

BIG DATA PROCESSING

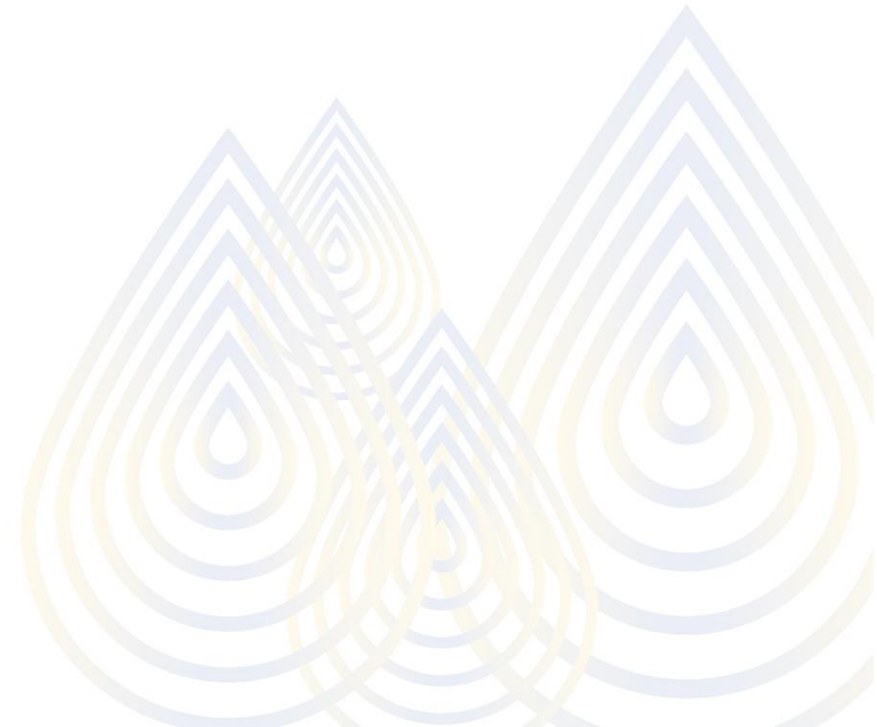
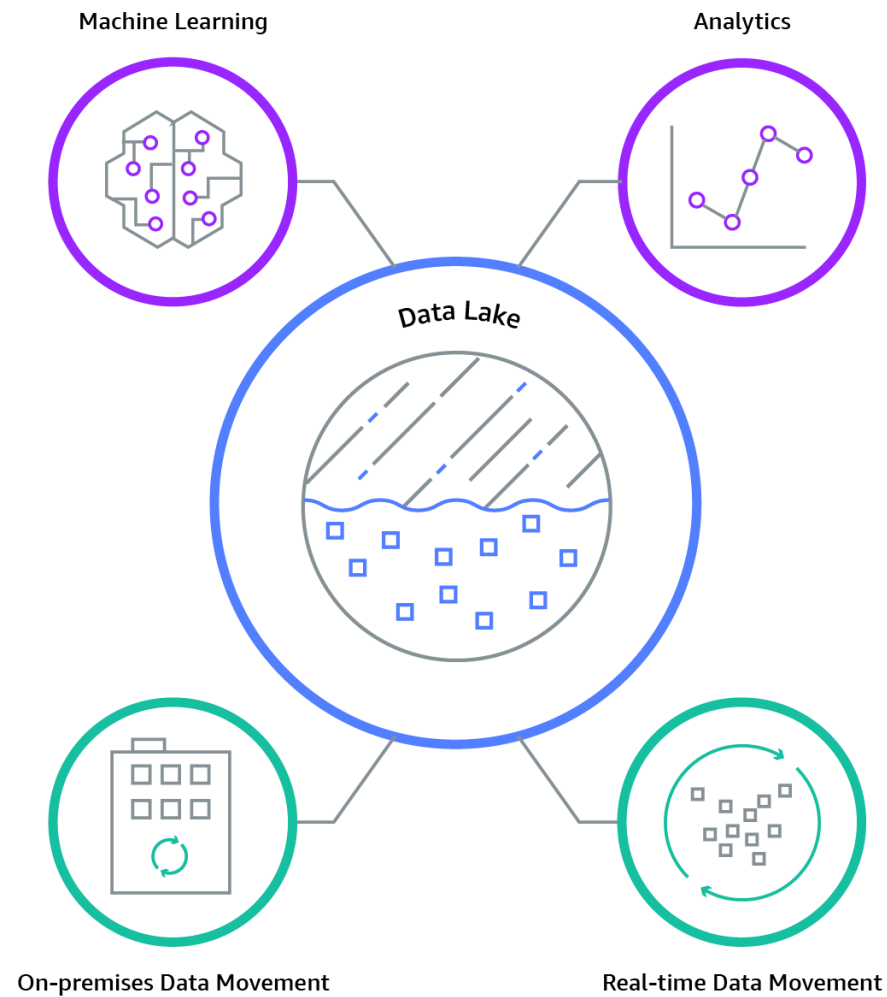
Week 7 Pipeline



Hadoop as a Big Data Platform

- Hadoop is one of **the most favourite**
 - Storage (HDFS)
 - Processing/Analysis (Hive, Spark, etc.)
 - Ecosystems (Sqoop, KafKa, AirFlow)
- Good choice for a data lake
 - Cheap storage compares with others
- On premise => **need large investment**
- **Need installation/administration**
- Con=> tie storage & processing together

Data Lake



Data Lakes

- Infrastructure that many stream can flow into
- Stored for processing in the original form
- Massive storage with huge processing power

“If you think of a data mart as a store of bottled water, cleansed and packaged and structured for easy consumption, the data lake is a large body of water in a more natural state. The contents of the data lake stream in from a source to fill the lake, and various users of the lake can come to examine it, dive in, or take samples.”-- James Dixon(2010), the Pentaho Corporation's CTO

Data Lakes Process

Load data from Source



Store Raw data

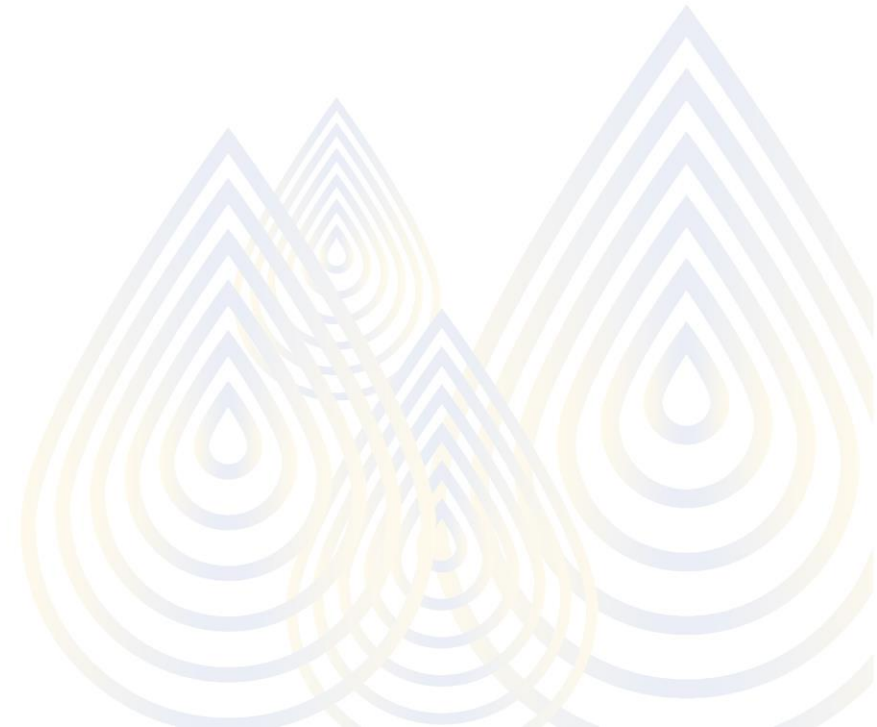


Add data model on read

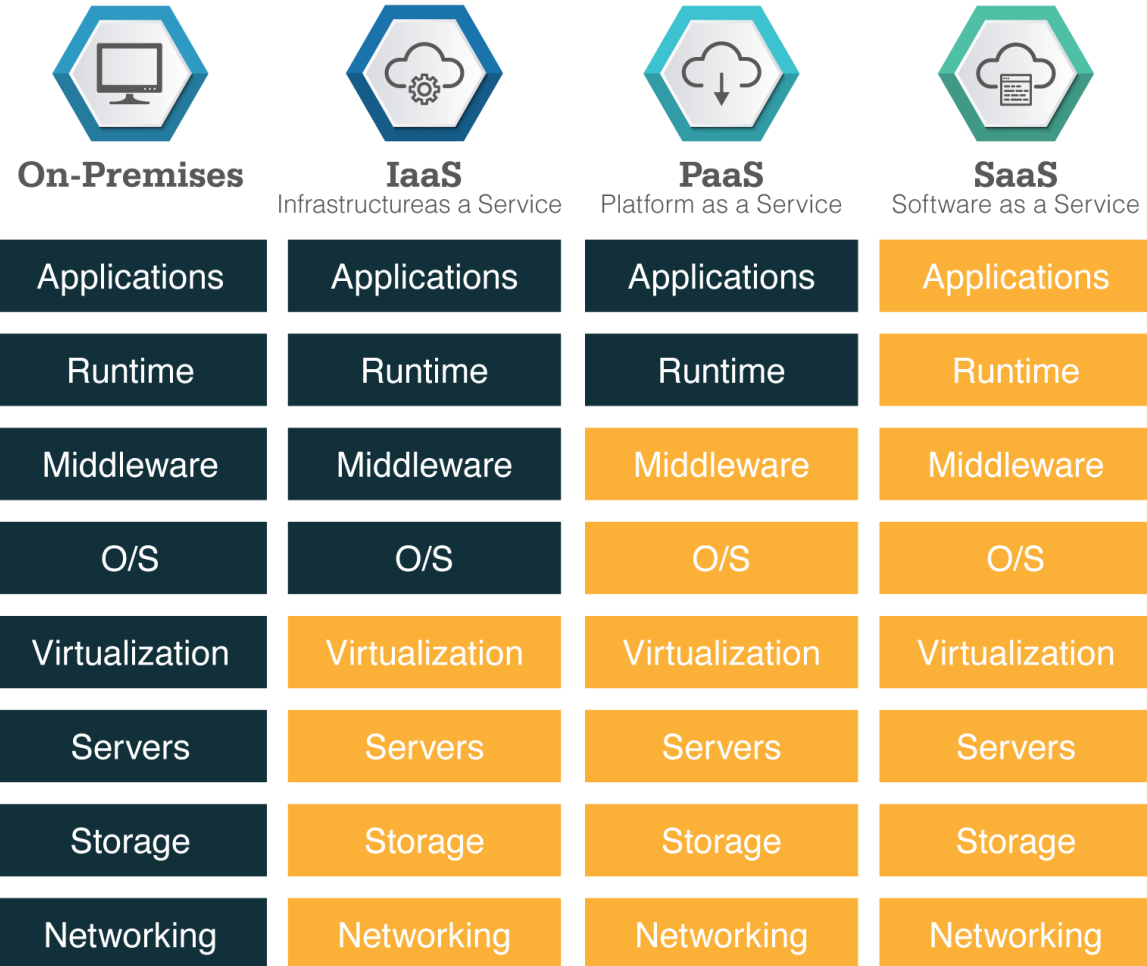
Schema on read

Use cloud storage as a data lake

- Cheap storage for data archive
- Cost per GB/month
- High availability
- High durability
- Focus on warm/cold data
- Transaction data
- Encrypt data for security

