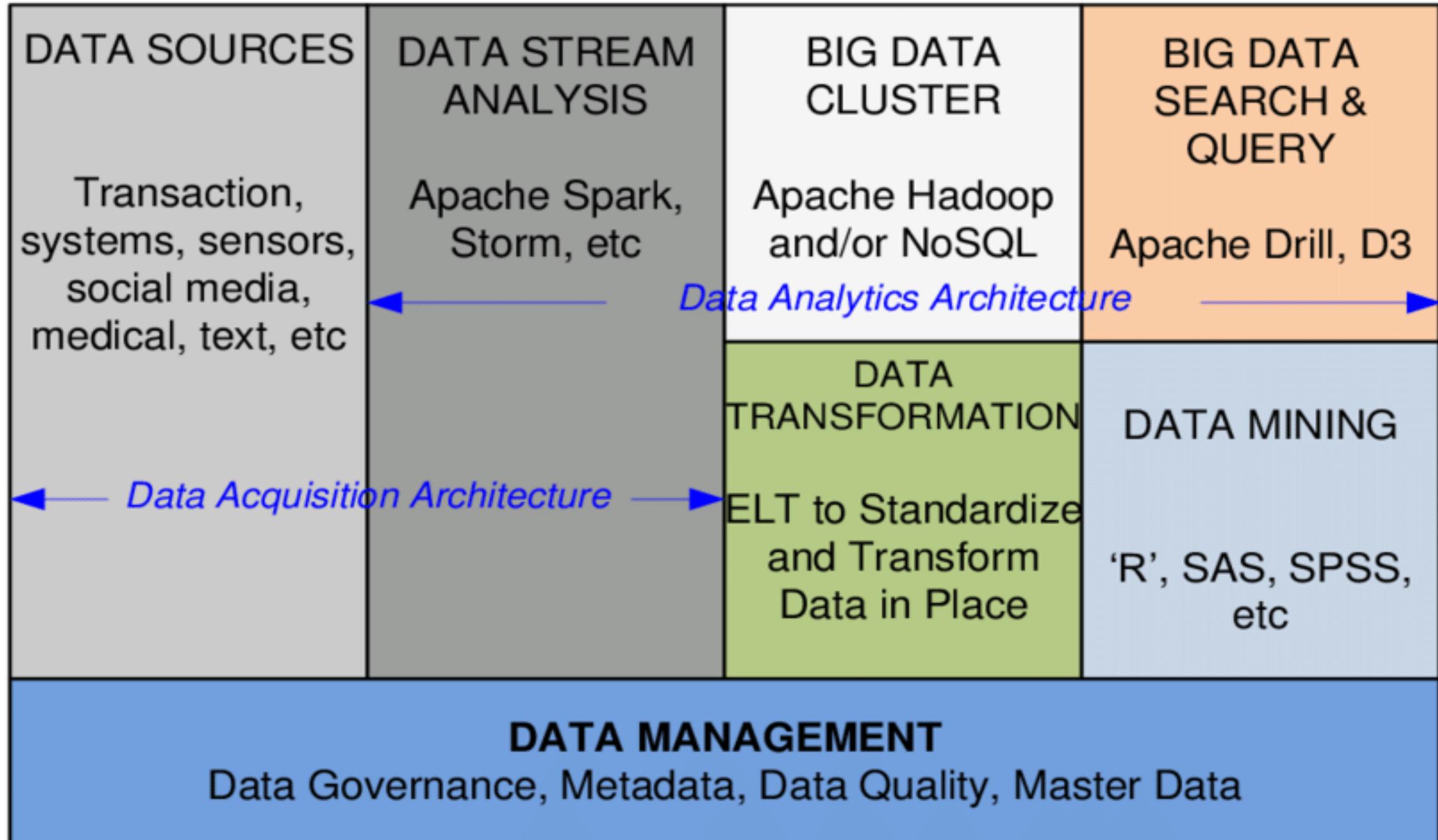
Architectures





Evolving Big Data Architectures

