

# Data Lake

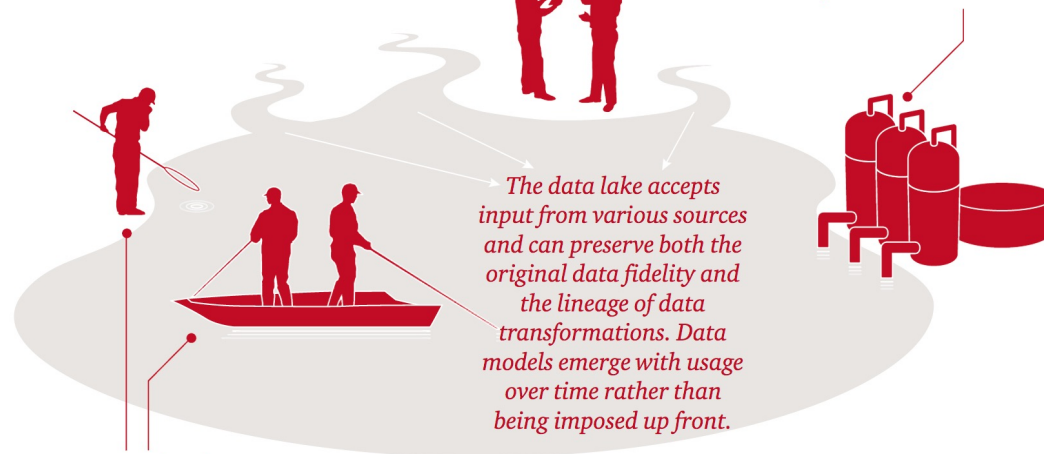
Enterprise Architectures for Big Data

# What is a Data Lake?

A repository for large quantities and varieties of data, both structured and unstructured.

Data generalists/  
programmers can tap  
the stream data for  
real-time analytics.

The lake can serve as a staging  
area for the data warehouse,  
the location of more carefully  
“treated” data for reporting  
and analysis in batch mode.



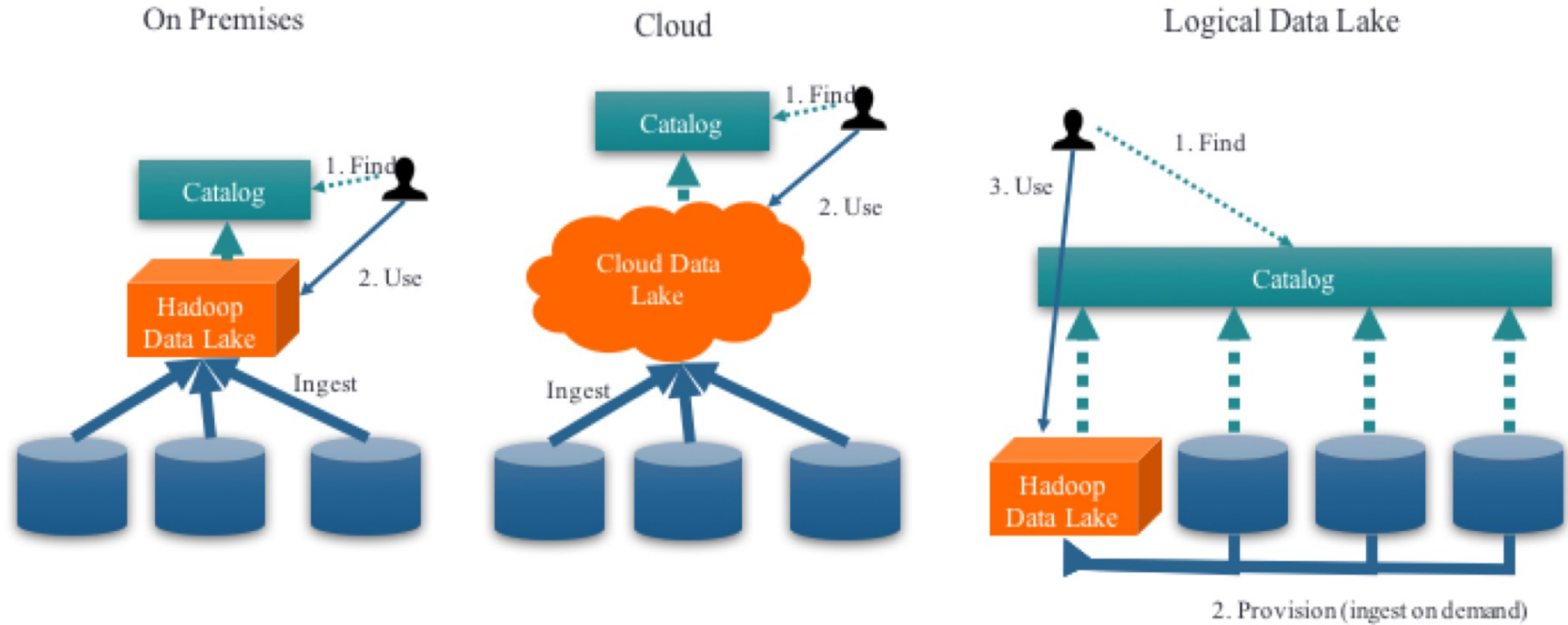
Data scientists  
use the lake for  
discovery and  
ideation.