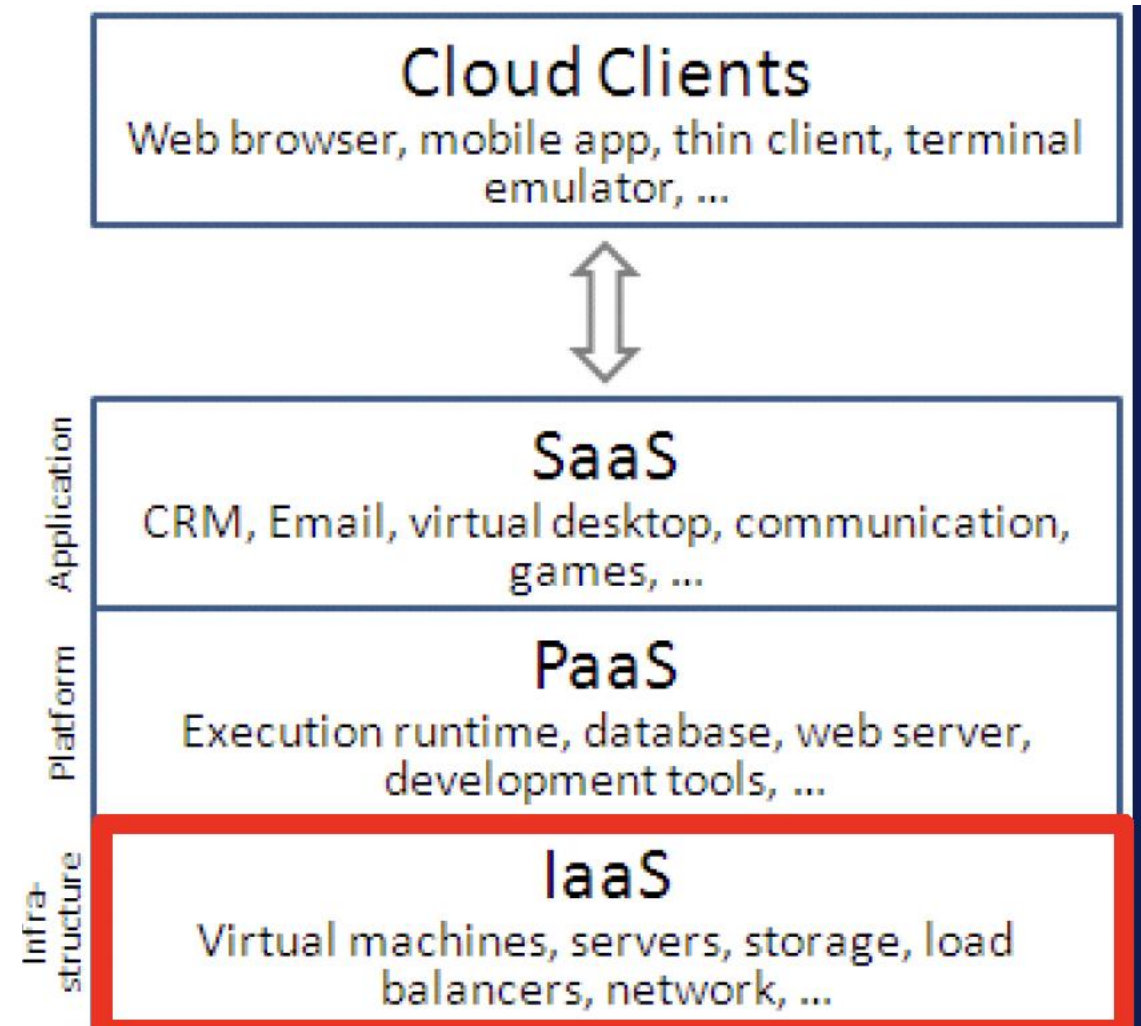


Cloud Service Model



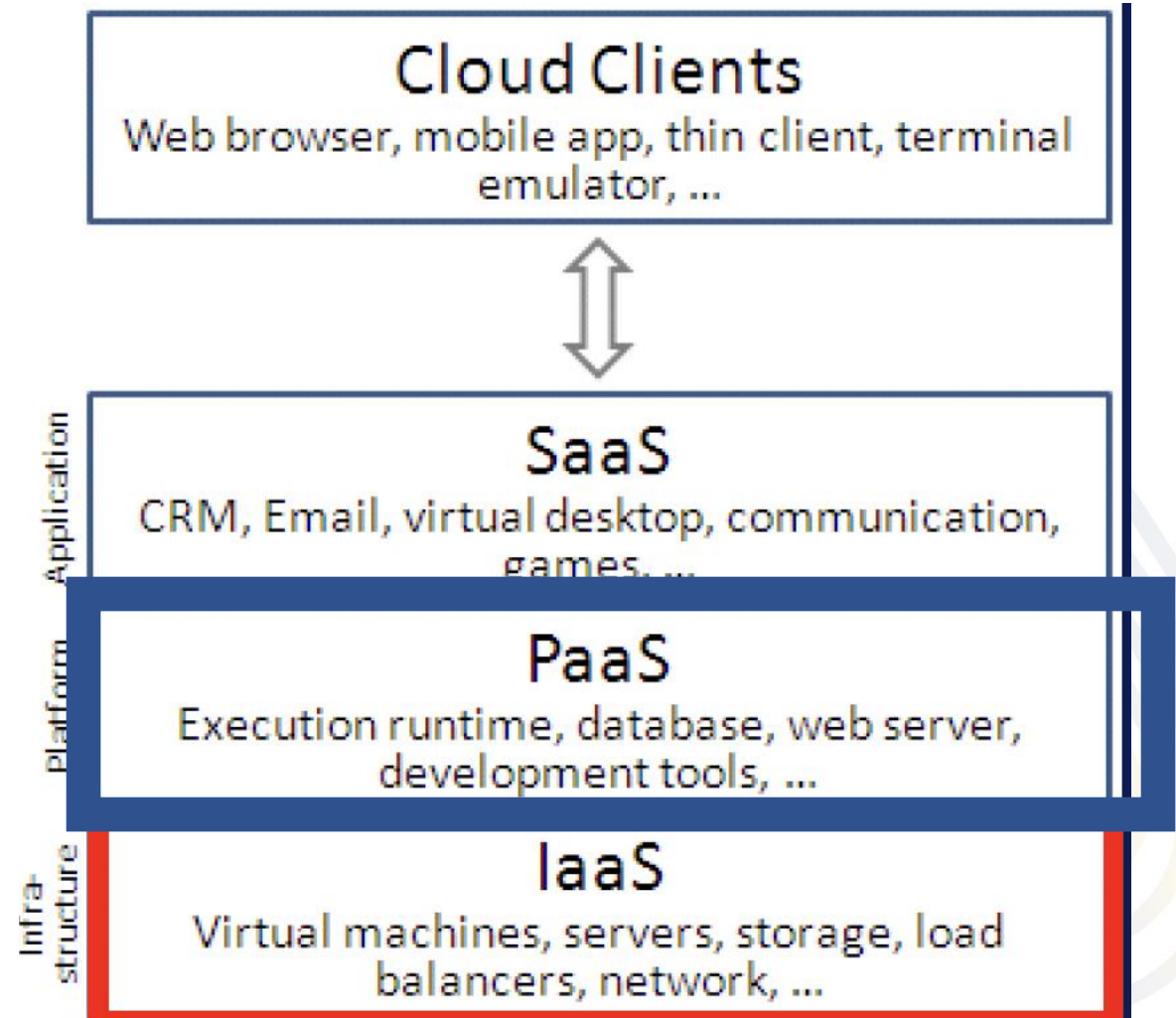
IaaS

- Infrastructure as a Service
- Install & maintain OS by yourself
 - Amazon EC2
 - Virtual Machine
 - Server
 - Load balancer
 - Network



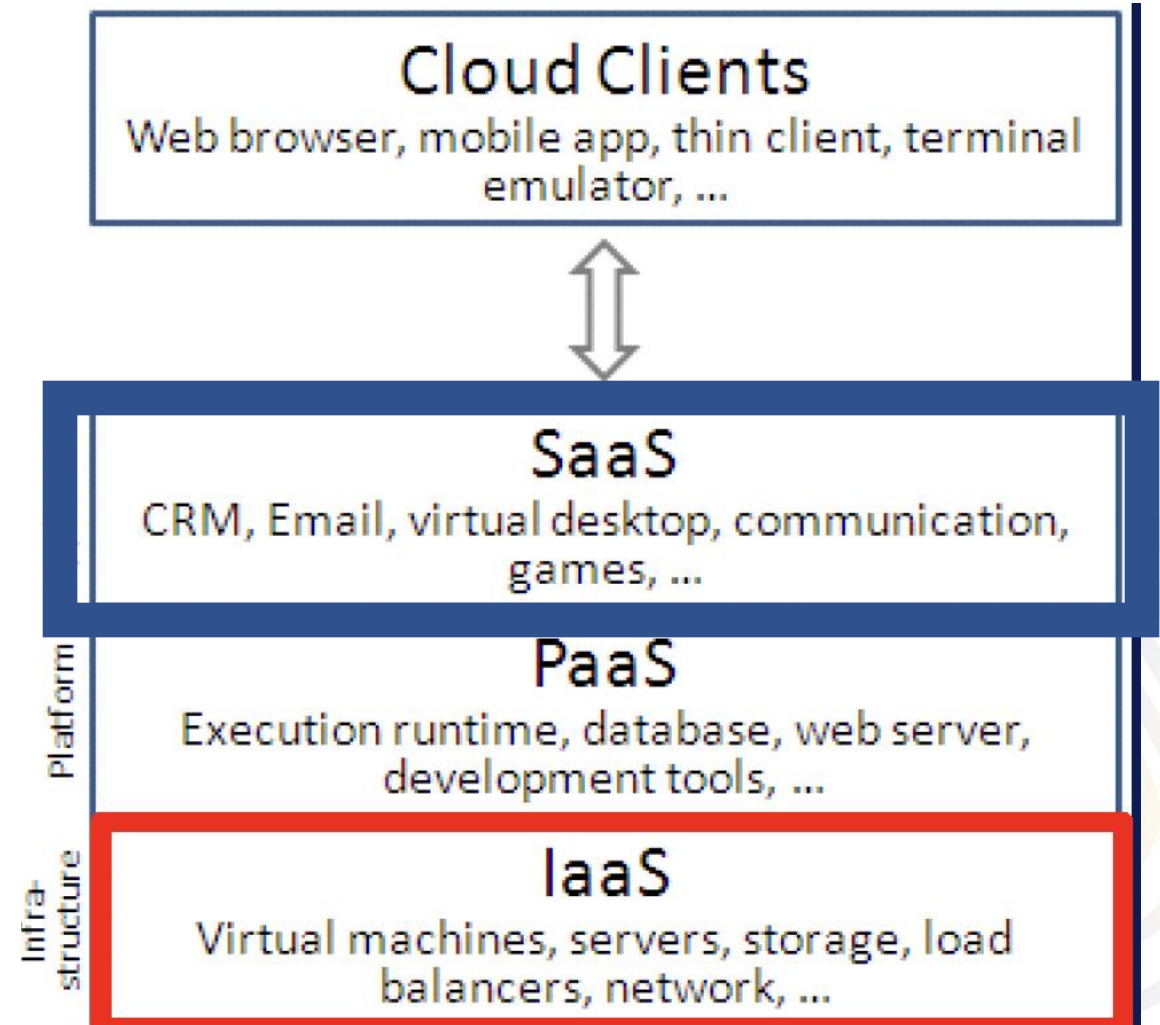
PaaS

- Platform as a Service
- Install Application by yourself
 - Google App, Microsoft Azure
 - Execution runtime
 - Database
 - Web server
 - Development Tools

