

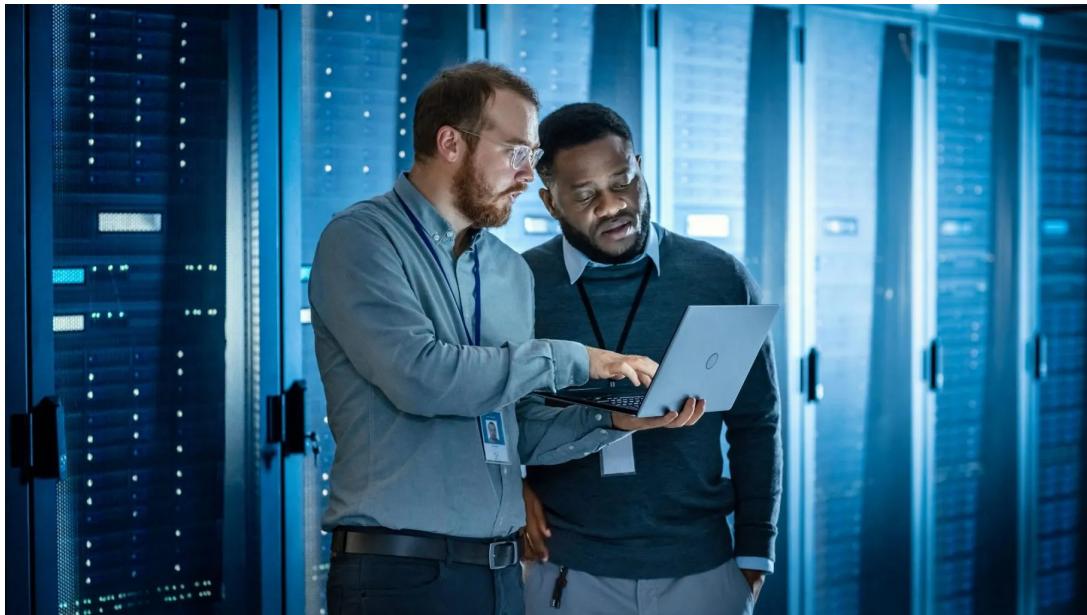
A photograph of two people, a man and a woman, looking at a large screen displaying various data visualizations like charts and graphs. The scene is dimly lit with a warm orange glow.

Engenharia de Dados do Zero



Stack Academy

O que é Engenharia de dados



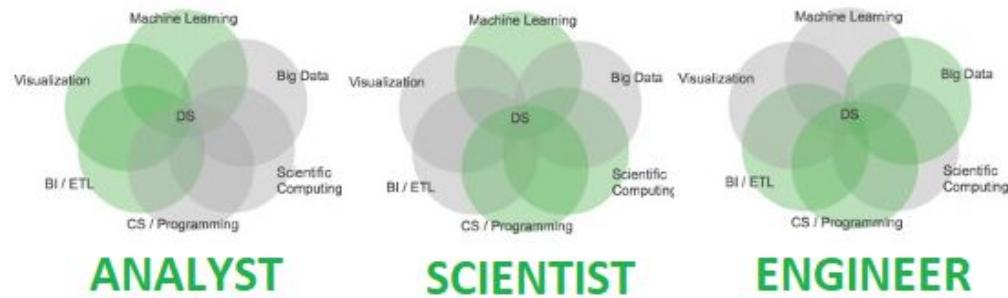
Para quem é esse curso?

1. Engenheiros de Dados.
2. Cientistas de Dados
3. Analista de Dados.
4. Gestores de Projeto.



Para quem é esse curso?

1. Engenheiros de Dados.
 2. Cientistas de Dados
 3. Analista de Dados.
 4. Gestores de Projeto.
-



Responsabilidades do Engenheiro de Dados



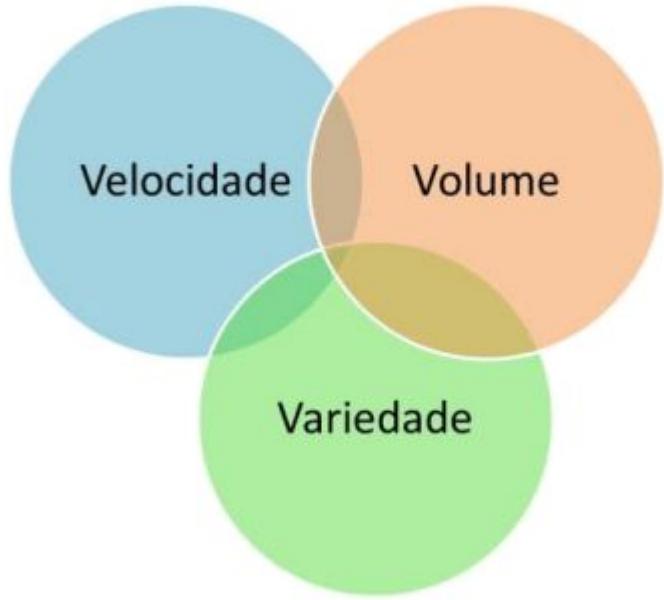
Responsabilidades do Cientista de Dados



O que é Big Data?



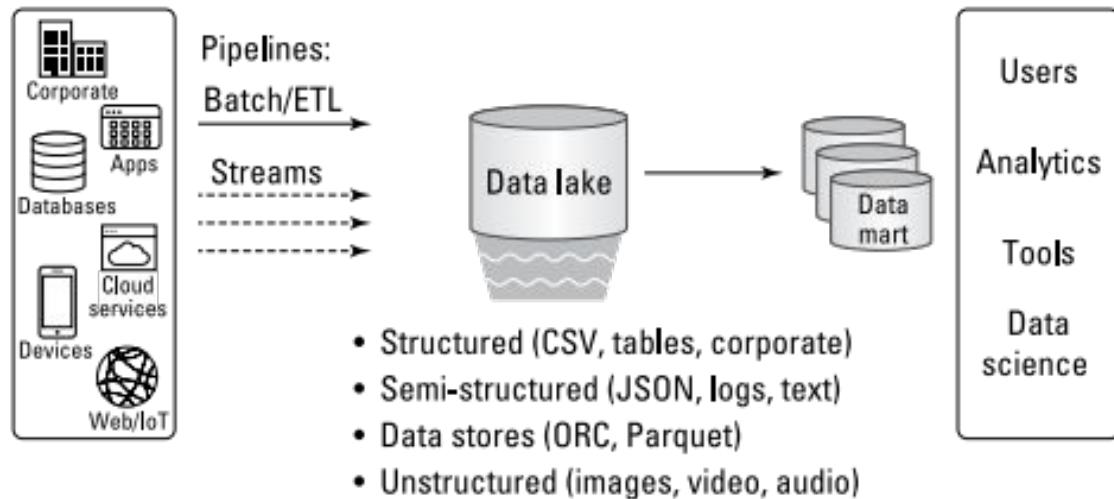
O que é Big Data?



O que são Data Lakes?

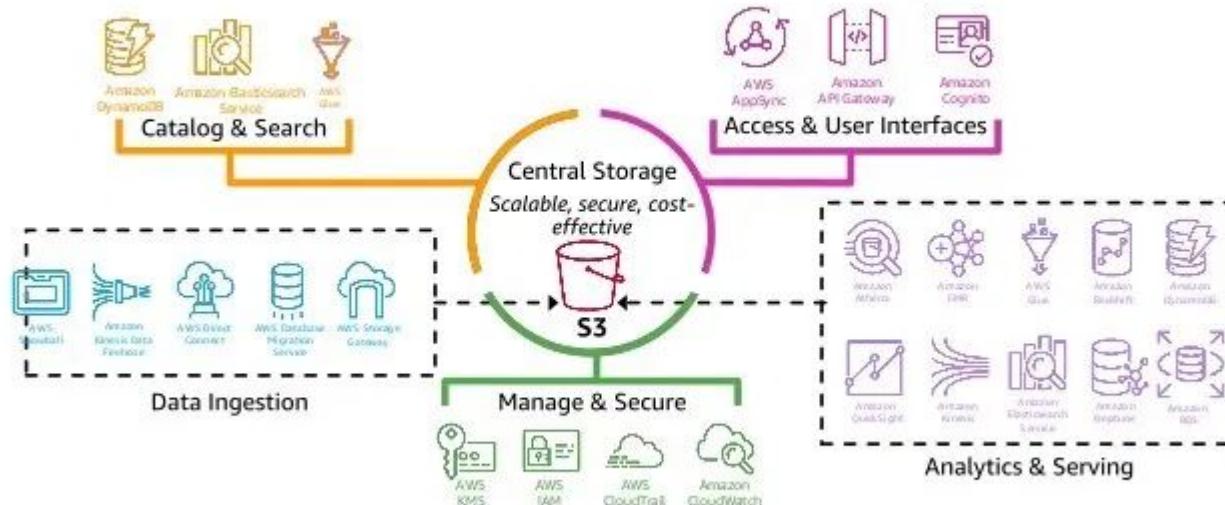
The Traditional Data Lake

Capture	Store	Deliver
batch or streaming data from many sources	single repository for all formats of data	rapid insights and business results from all data