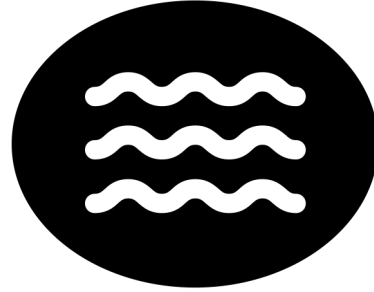
Data lakes take advantage of commodity cluster computing techniques for massively scalable, low-cost storage of data files in any format.

# Data Lake

- Repository for raw data
- Includes structured, semi-structured and unstructured data
- Often Hadoop / HDFS based
- No upfront schema (schema at read vs. schema at write)
- Users often Data Scientists
- Could be a source for downstream systems like data warehouses
- Flexible usage and access

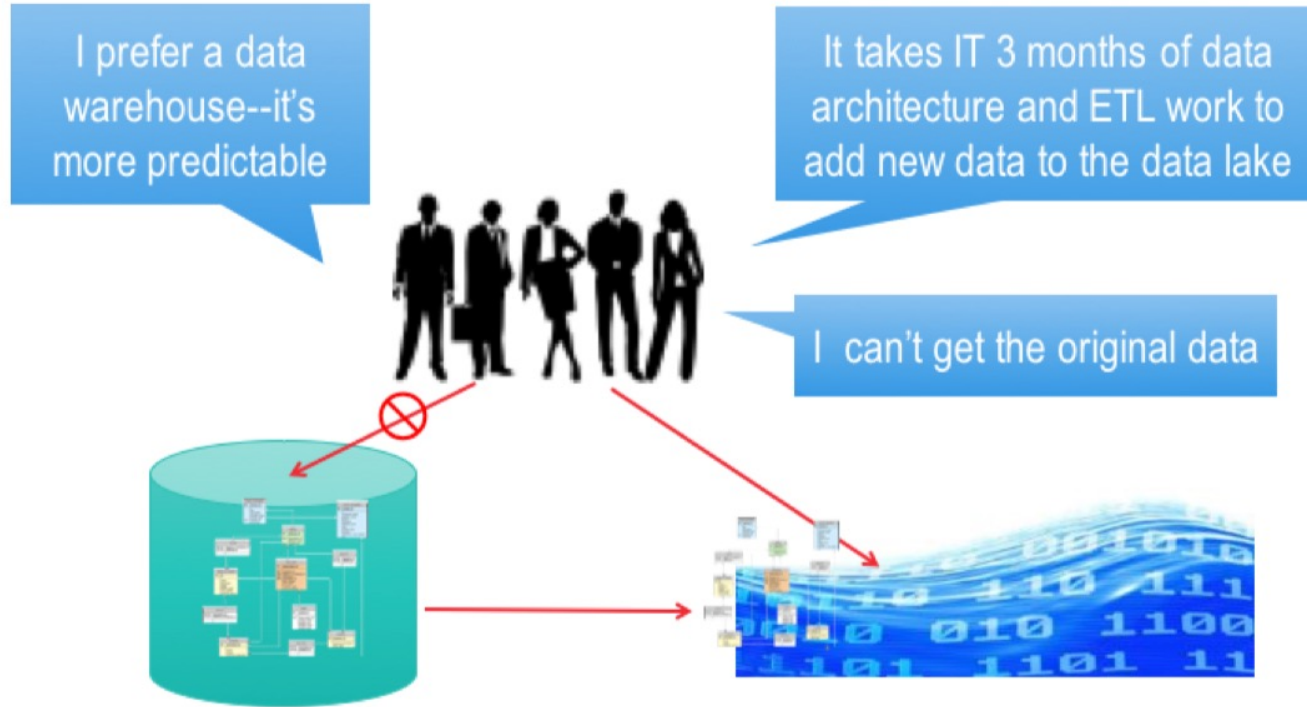
# Different Data Lake Architectures



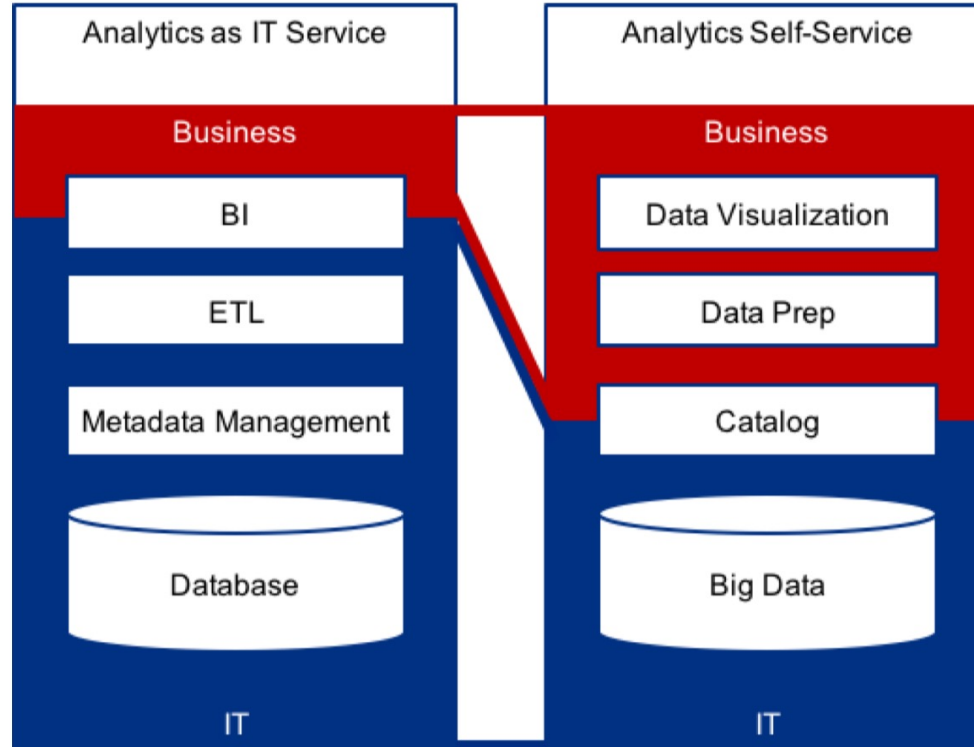


# Data Warehouse vs. Data Lake

# Drawbacks of data warehouse offloading



# Enabling analysts and reducing the load on IT with self-service analytics



# Data Warehouse Solves two main problems

1. Integrated Data Schema
2. Data Locality



# Data Warehouses vs. Data Lakes

## Data Warehouses

- Schema on Write
- ETL > Analytics/OLAP
- Queries mostly Interactive
- Structures Data
- Often RDBM based

## Data Lakes

- Schema on Read
- EL > T > Analytics/OLAP
- Queries mostly Batch
- Structured, Semi-Structured and Unstructured
- Often Hadoop Based