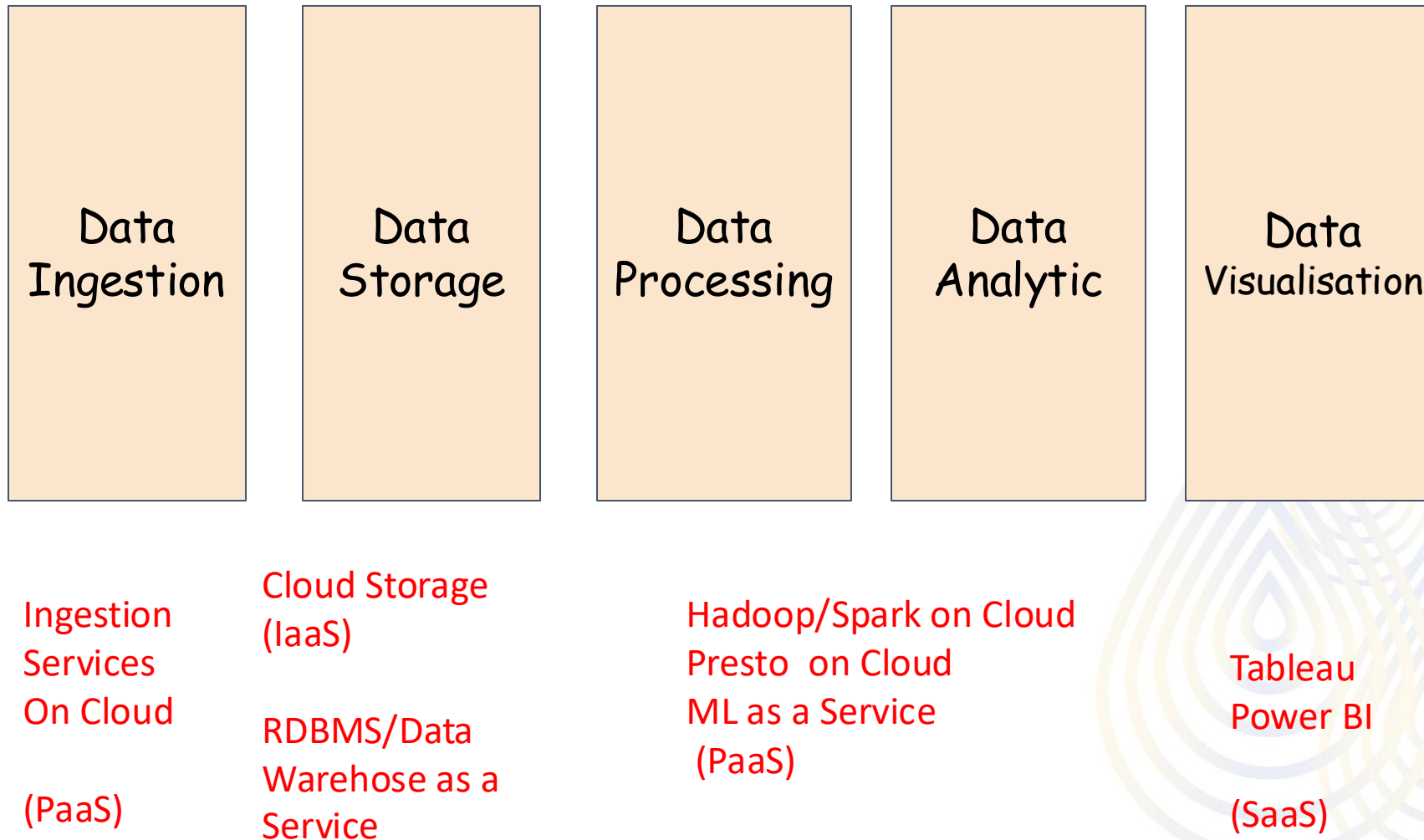


SaaS

- Software as a Service
- Just use an application
 - Dropbox, Google App

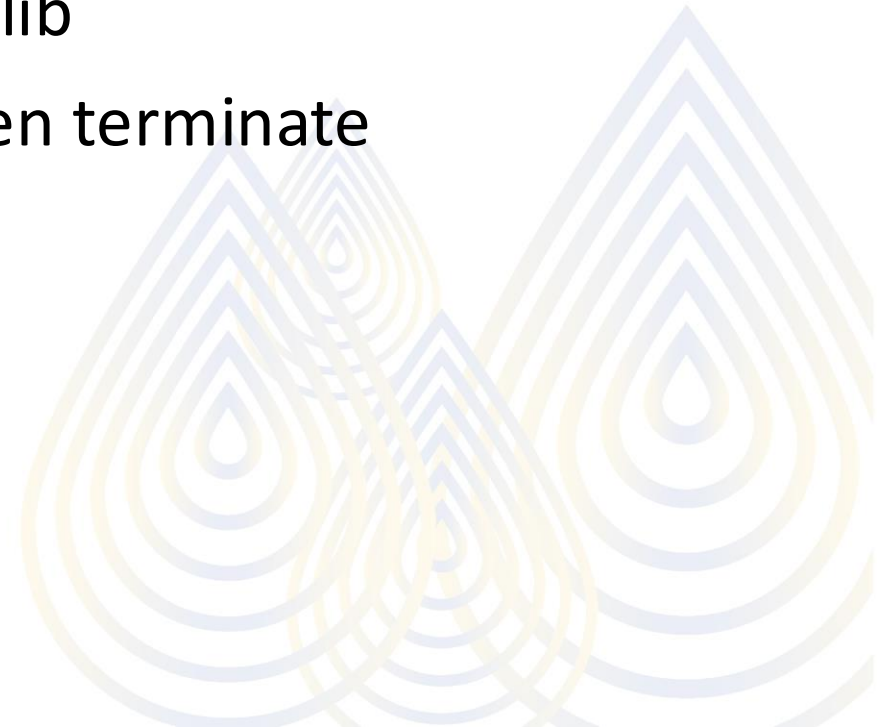


Typical Big Data Batch Pipeline on Cloud



Spark/Hadoop as a service for analytics

- Separate process layer from storage layer
- No need to install/admin a cluster
- Use processing tools: Hive, Spark, MLlib
- Start cluster only processing time, then terminate cluster when finish.
- Scalable CPU powers
- Pay only processing (CPU) time.



Analytics as a service



Amazon EMR



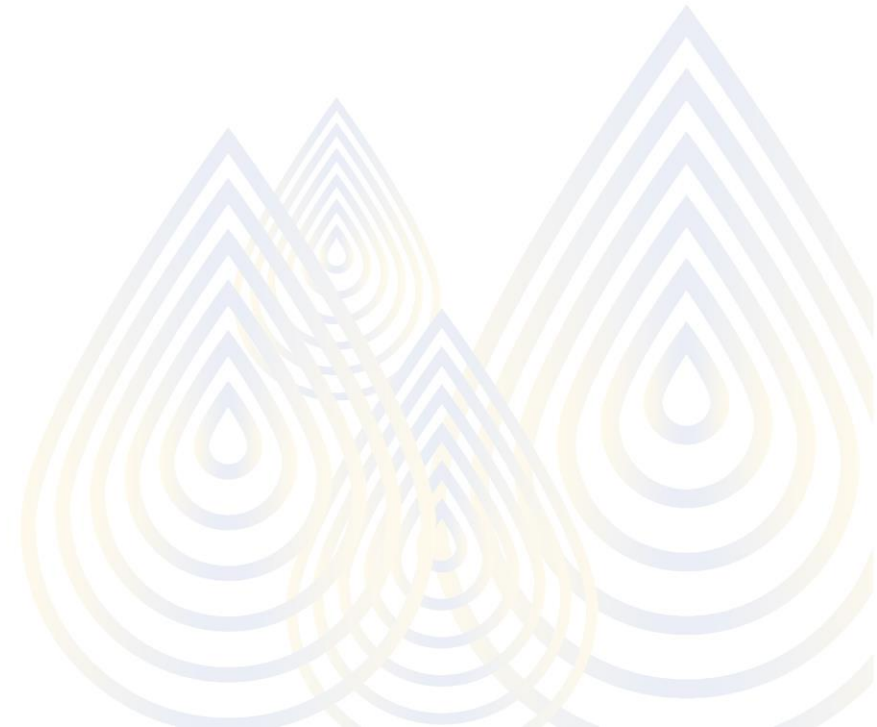
Cloud Dataproc



Big Data 

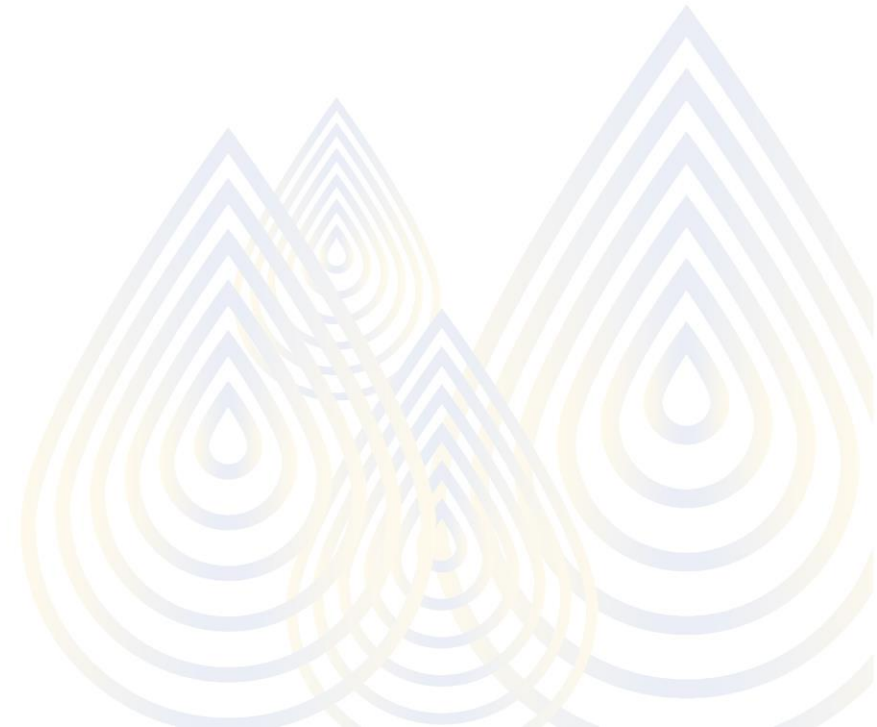
Data Warehouse as a Service

- A fully-managed and cloud-based interactive query service for massive datasets
- Petabyte Data Warehouse
- Example Cloud Services
 - Google BigQuery
 - Amazon Redshift
 - Azure Synapse