Soluções de Data Lakes

Data lake on AWS

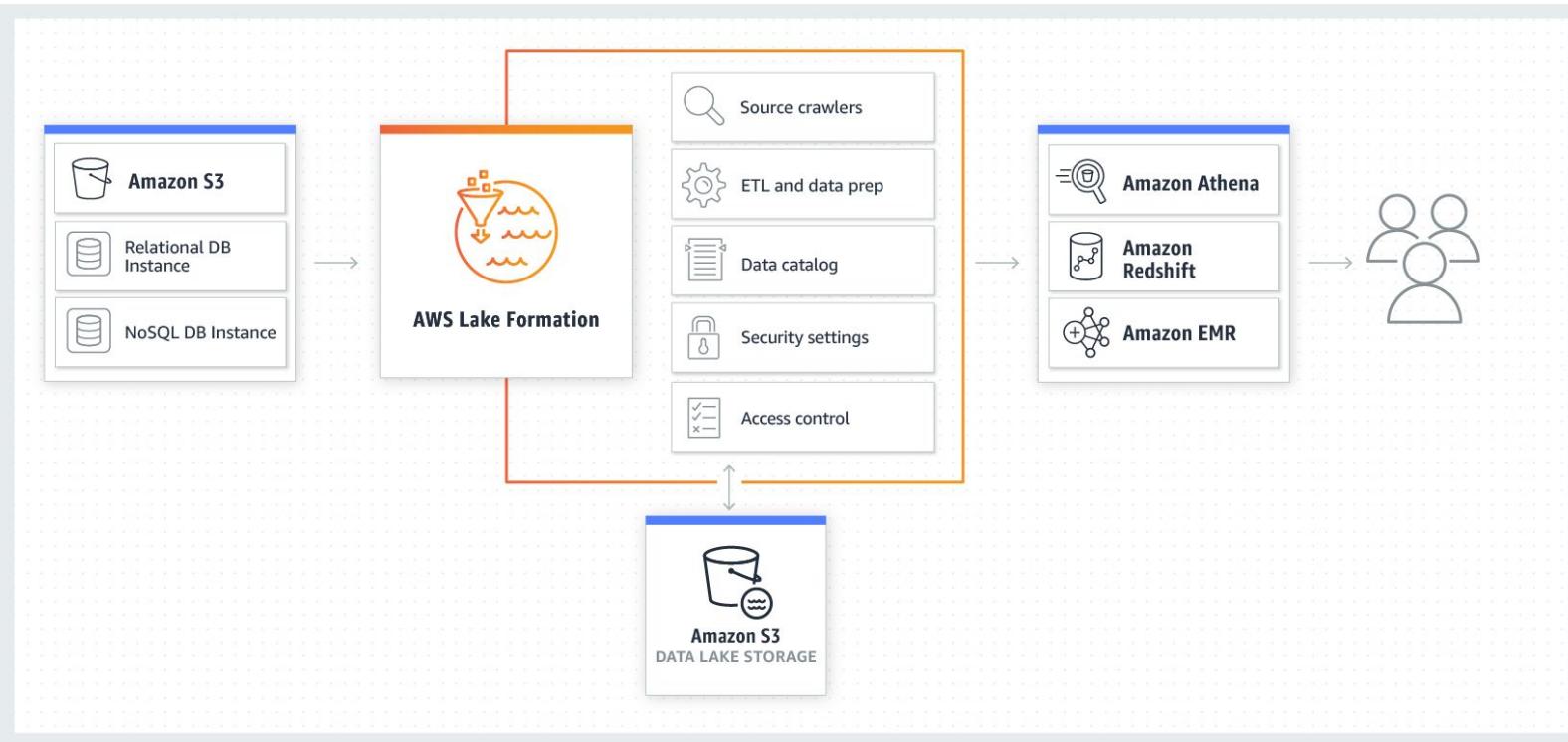


Invent

© 2018, Amazon Web Services, Inc. or its affiliates. All rights reserved.