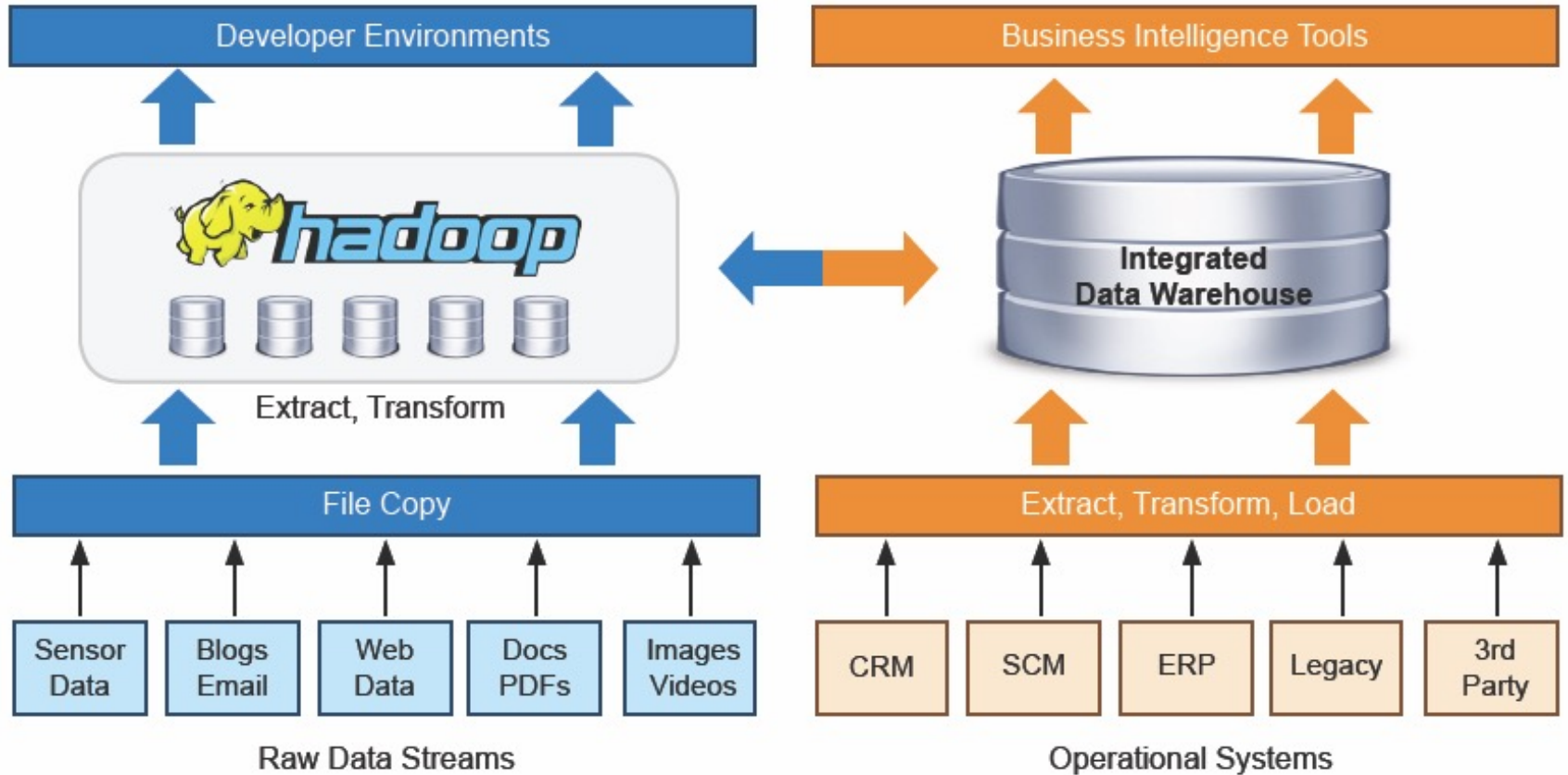
# Schema on-write vs. Schema on-read

Attribute	Schema On-Write (Early Binding) Data Warehousing	Schema on-read (Late Binding) Data Lake
Data Provider	<ul style="list-style-type: none"><li>• Evaluate Data</li><li>• Define Data Structure</li><li>• Collect Data &amp; Ingest</li><li>• Apply Structure</li></ul>	<ul style="list-style-type: none"><li>• Collect Data &amp; Ingest</li></ul>
Data Consumer	<ul style="list-style-type: none"><li>• Answer Questions</li></ul>	<ul style="list-style-type: none"><li>• Evaluate Data</li><li>• Define Data Structure</li><li>• Apply Structure</li><li>• Answer Questions</li></ul>
Ideal for	<ul style="list-style-type: none"><li>• Reused &amp; Known Data</li><li>• Consistent Results</li><li>• The Masses</li></ul>	<ul style="list-style-type: none"><li>• Unfamiliar Data</li><li>• Infrequent usage</li><li>• Unstable source schema</li></ul>