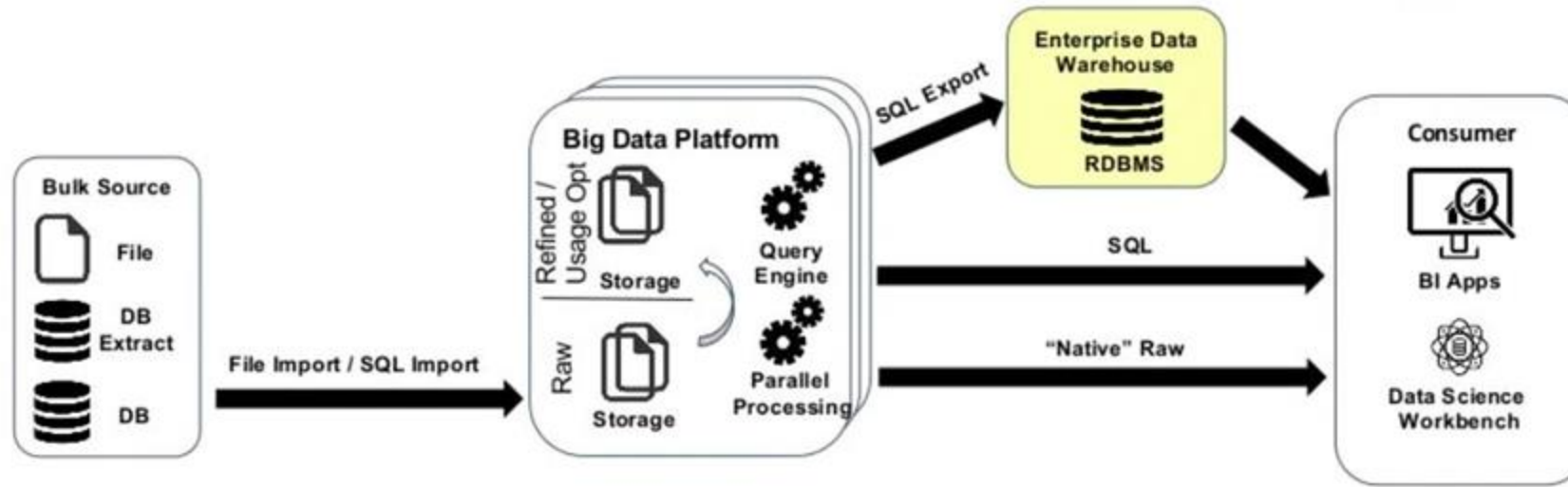


Fast SQL as a Service

- MPP SQL
- Example Cloud Services
 - Google BigQuery
 - Amazon Athena
 - Azure Data Lake Analytics