Soluções de Data Lakes



Soluções de Data Lakes

Azure Data Lake



No limits Data Lake

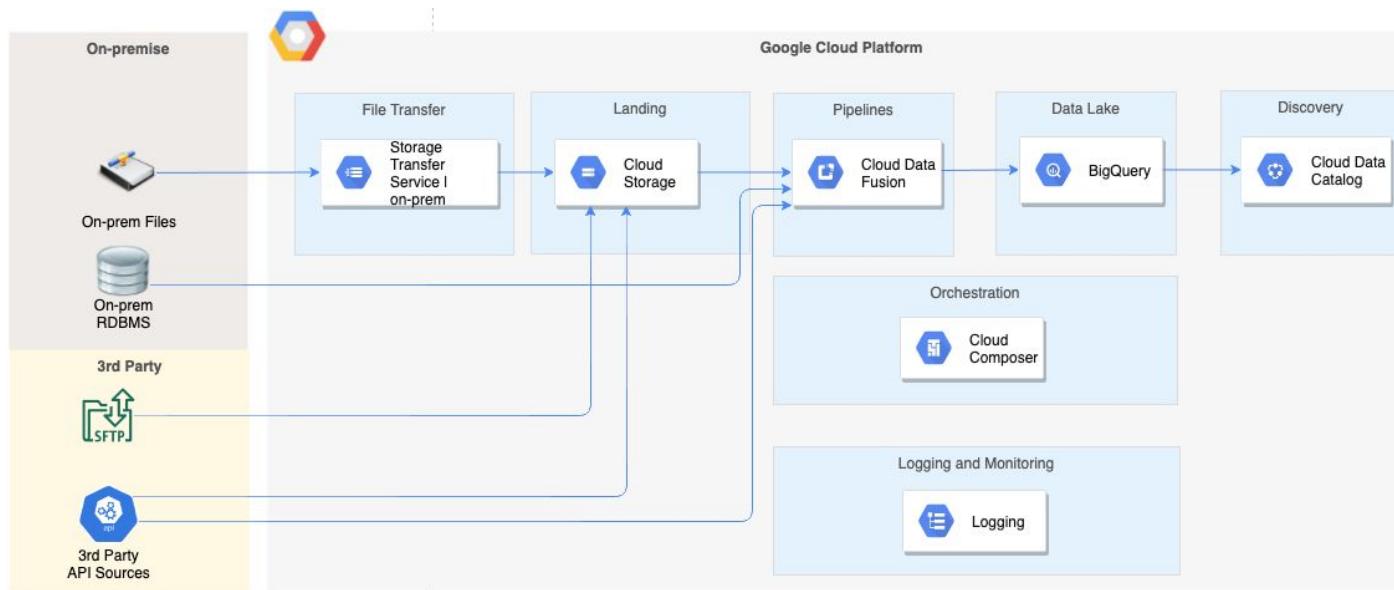


Analytics job service

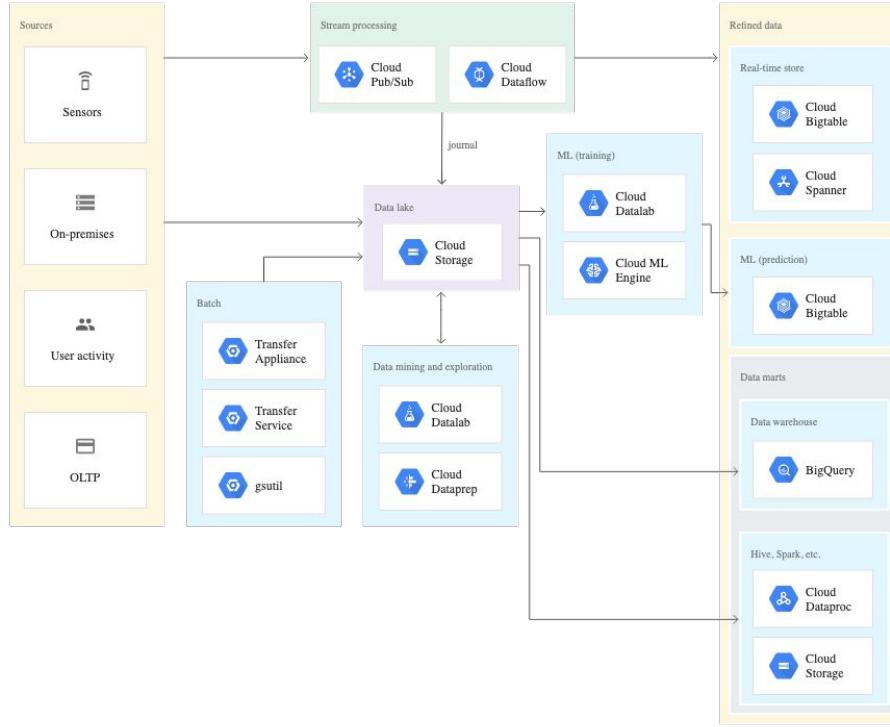


Managed Clusters

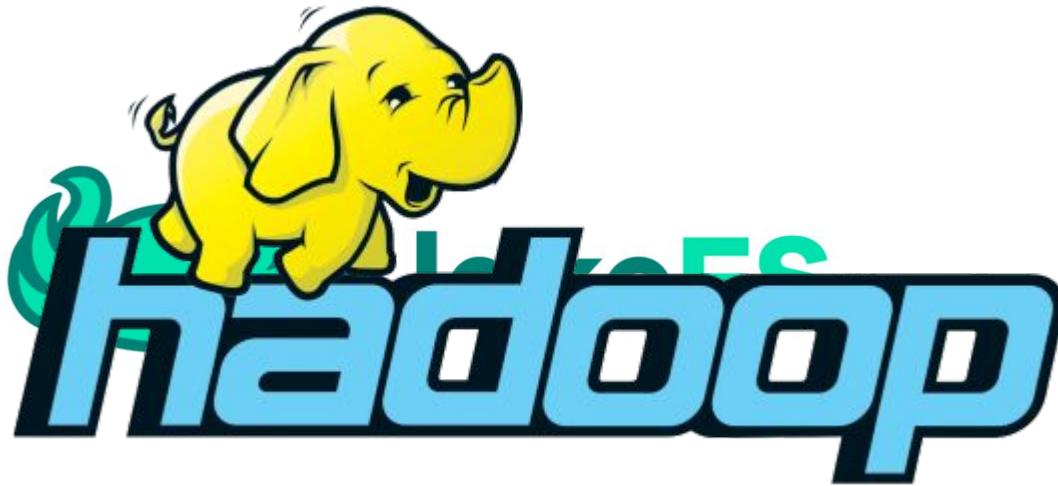
Soluções de Data Lakes



Soluções de Data Lakes



Soluções de Data Lakes Open Source



Data Swamp

Data Swamp



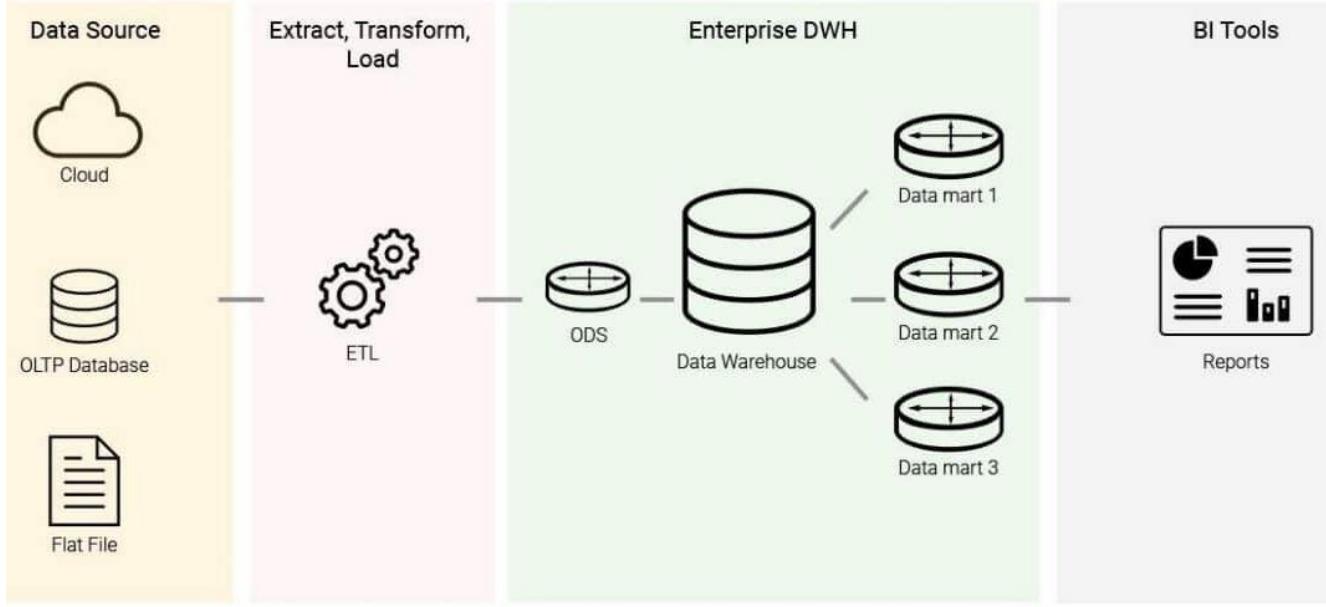
- ✗ No metadata
- ✗ Broken metadata management
- ✗ No data governance
- ✗ Broken ingestion process

Data Lake

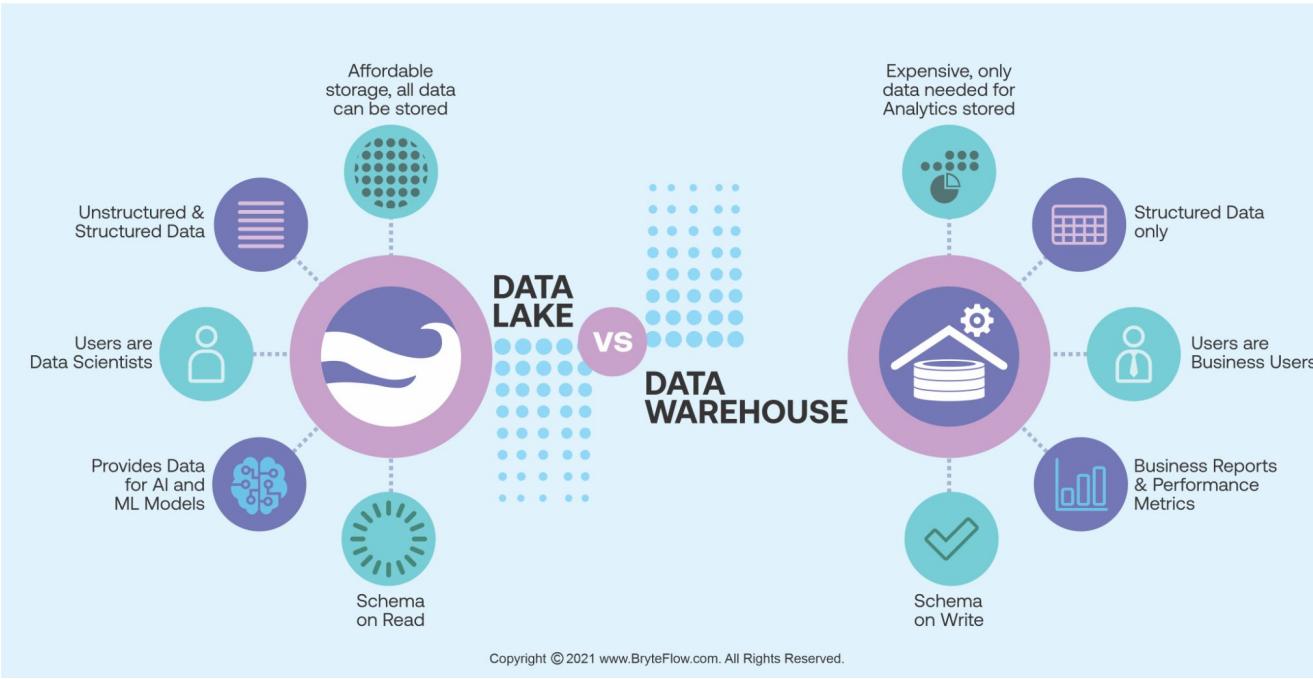


- ✓ Metadata
- ✓ Information is in rows and columns
- ✓ Easily ordered and processed with data mining tools
- ✓ Has data context
- ✓ Contains a data set for running analytics
- ✓ Has directories and sub-directories

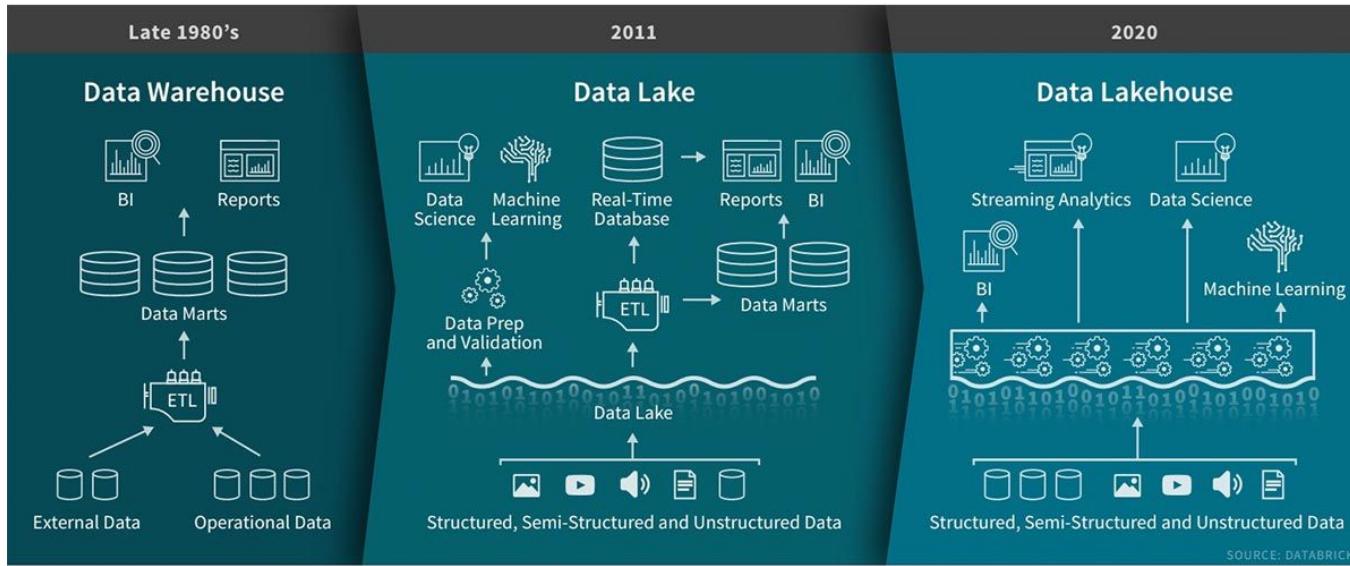
O que são Data Warehouses?