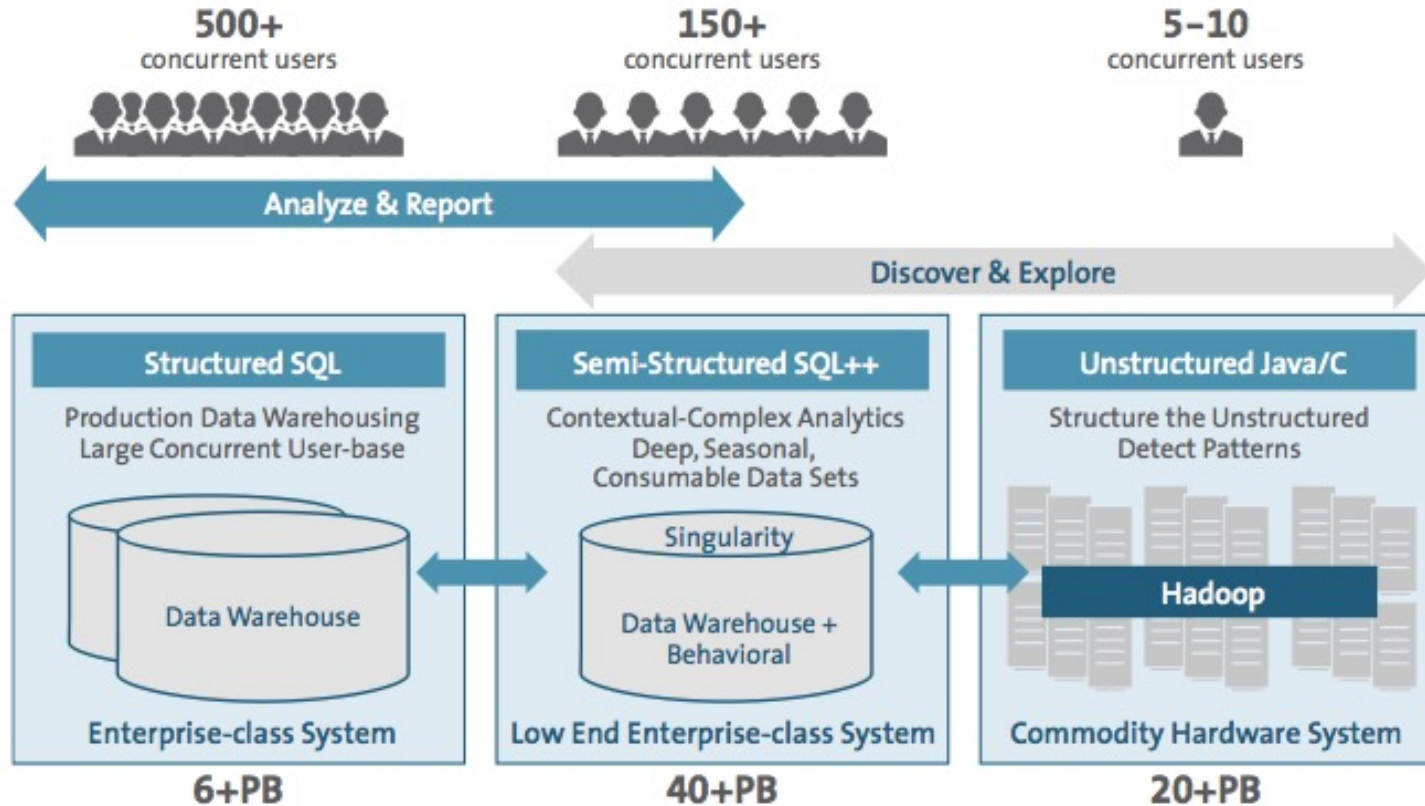
# Schema on-write vs. Schema on-read

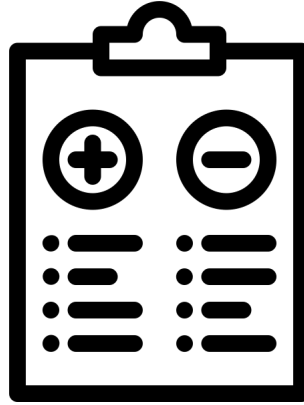
Attribute	Data Warehousing	Data Lake
Workload	<ul style="list-style-type: none"><li>• Hundreds to thousands of concurrent users</li><li>• Performing online (interactive) analytics</li><li>• Advanced workload management capabilities</li><li>• Batch processing</li></ul>	<ul style="list-style-type: none"><li>• Batch processing of data at scale</li><li>• Currently improving its capabilities to support more interactive users</li></ul>
Schema	<ul style="list-style-type: none"><li>• Typically schema is defined before data is stored</li><li>• Requires work at the beginning of the process, but offers performance, security and integration</li></ul>	<ul style="list-style-type: none"><li>• Typically schema is defined after data is stored</li><li>• Offers extreme agility and ease of data capture, but requires work at the end of the process.</li><li>• Works well for data types where data value is not known</li></ul>
Scale	<ul style="list-style-type: none"><li>• Large data volumes at moderate cost</li></ul>	<ul style="list-style-type: none"><li>• Extreme data volumes at low cost</li></ul>