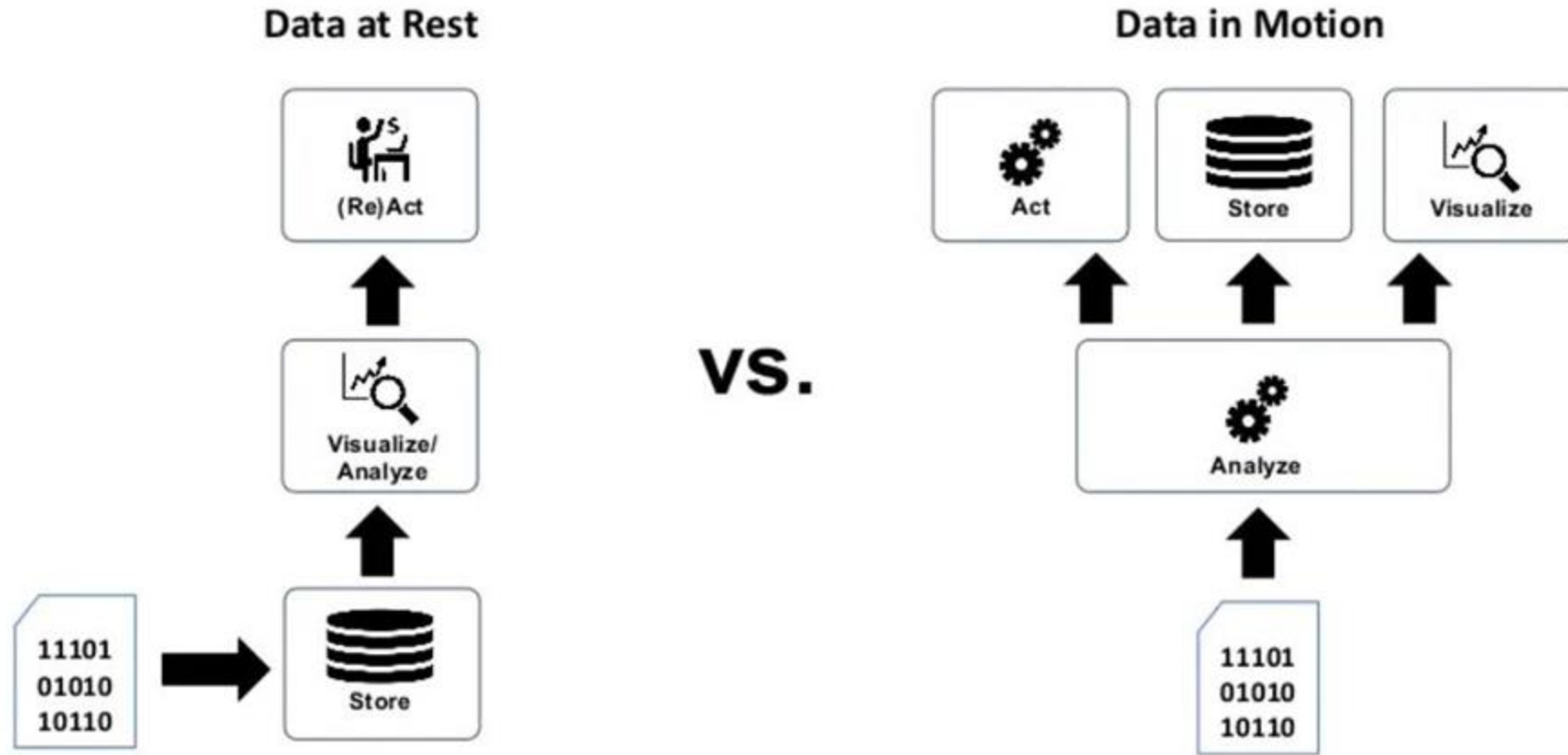


1) Data Lake + DW Architecture

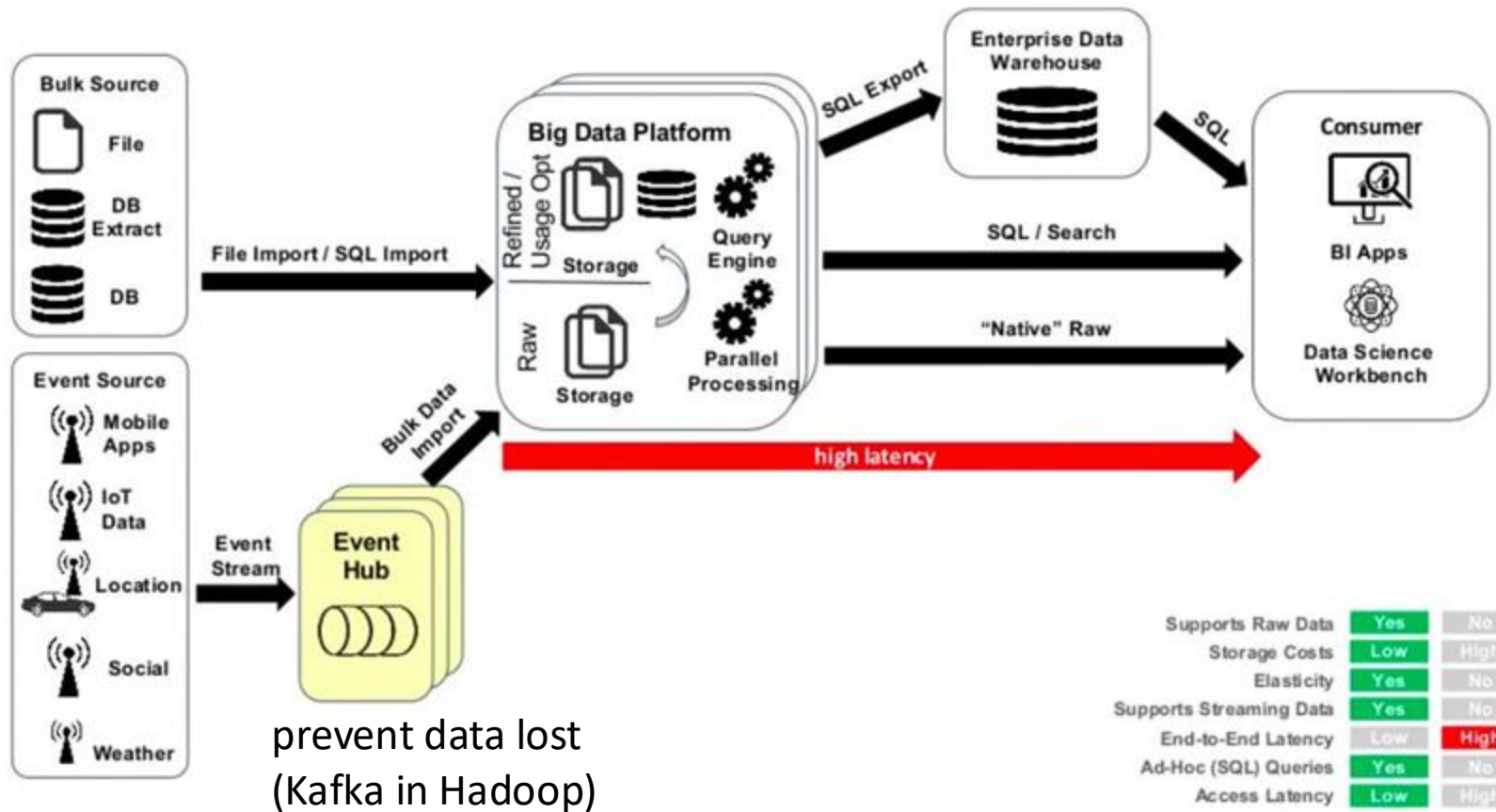


Low latency

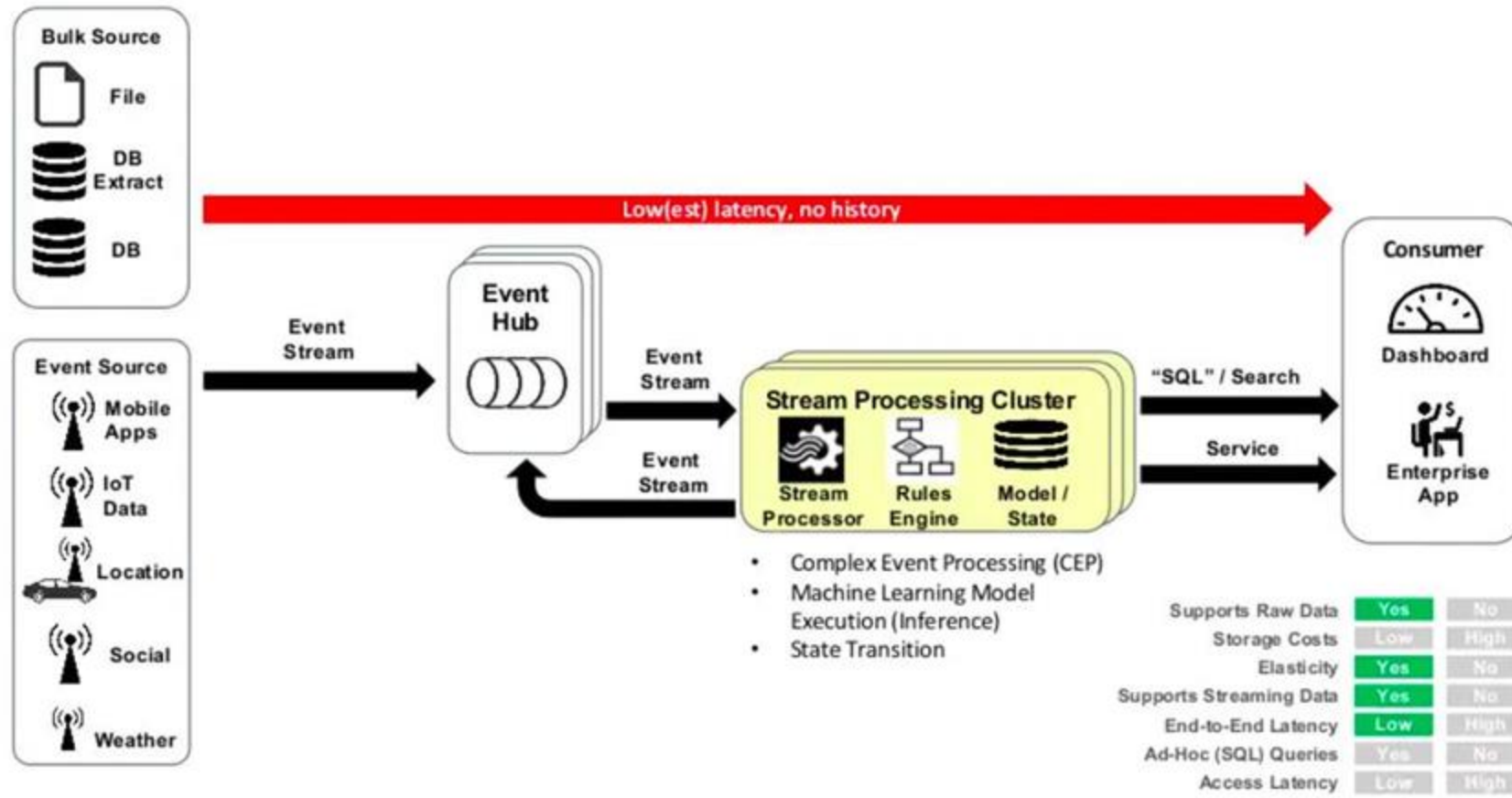
Supports Raw Data	Yes	No
Storage Costs	Low	High
Elasticity	Yes	No
Supports Streaming Data	Yes	No
End-to-End Latency	Low	High
Ad-Hoc (SQL) Queries	Yes	No
Access Latency	Low	High