Data Warehouses Vs Data Lakes

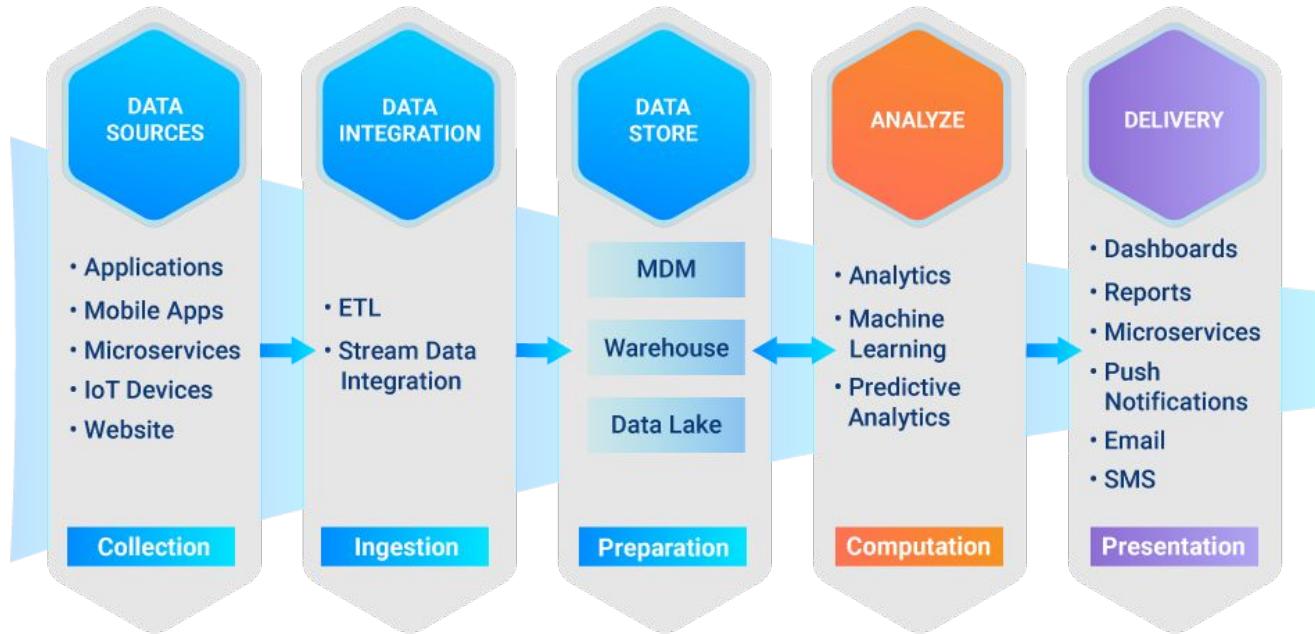


Data Warehouse vs Data Lake vs Data Lakehouse

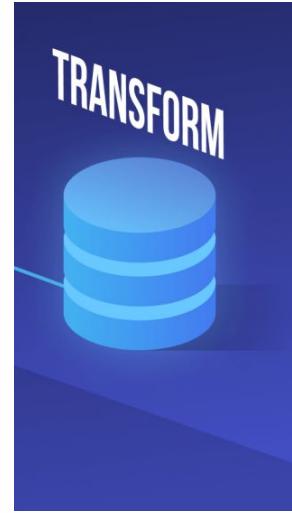
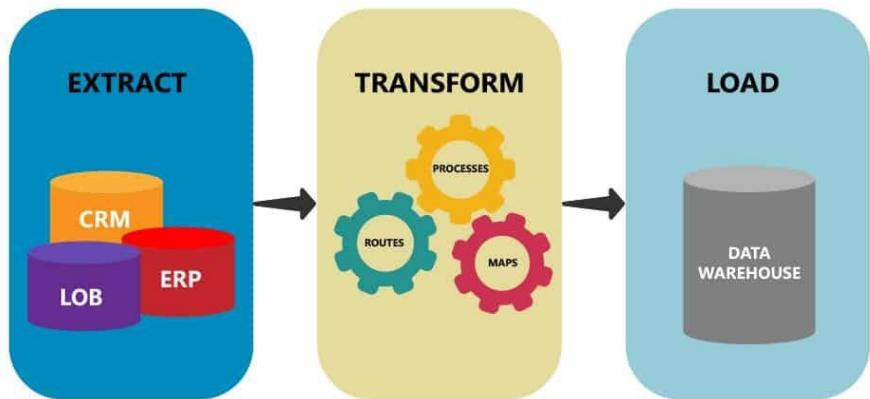