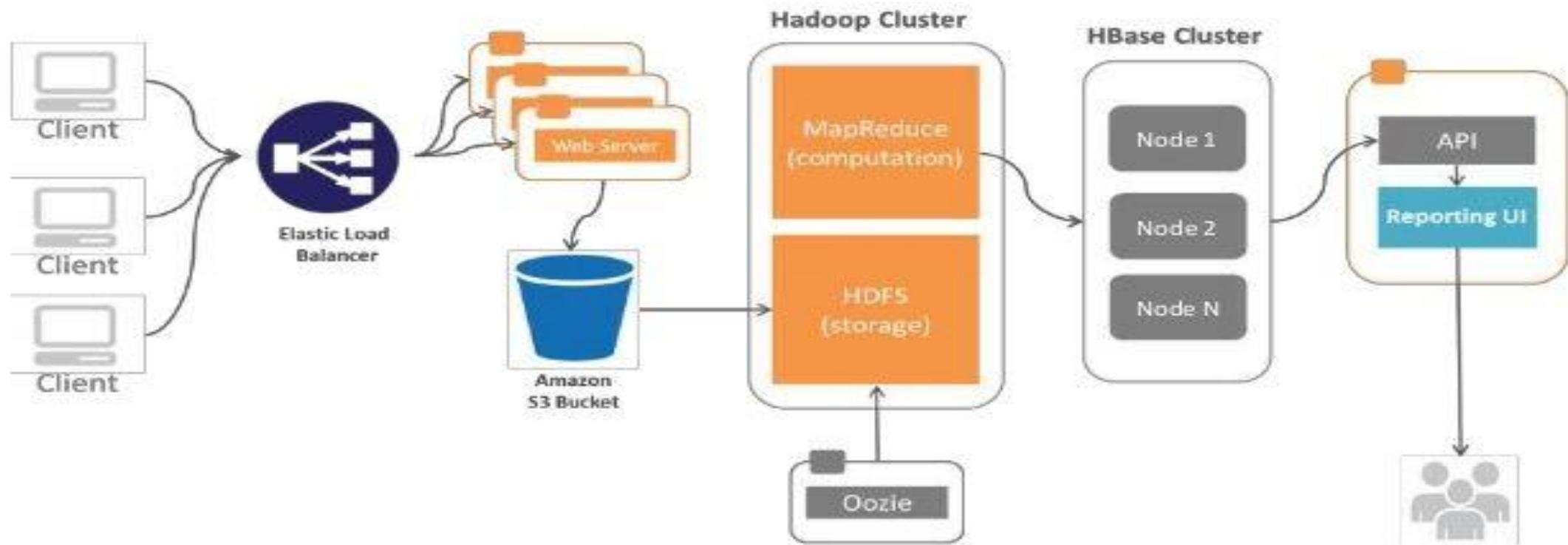




Solution Architecture

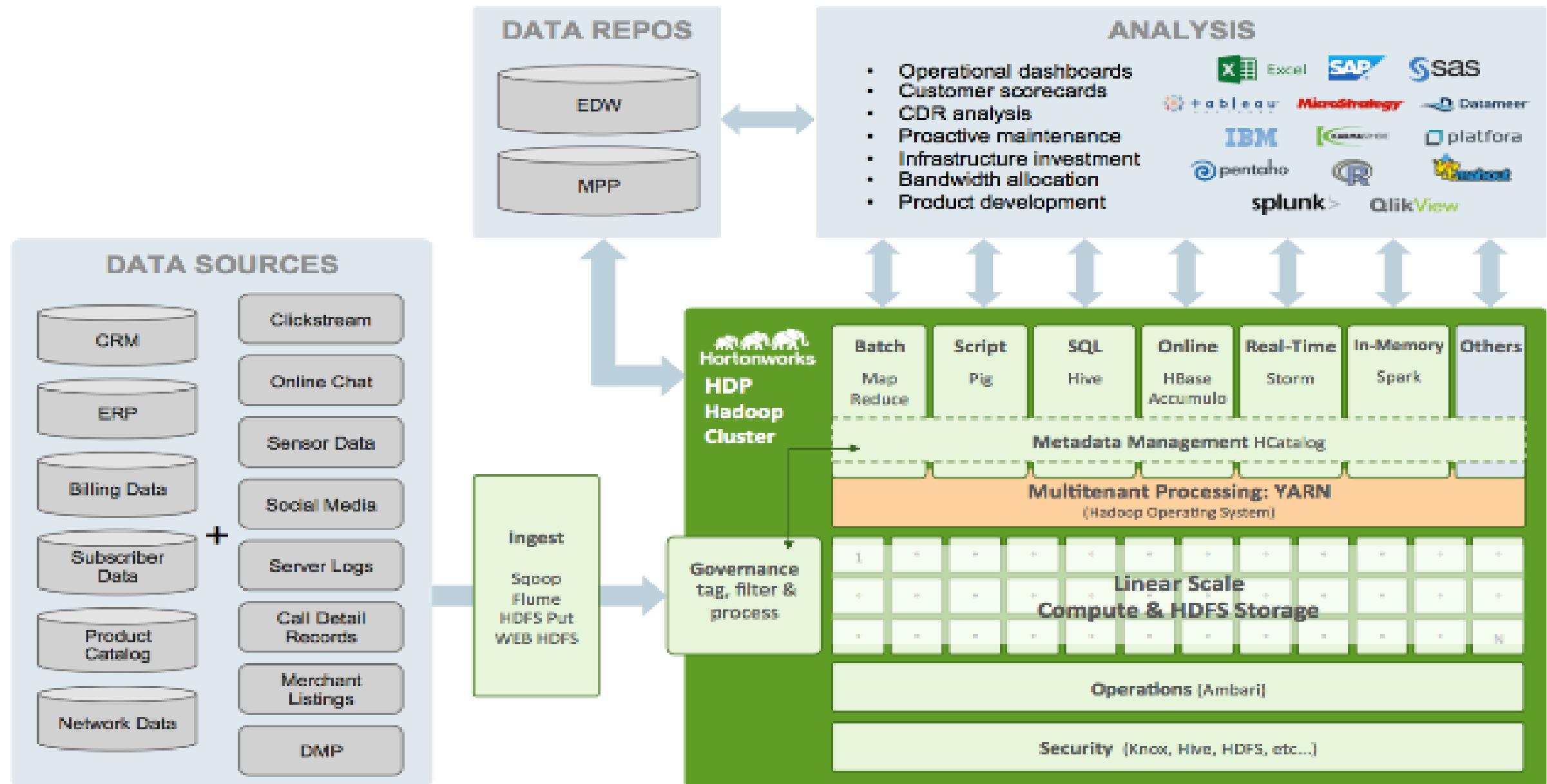
Technologies:

- Amazon S3
- Flume
- Hadoop/HDFS, MapReduce
- HBase
- Oozie
- Hive



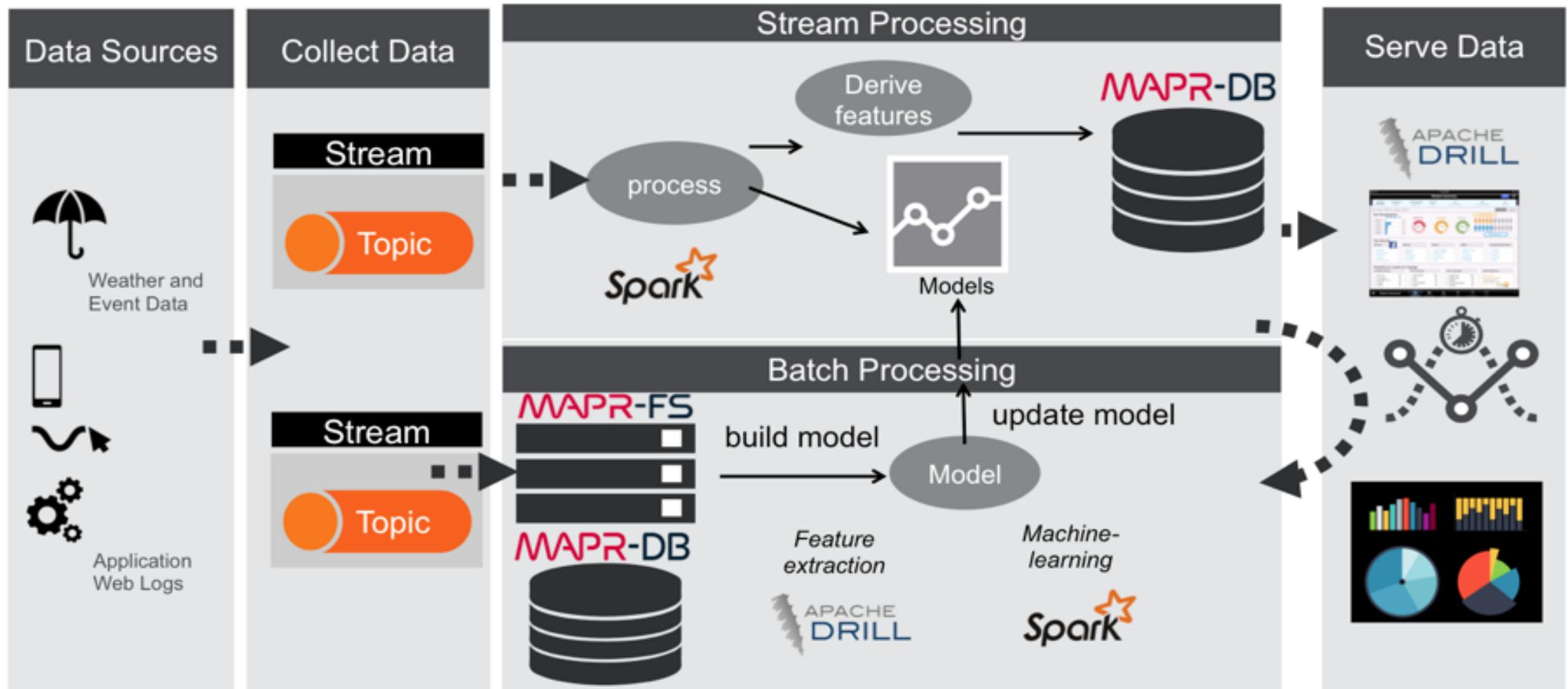


Big Data Process for Telecom



Big Data Process for Retail

Data-Driven Supply Chain





Location Tracking



Big Data Pipeline on AWS