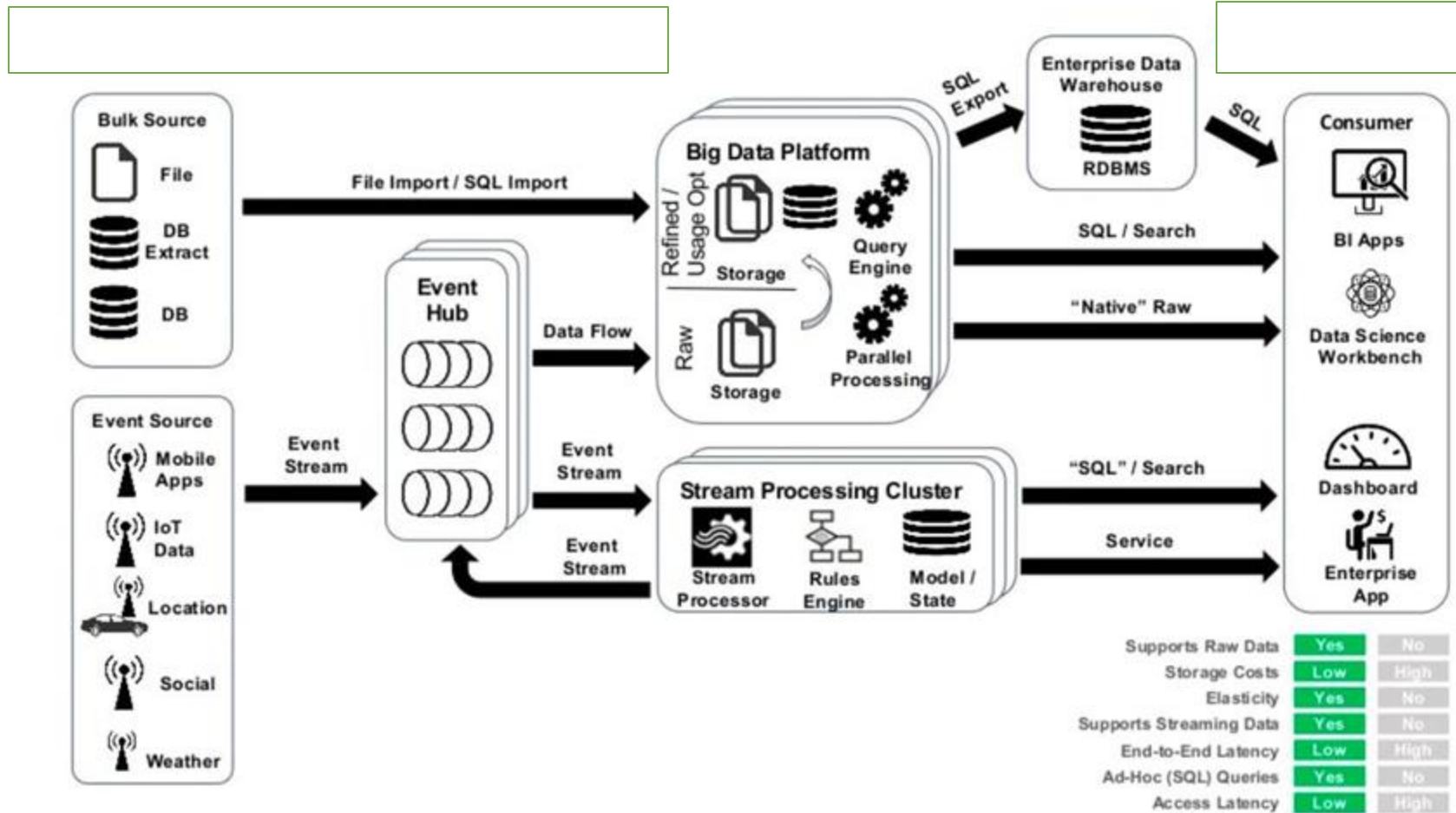
2) Big Data Architecture with Streaming



Event Processing Architecture



3) Data Lake + Event Processing



mu

Data Visualization



Data Processing / Data Analytics

Batch



Scriptin



SQL



In-



NoSQL



Streamin



Machine Learning



Data Storage



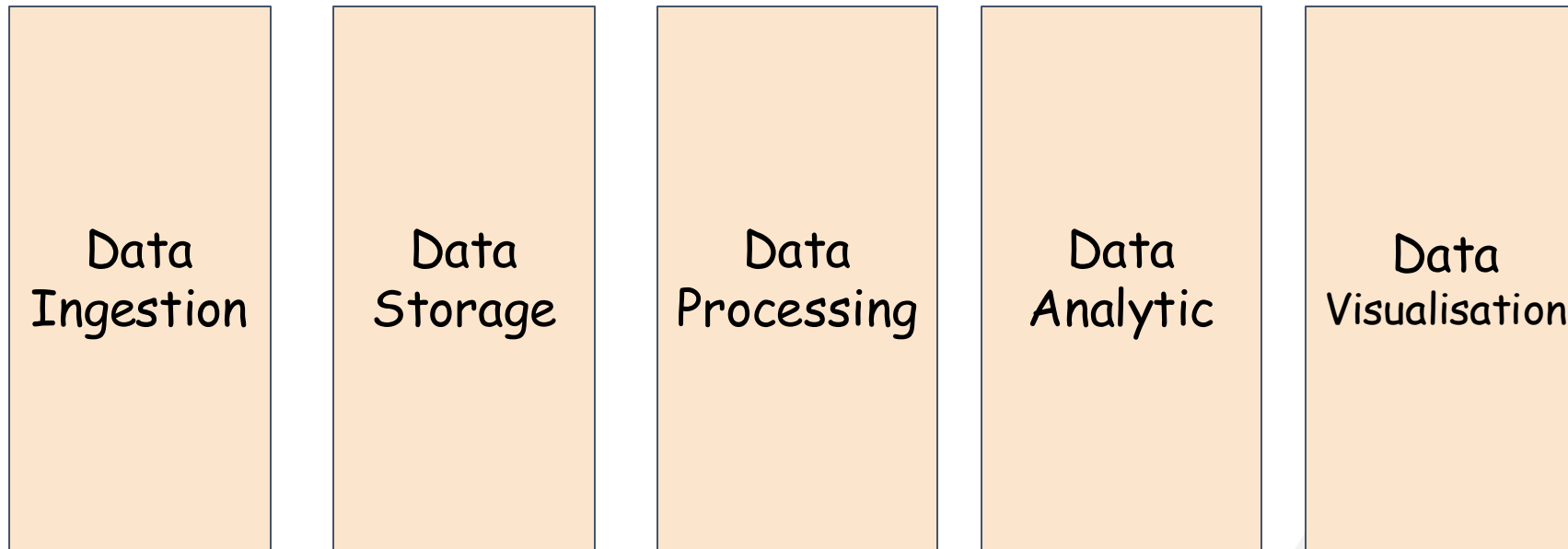
Data Ingestion





Mahidol University
Wisdom of the Land

Typical Big Data Batch Pipeline



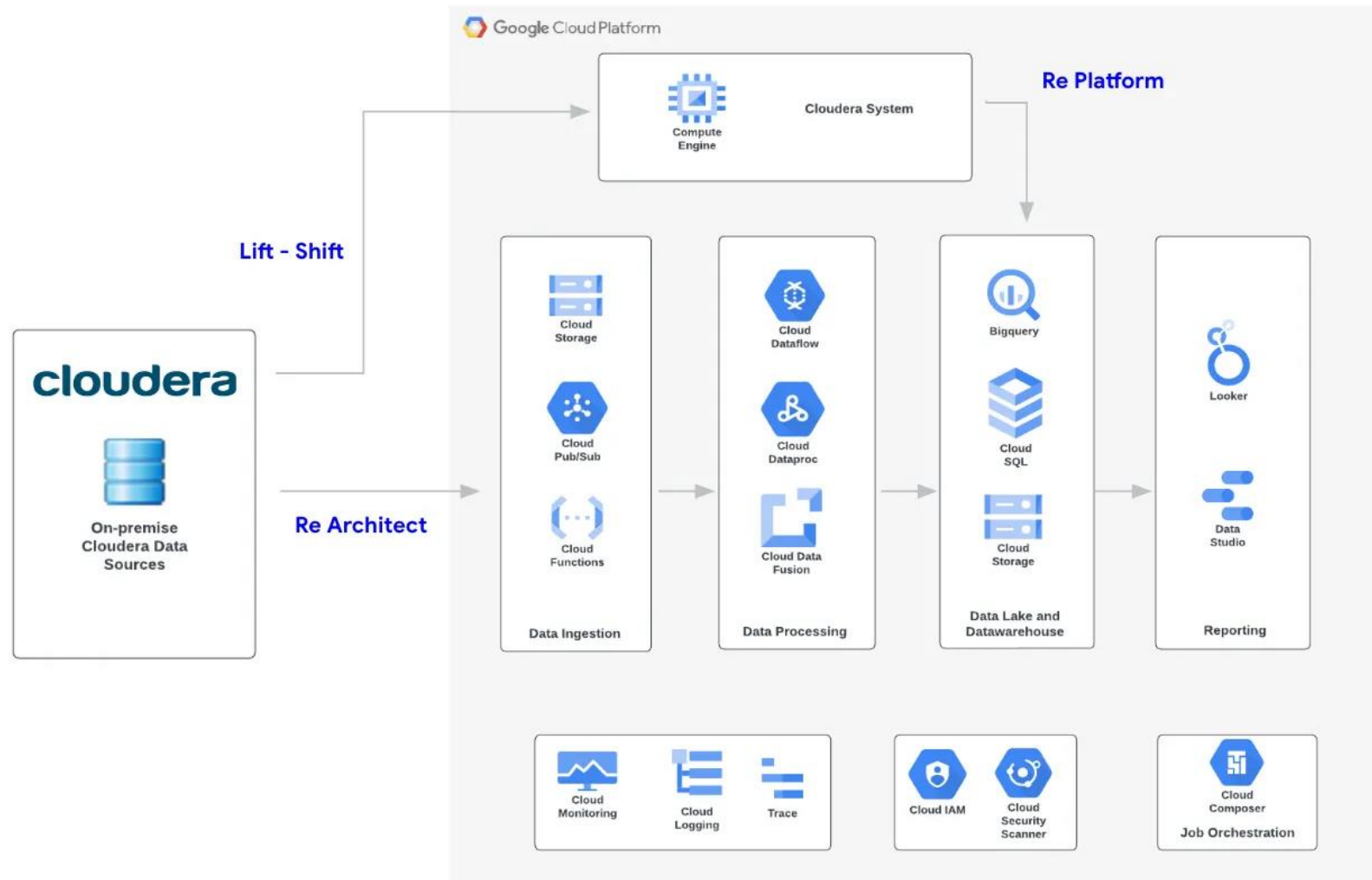
Sqoop
KafKa
Aegisthus

Hadoop HDFS
S3
Cassandra
MySQL Cluster
Elastic search
Redis

Spark
Storm
Spark Streaming
Hive, Pig
Presto

Tableau

On-Premise v.s. Public Cloud Services



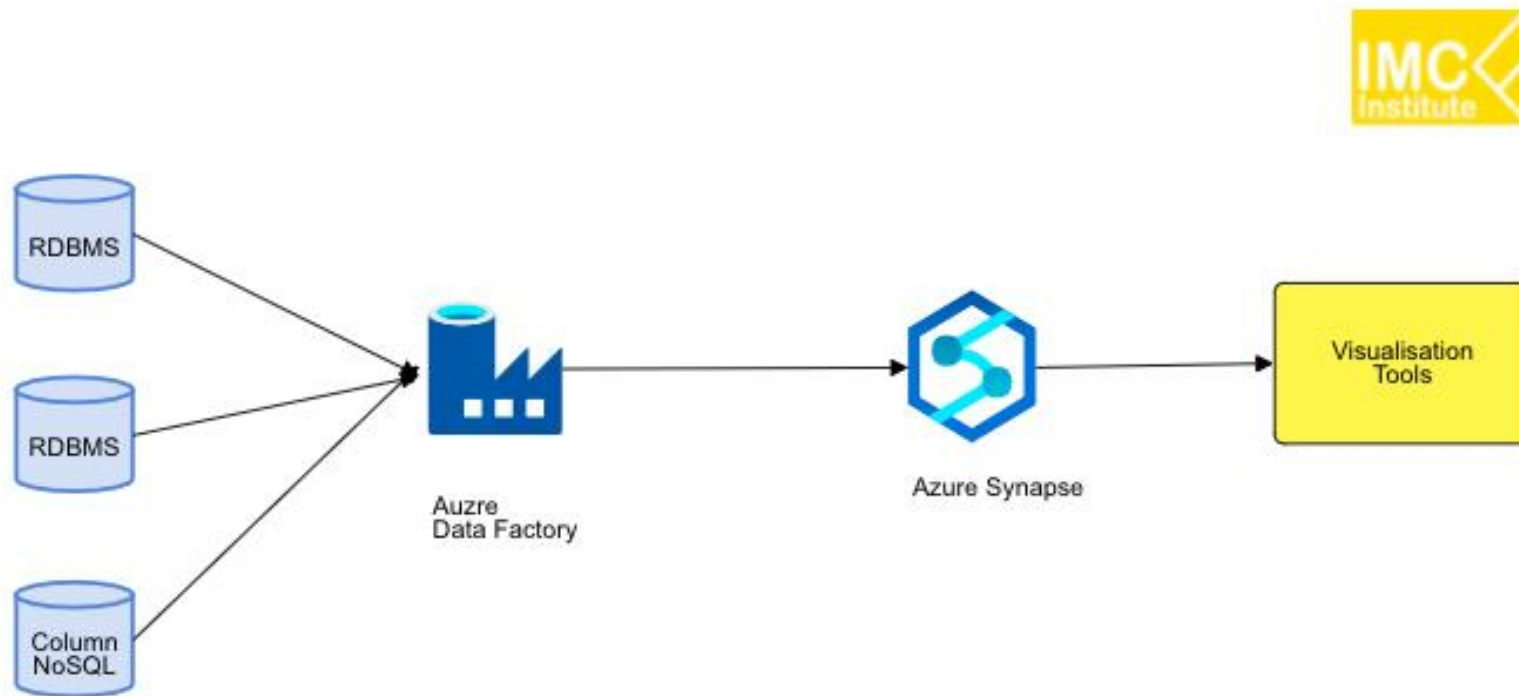
On-Premise v.s. Public Cloud Services

Provided by: IMC Institute		On-Premise	GCP	AWS	Azure	Huawei Cloud	Alibaba Cloud	Oracle Cloud
Ingestion	Event	KafKa	Cloud Sub/Pub	Managed Streaming for KafKa	Apache Kafka on HDInsight	Distributed Message Service (DMS) for Kafka	DataHub	Streaming
		Flume		Kinesis	Event Hub	Data Ingestion Service (DIS)		
	ETL	Sqoop / Nifi	Cloud DataFusion	Glue	Data Factory	Data Replication Service (DRS)	Data Integration	Data Integration
Storage	Data Warehouse	Data Warehouse	Google BigQuery	Redshfit	Synapse	Data Warehouse Service (DWS)	AnalyticDB	Autonomous Data Warehouse
	Data Lake	Hadoop HDFS	Cloud Storage	S3	Azure Data Lake	Object Storage Service (OBS)	Object Storage Service	Object Storage
	Cold Data Archive		Cold Line Cloud Storage/ Arcive storage	S3 Glacier	Archive Storage			Archive Storage
	Hot data	RDBMS	Cloud SQL	RDS	Azure Database	RDS	ApsaraDB RDS	Autonomous Database
		NoSQL	Cloud Datastore/ Cloud BigTable	DynamoDB	Cosmos DB	Gauss DB	Tablestore	NoSQL
Processing / Analytics	Generic	Hadoop/Spark	Cloud DataProc	EMR	HDInsight/ Azure Databrick	MapReduce Service / Data Lake Insight (DLI)	E-MapReduce	Big Data Service
	Streaming	Spark Streaming	Cloud Dataflow	Kinesis Analytics	Stream Analytics	Cloud Stream Service (CS)	Realtime Compute for Apache Flink	Data Flow
	Fast SQL	Hive	Google BigQuery	Athena	Data Lake Analytics	Data Lake Insight (DLI)	Data Lake Analytics	Autonomous Database
	Machine Learning	Spark MLlib	Cloud AI Platform	SageMaker	Azure ML	ModelArts	Machine Learning Platform For AI	Data Science
Visualisation	Dashboard	Tableau, Power BI etc.	Google Data Studio, Looker	QuickSight	Power BI	Data Lake Visualisation (DLV)	Quick BI, Data V	Oracle Analytic Platform