SOURCE: DATABRICKS

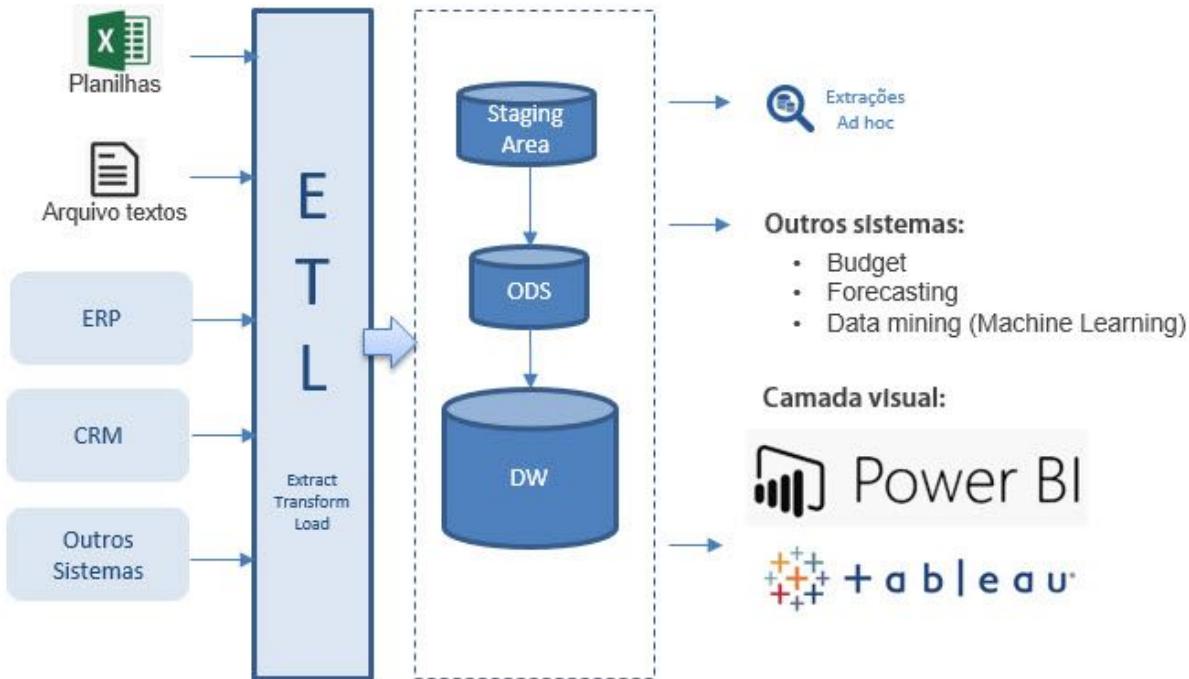
BIG DATA PIPELINE



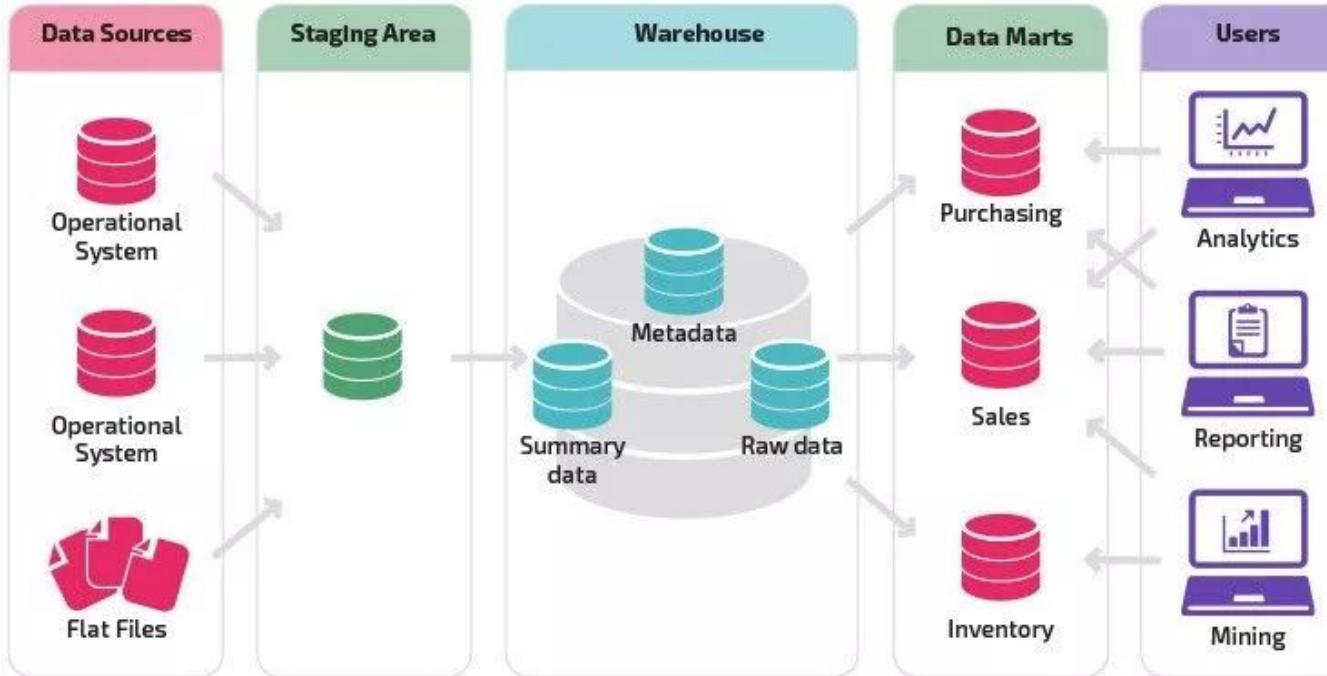
ETL vs ELT



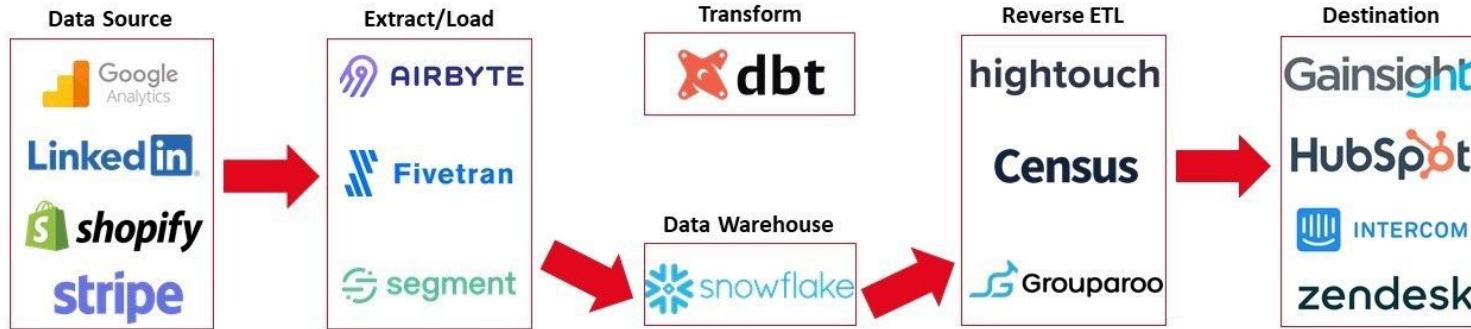
Extract Transform Loading



Extract Transform Loading - Data Marts

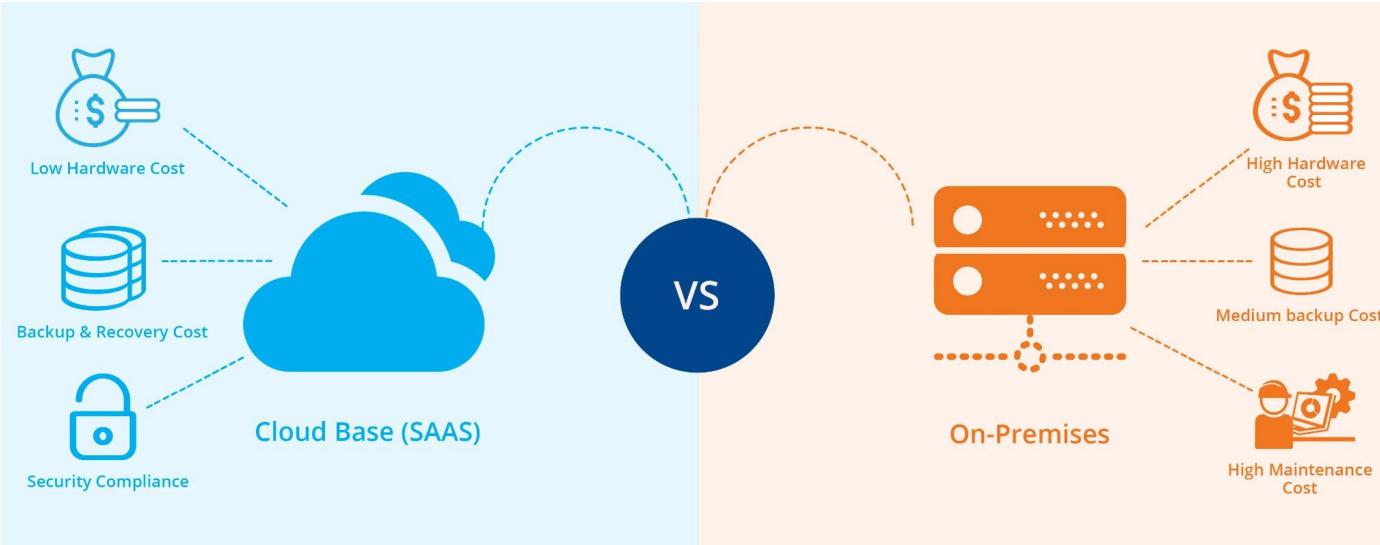


Soluções para ETL's



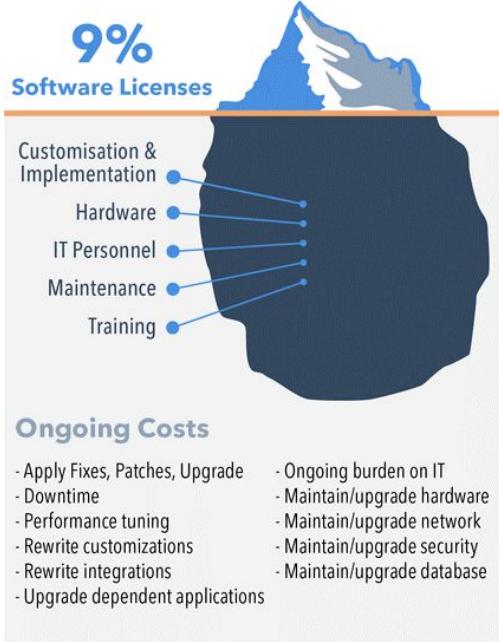
@AstasiaMyers

On-Premises vs Cloud



On-Premises vs Cloud

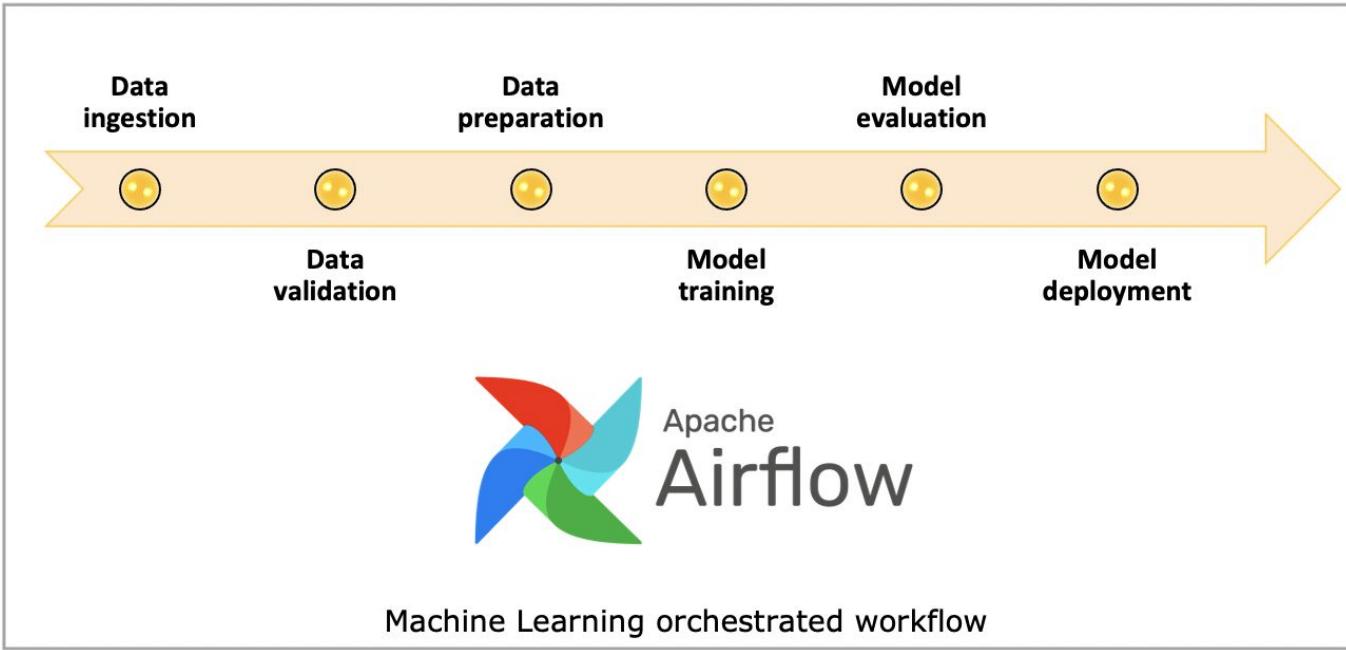
On-Premises



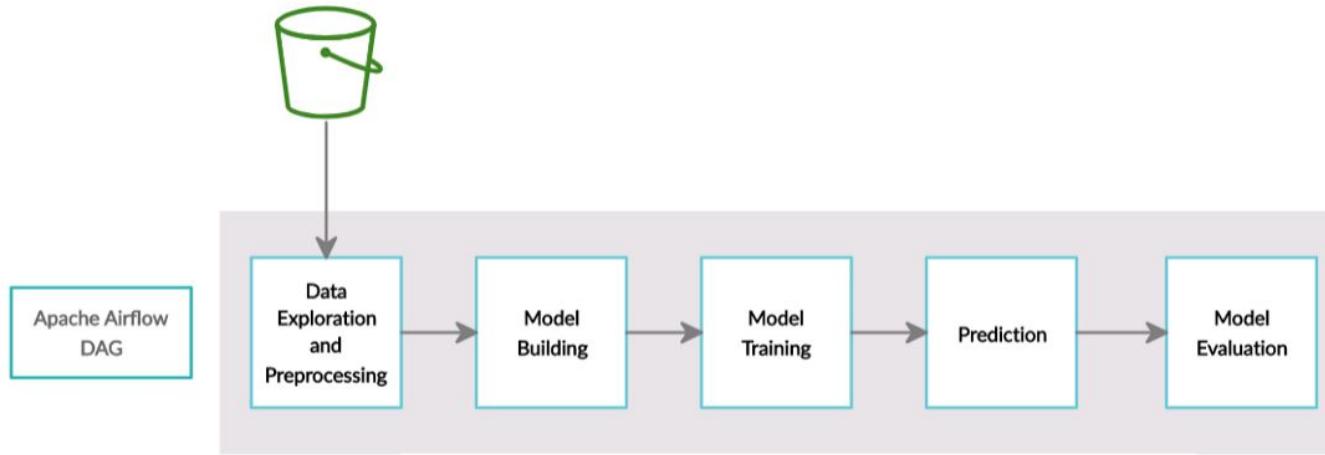
Cloud Computing



Soluções para Data Pipelines



Soluções para Data Pipelines

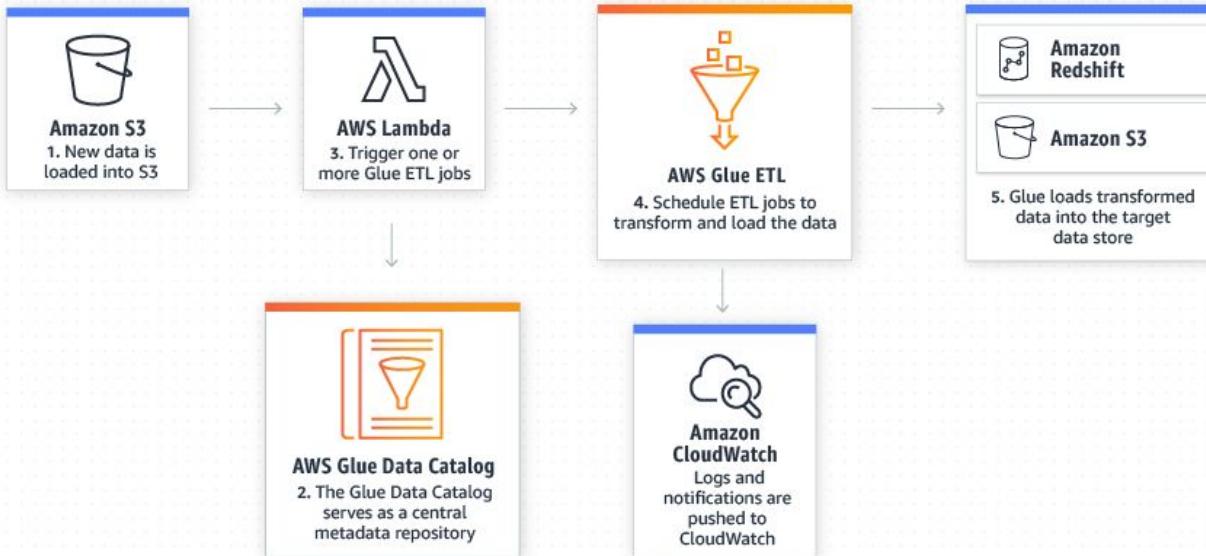


Soluções para Data Pipelines



Azure Data Factory

Soluções para Data Pipelines



Soluções para Data Pipelines



Google
Cloud Composer

Responsabilidades do Engenheiro de Dados para gerenciar um Data Pipeline

- Definir ambiente e arquitetura.
 - On-Premisse, On Cloud
- Especificar tecnologias para armazenamento.
 - Data Lakes, Data Warehouse.
- Definir formas de processamento.
 - Streaming, Clusters.
- Definir e implementar padrões e segurança e integração com outros sistemas.
- Monitoramento.