# Coexistence of Hadoop and DW



# Big Data & Hadoop at eBay





# Pro and Cons

# Advantages & Disadvantages of Data Lakes

## Advantages

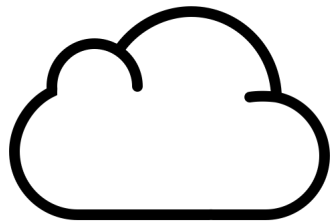
- Size
- Low cost
- Flexibility
- Easy accessibility
- No Data Silos
- “let’s store everything” strategy
- Self-Service Analytics
- Data Catalog of all data

## Disadvantages

Often Problems with

- Data Quality
- Data Lineage (origin of the data)
- Missing Metadata
- Different Schemas
- Not Integrated
- Risk of a “data graveyard” / “Data Swamp”

(However, Data Catalog and ELT can help with lineage, meta data and quality)

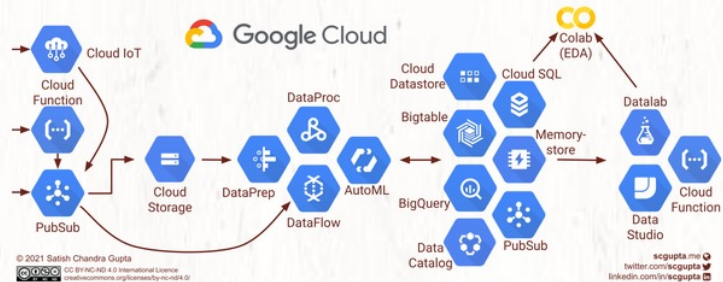
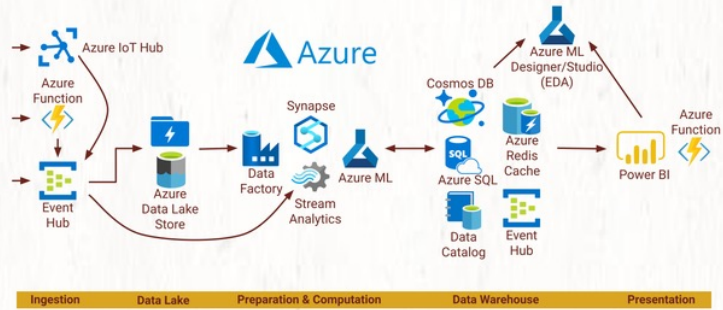
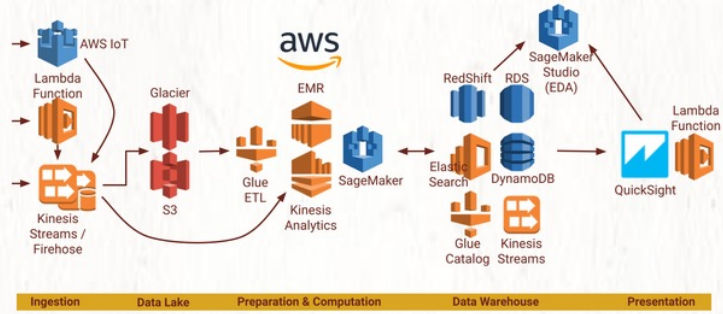


# Data Lake in the Cloud



# Big Data Pipelines on AWS, Microsoft Azure, and GCP

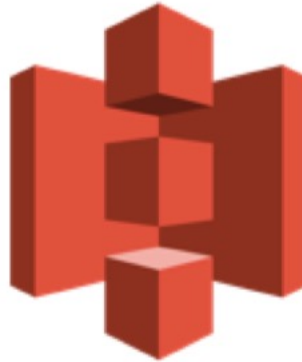
[scgupta.link/big-data-pipeline](https://scgupta.link/big-data-pipeline)



# Amazon AWS Data Lake Offerings



EC2 – Elastic, On-Demand Cluster



S3 – Unlimited Storage



EMR – Scale-Out Computing

# Microsoft Azure Data Lake