Mahidol University
Wisdom of the Land

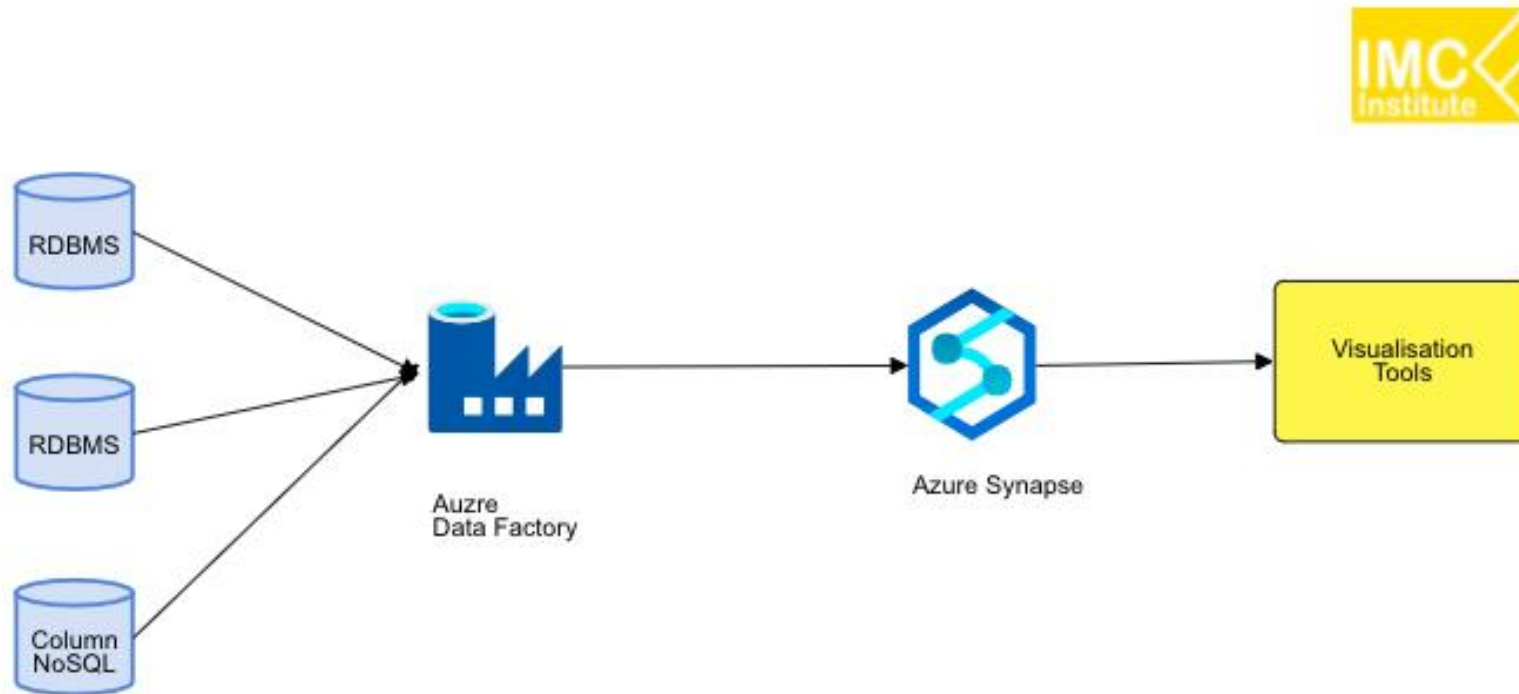
Data warehouse architecture on Azure



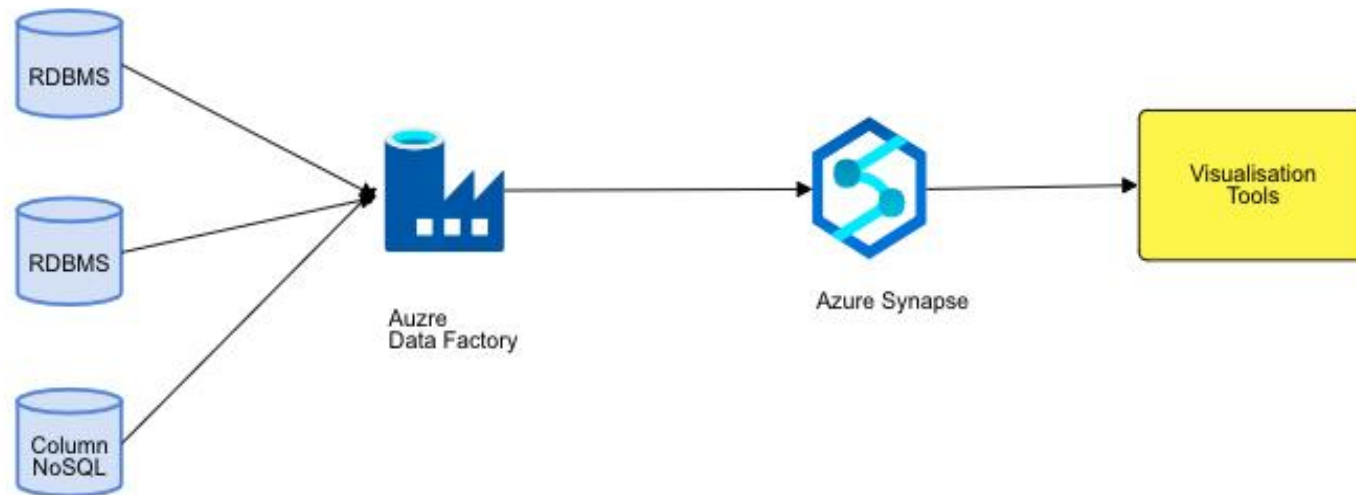


Mahidol University
Wisdom of the Land

Data warehouse architecture on Azure



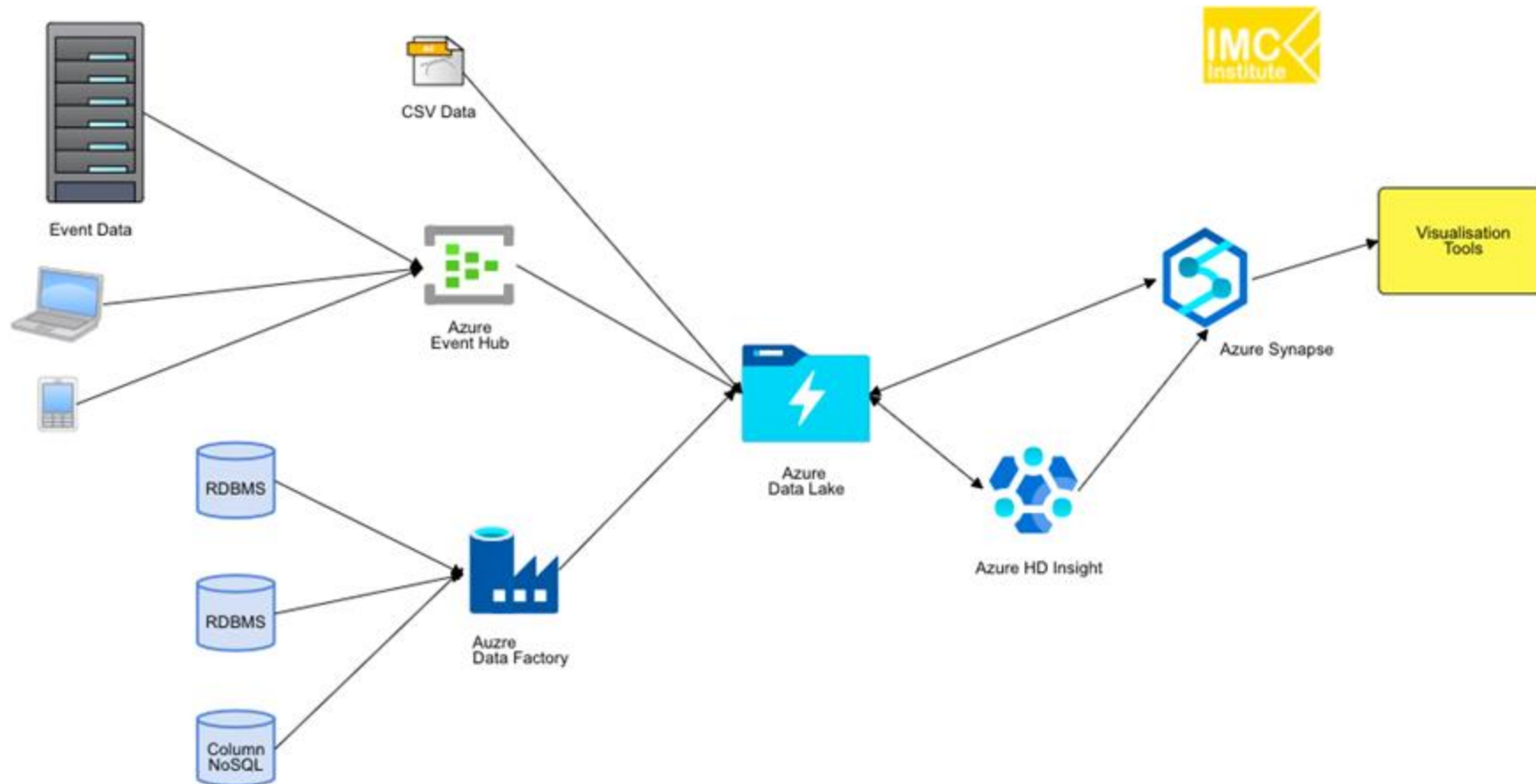
GROUP WORK 1



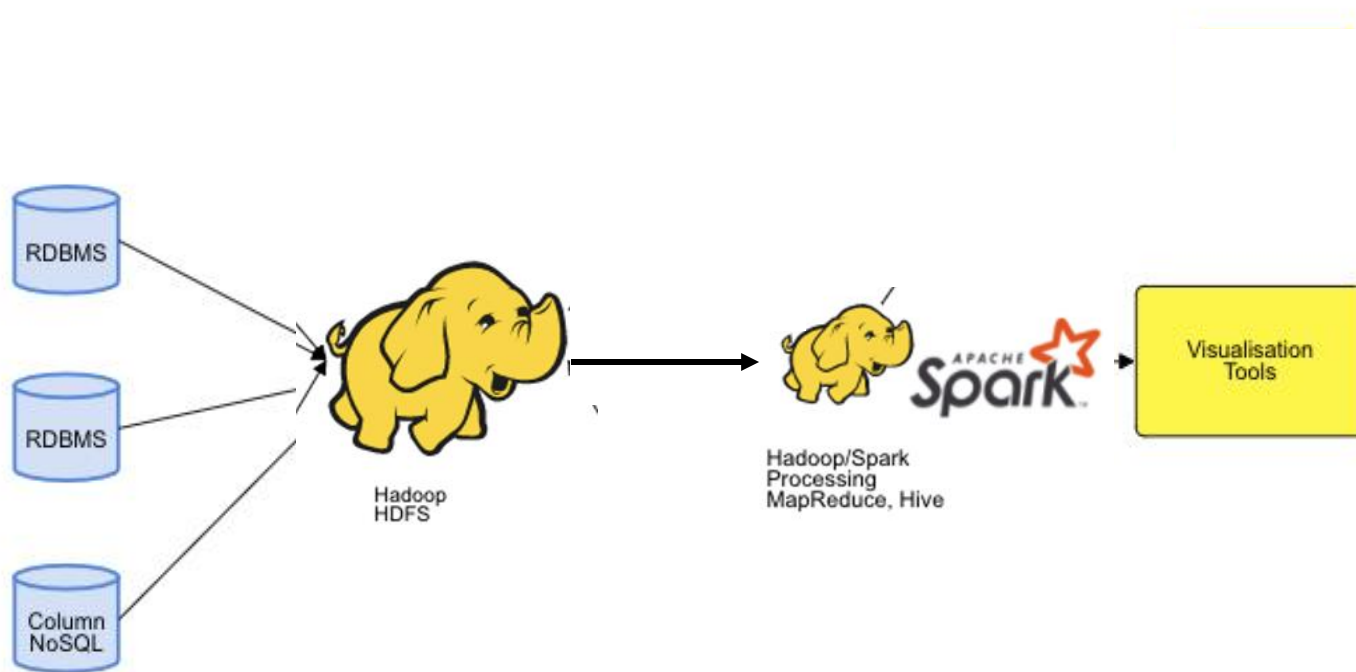
Compare it with

- GCP
- AWS

GROUP WORK 2



HDFS Architecture

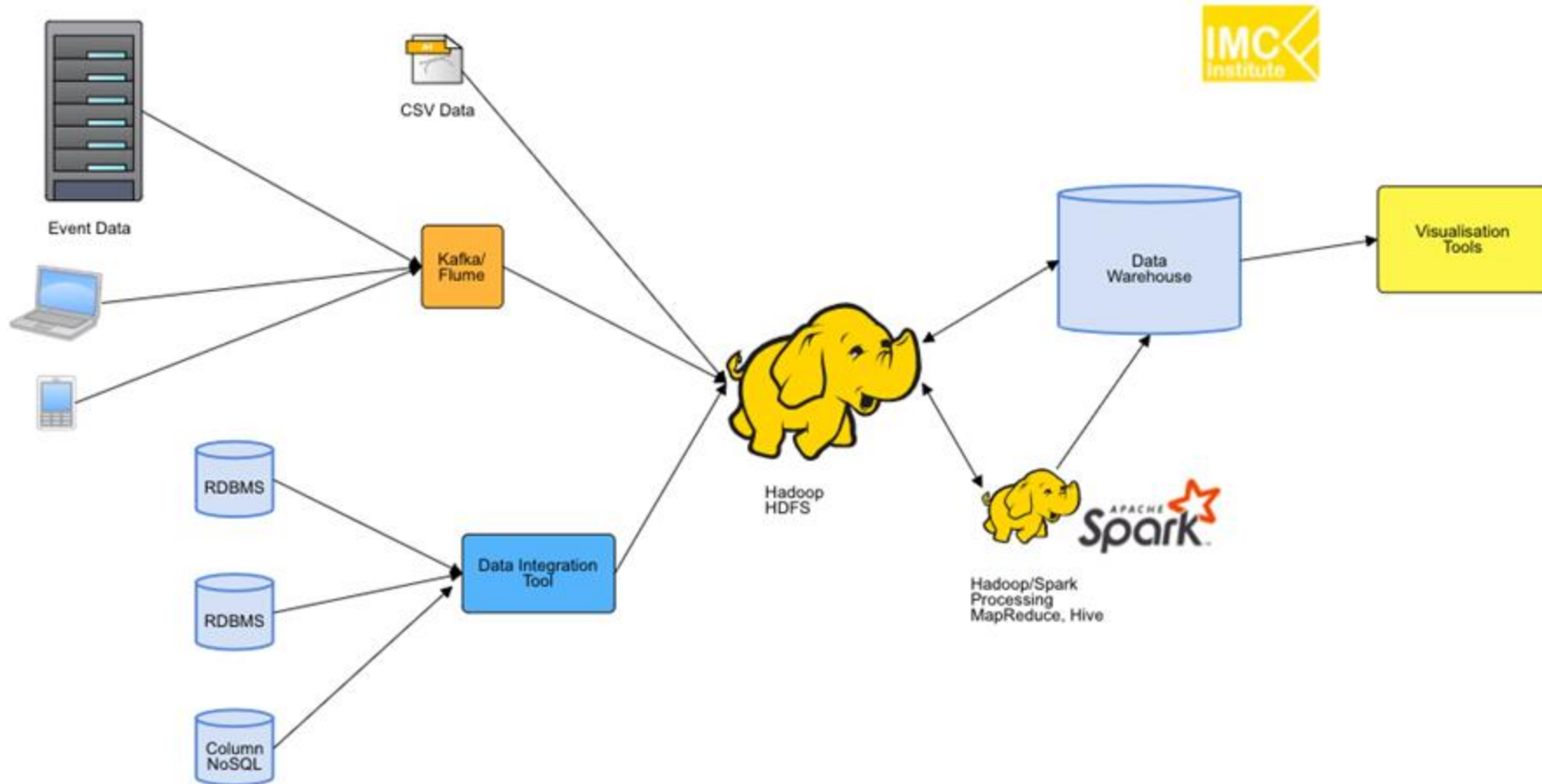


GROUP WORK 1

Compare it with

1. Replace with other platforms
 - GCP
 - Azure
 - AWS
 - Huawei
 - Alibaba
 - Oracle
2. Check the price (assume usage by yourself)
 - GCP <https://cloud.google.com/products/calculator?hl=en>
 - Azure <https://azure.microsoft.com/en-us/pricing/calculator/>
 - AWS <https://calculator.aws/#/>
 - Huawei Price calculator <https://www.huaweicloud.com/intl/en-us/pricing/calculator.html#/ecs>
 - Alibaba https://www.alibabacloud.com/en/pricing-calculator?_p_lc=1#/
 - Oracle <https://www.oracle.com/cloud/pricing/>

Group work2

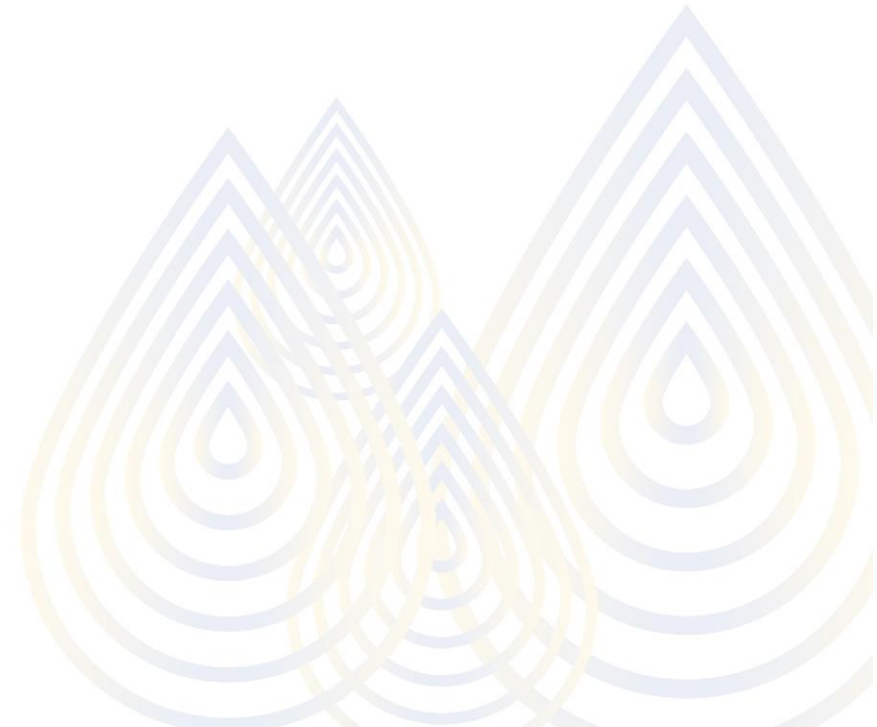


Price Check

- GCP <https://cloud.google.com/products/calculator?hl=en>
- Azure <https://azure.microsoft.com/en-us/pricing/calculator/>
- AWS <https://calculator.aws/#/>
- Huawei Price calculator <https://www.huaweicloud.com/intl/en-us/pricing/calculator.html#/ecs>
- Alibaba https://www.alibabacloud.com/en/pricing-calculator?_p_lc=1#/