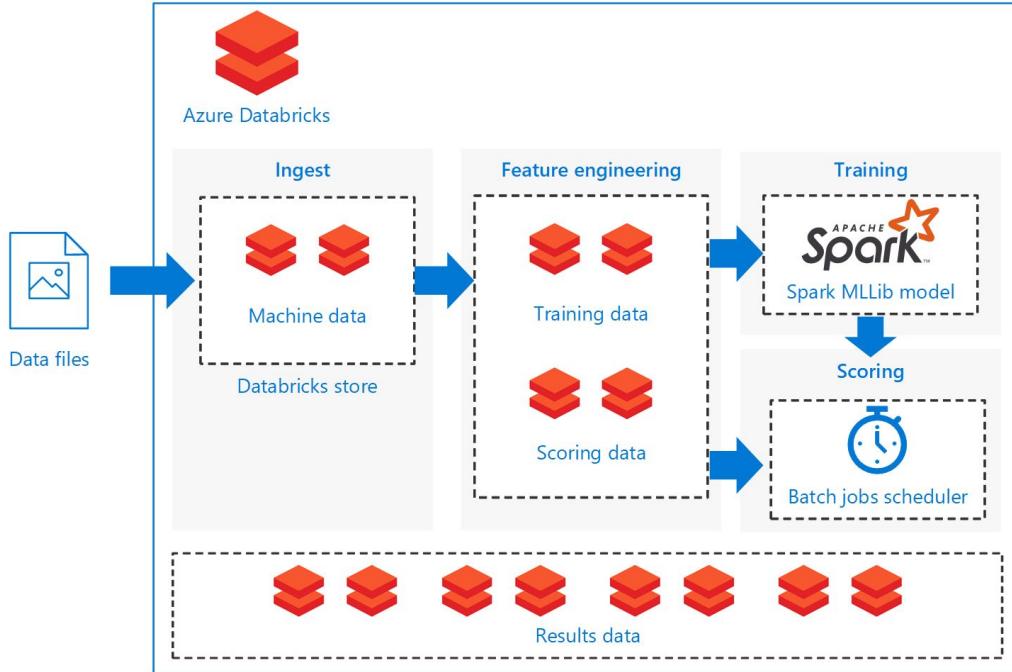


Data Engineering at World



databricks

Data Engineering at World

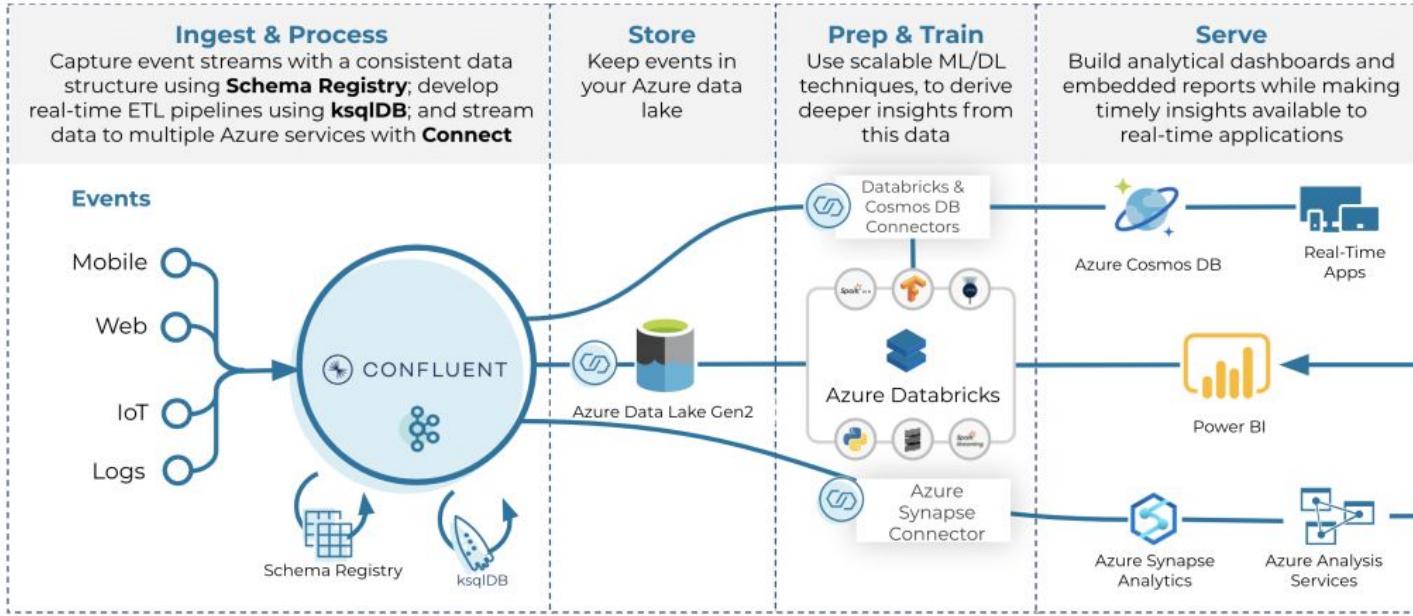


Data Engineering at World



CONFLUENT

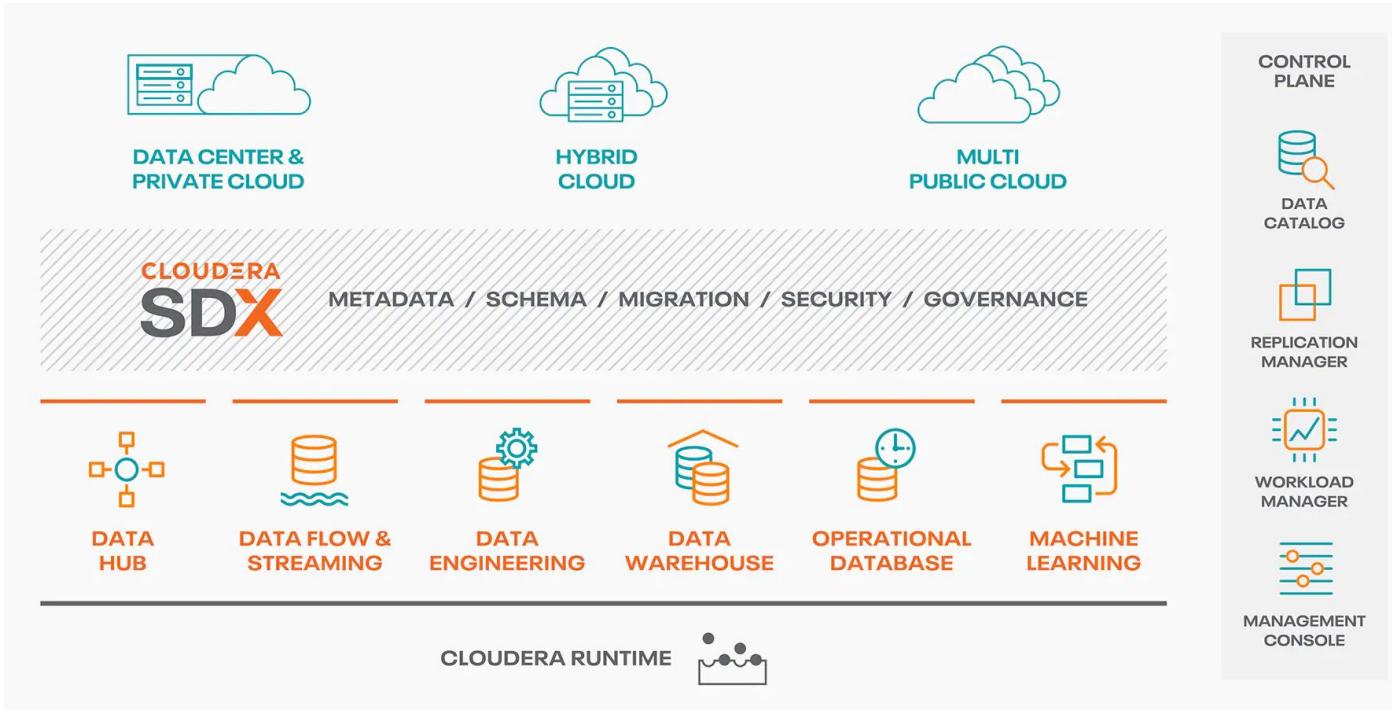
Data Engineering at World



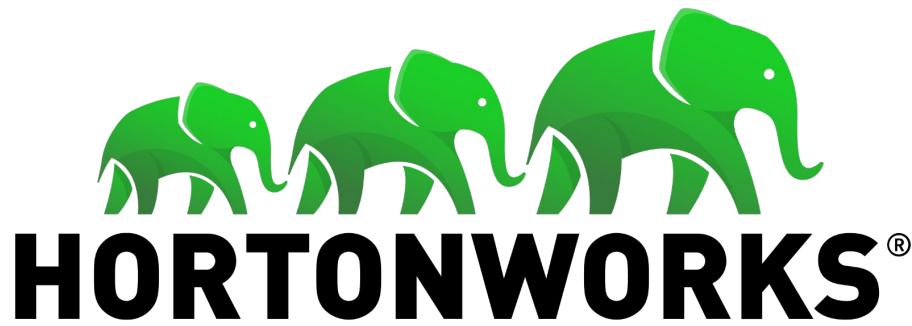
Data Engineering at World

CLOUDERA

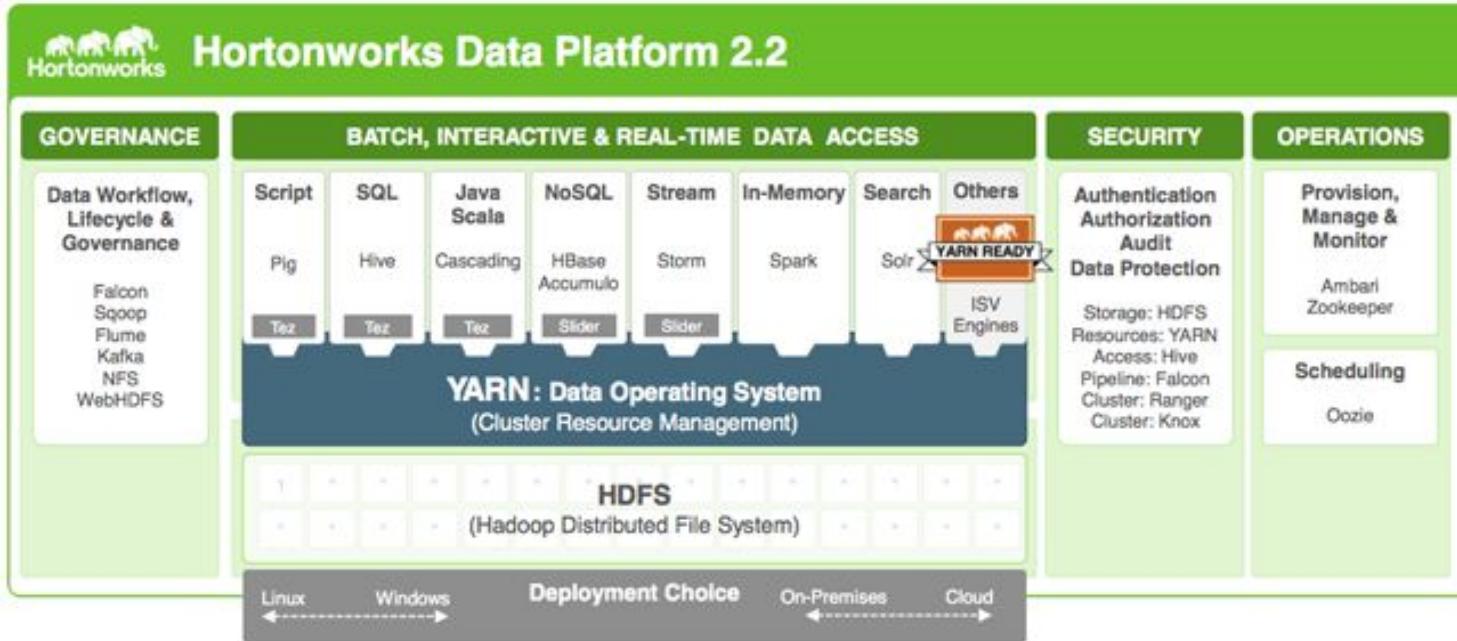
Data Engineering at World



Data Engineering at World



Data Engineering at World



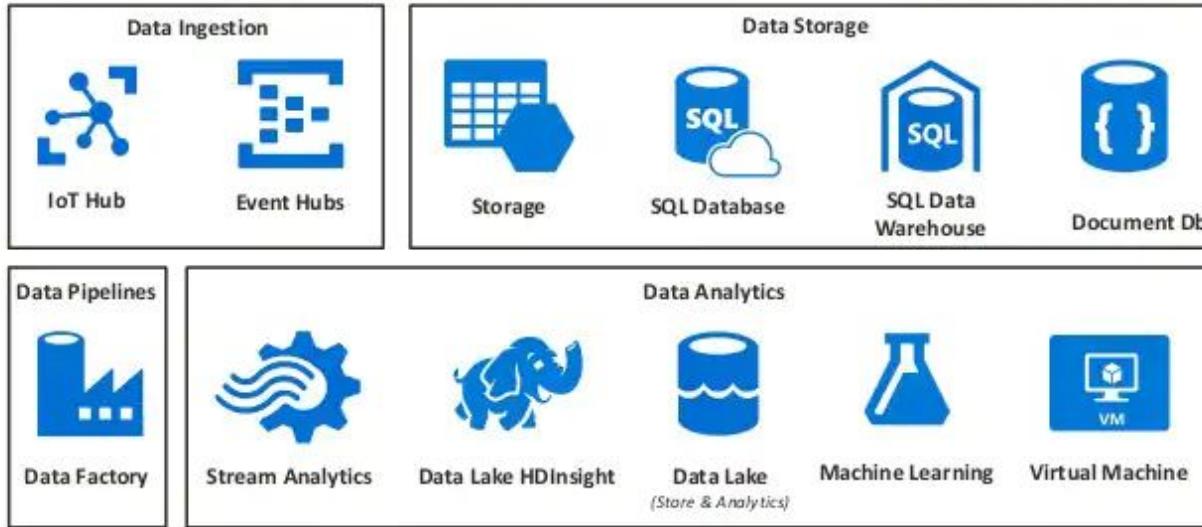
Data Engineering at World



Microsoft Azure

Data Engineering at World

Overview in Azure



codit

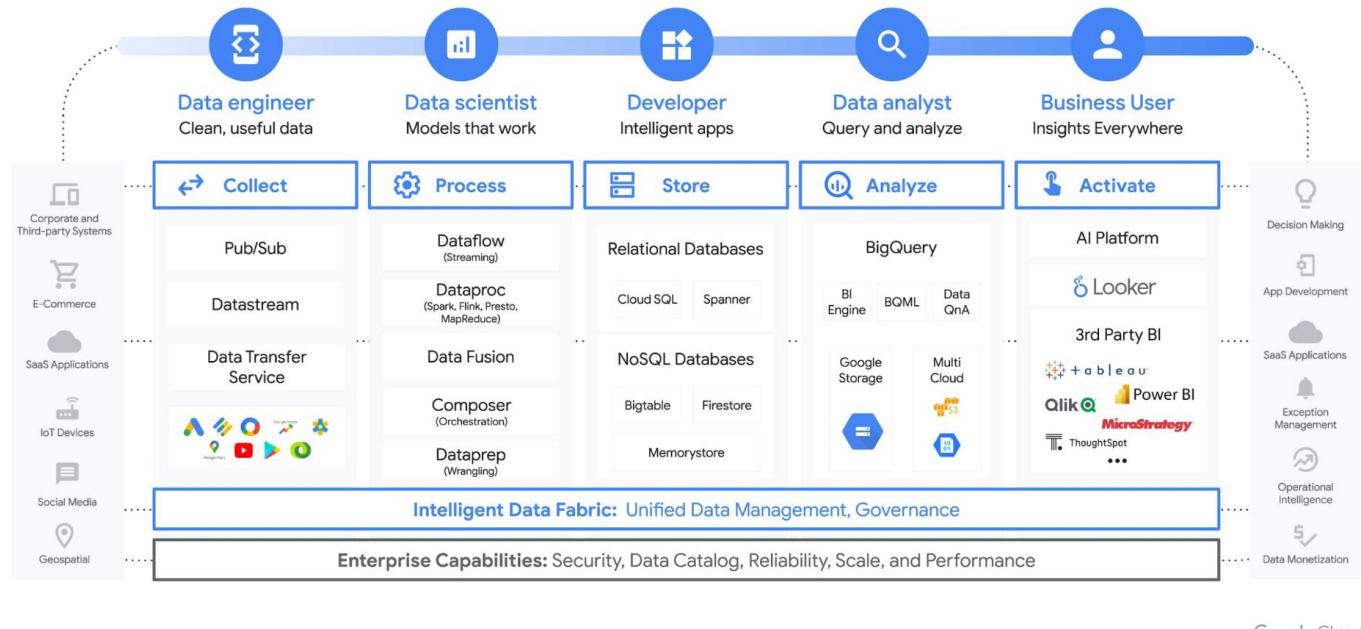
13

Data Engineering at World



Google Cloud Platform

Data Engineering at World



Data Engineering at World



Data Engineering at World

Data Processing



Amazon EMR
Managed Hadoop Applications



Amazon Redshift
Petabyte-scale Data Warehousing



Amazon Athena
Interactive Query

Serverless Compute



Amazon Kinesis Firehose
Real-Time Data Streaming



AWS Glue
ETL & Data Catalog



Amazon Redshift Spectrum
Fast @ Exabyte scale