- Oracle <https://www.oracle.com/cloud/pricing/>

Issue with Big data on cloud

- Vendor lock-in
- Latency
- Data privacy on cloud
- Regulation/Compliance

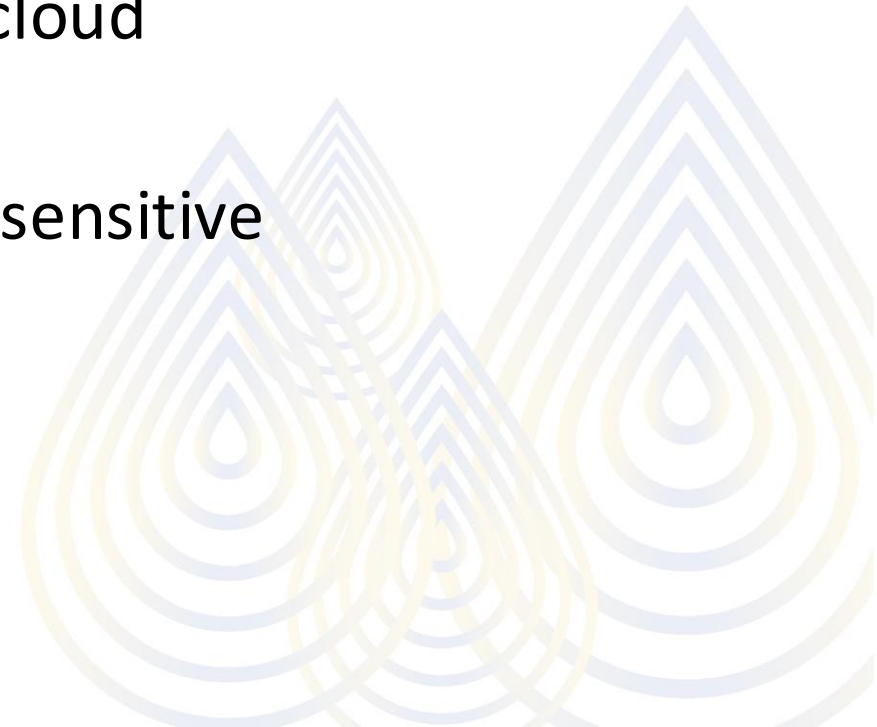


Hybrid & Multi Cloud Model

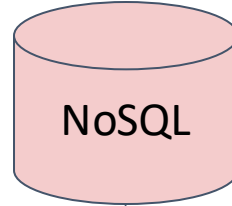
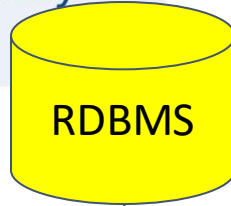
- Using both on-premise & public cloud services
- On-premise storage and processing for sensitive data
- Cloud storage for high scale data and also large archive data
- Scalable CPU for Hadoop as a services
- Hadoop on-premise distribution & Big data on cloud trend towards multi cloud model.

Recommendation

- Always start with big data on public cloud
- Cloud storage as a main data lake
- Process large scale data using public cloud services such as Hadoop as a Service
- Small on-premise big data cluster for sensitive data and local processing.



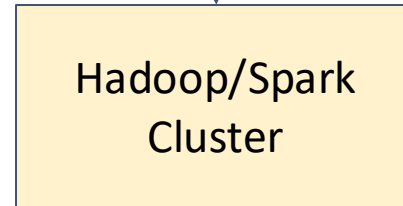
Hot Data



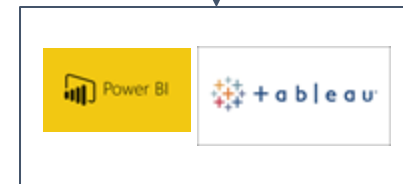
Warm/Cold
Data



Analytics

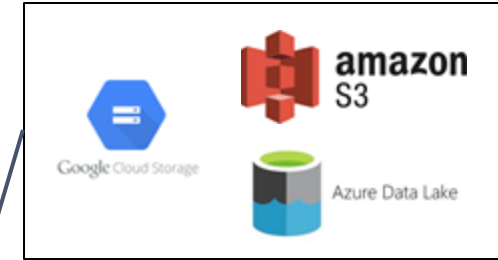


Visualisation



On-Premise

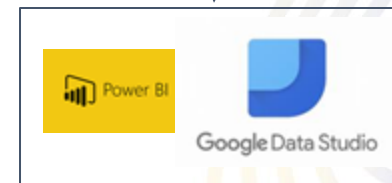
VPN
Tunnel



Unlimited
Low cost
Cloud storage



Scaleable
Hadoop/
Spark
Cluster



Cloud
Visualisation
Tools

Public Cloud

HDFS vs. Cloud Storage: Pros, cons and migration tips

<https://cloud.google.com/blog/products/storage-data-transfer/hdfs-vs-cloud-storage-pros-cons-and-migration-tips>