AWS Lambda
Trigger-based Code Execution

Storage

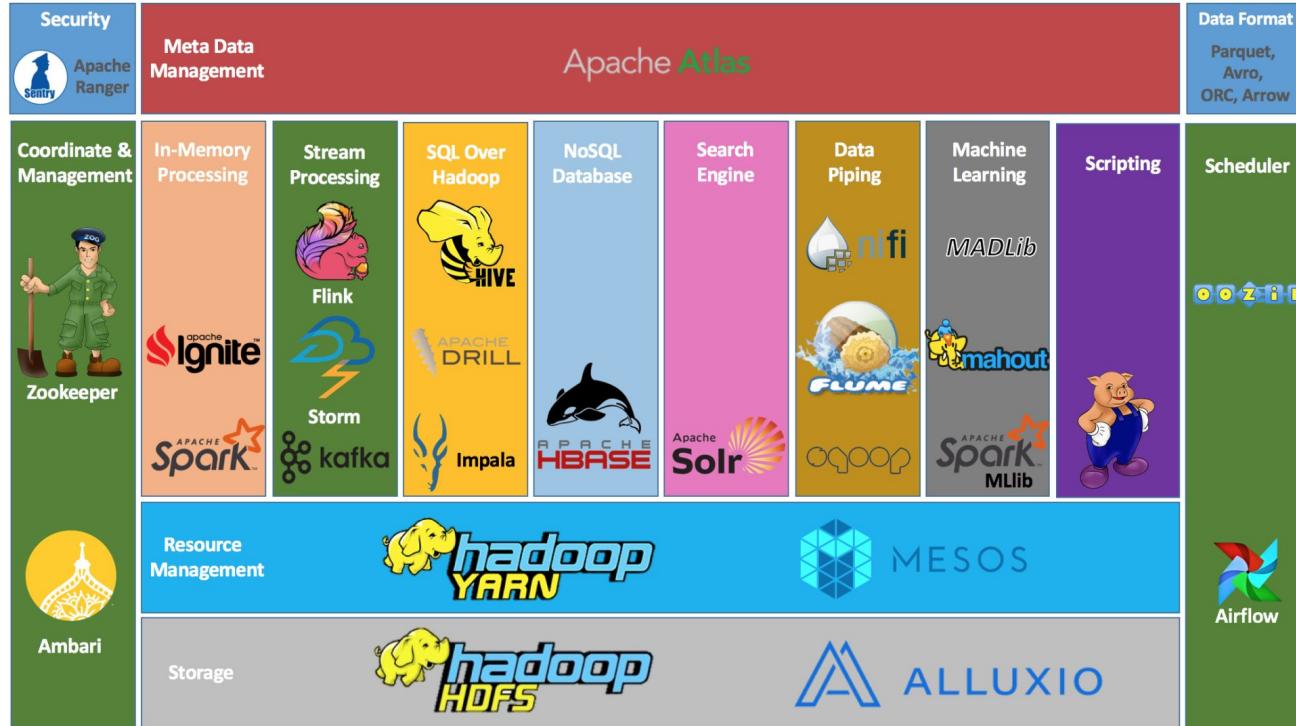


Amazon S3
Exabyte-scale Object Storage



AWS Glue Data Catalog
Hive-compatible Metastore

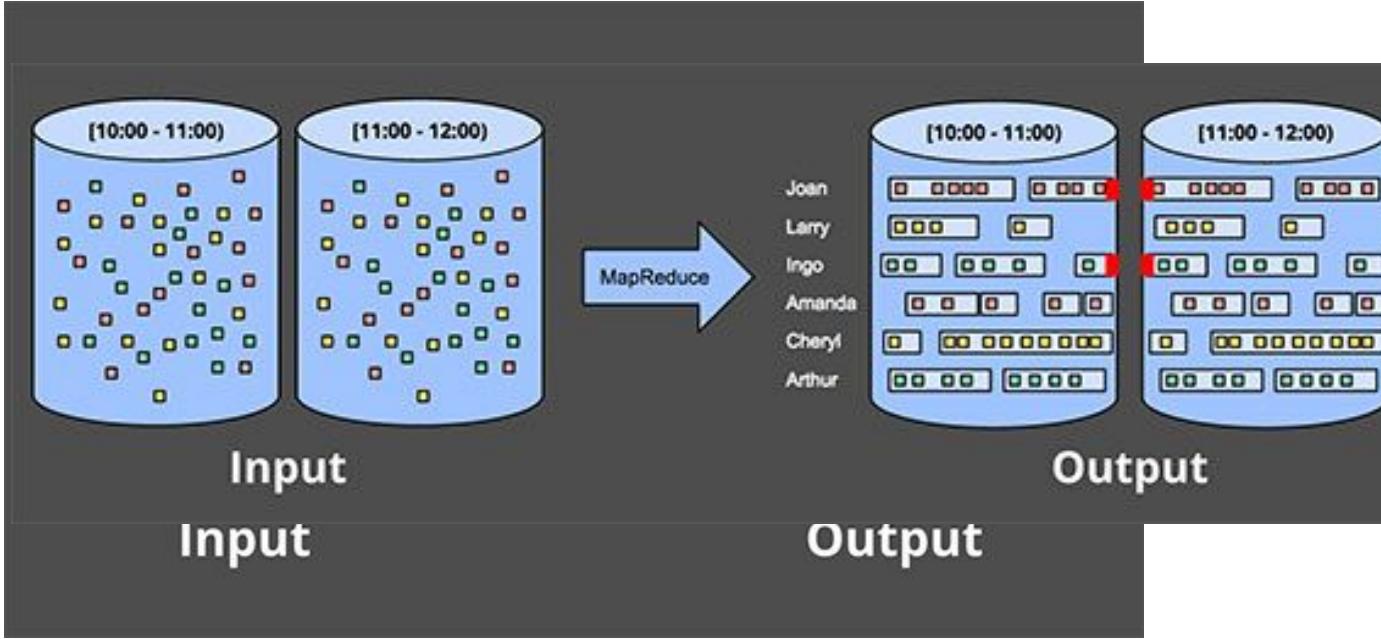
Ecossistema Hadoop



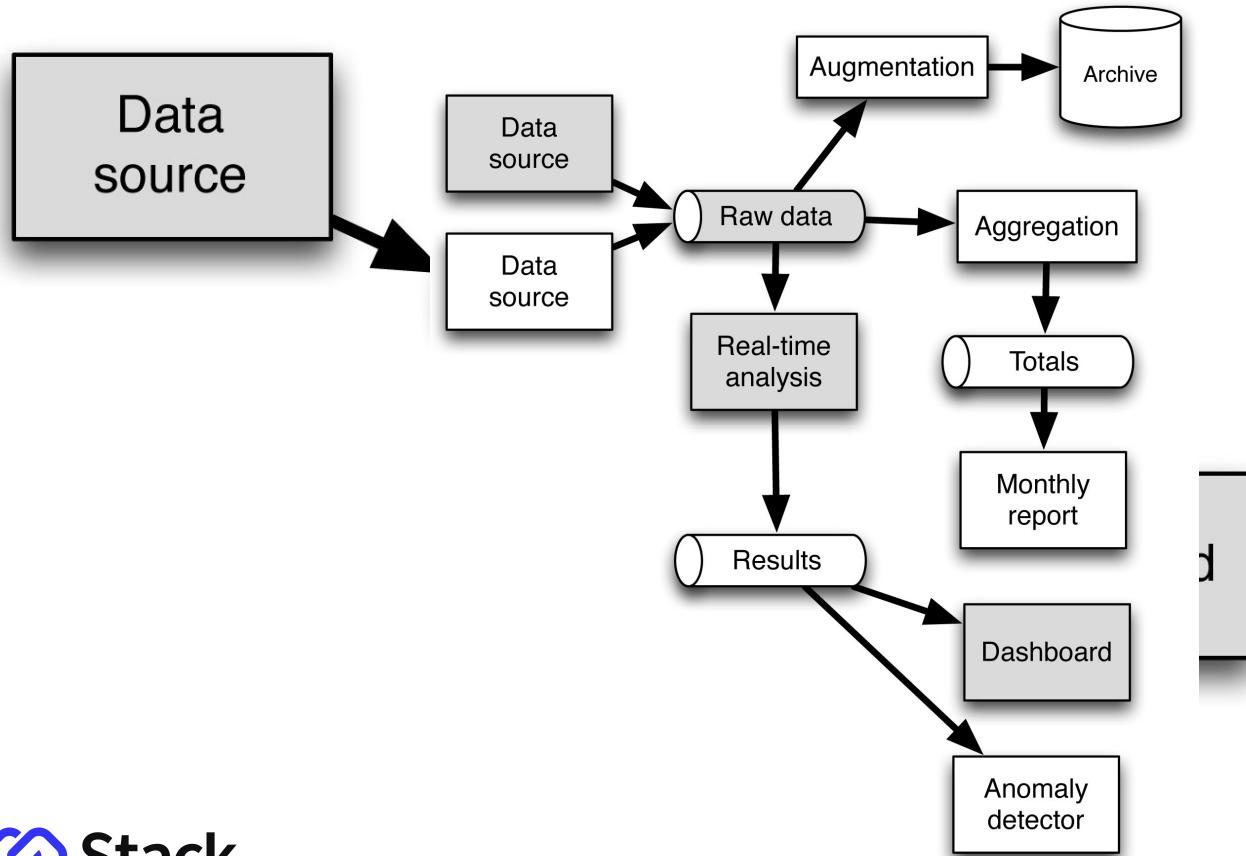
Processamento Batch



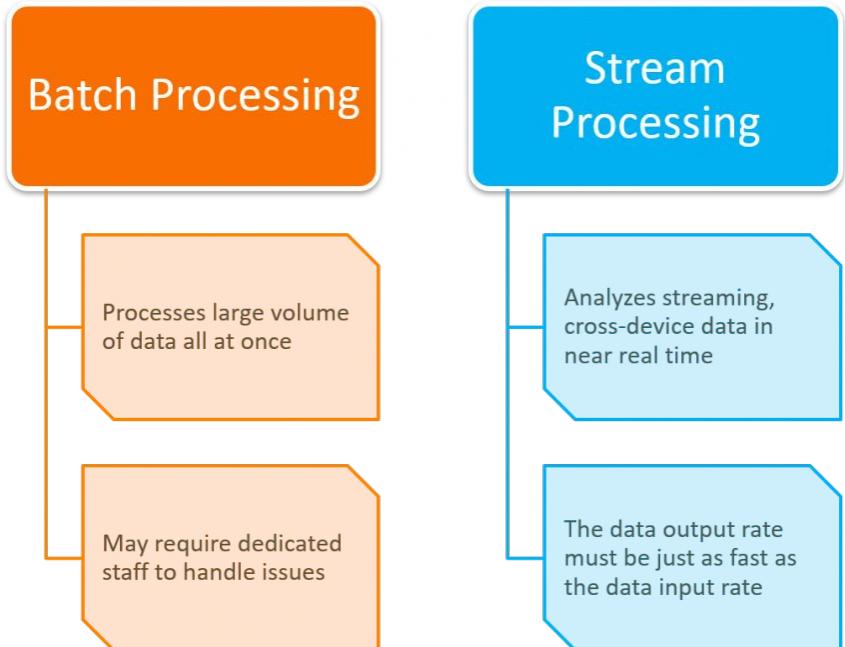
Processamento Streaming



Processamento Streaming



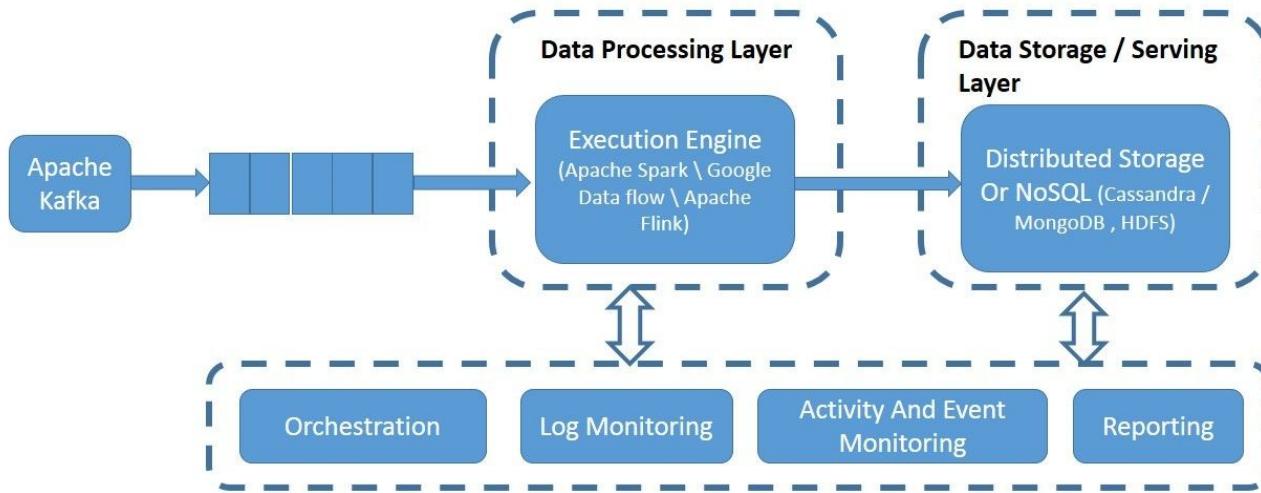
Batch Processing vs Streaming Processing



Benefits and drawbacks of common data processing types

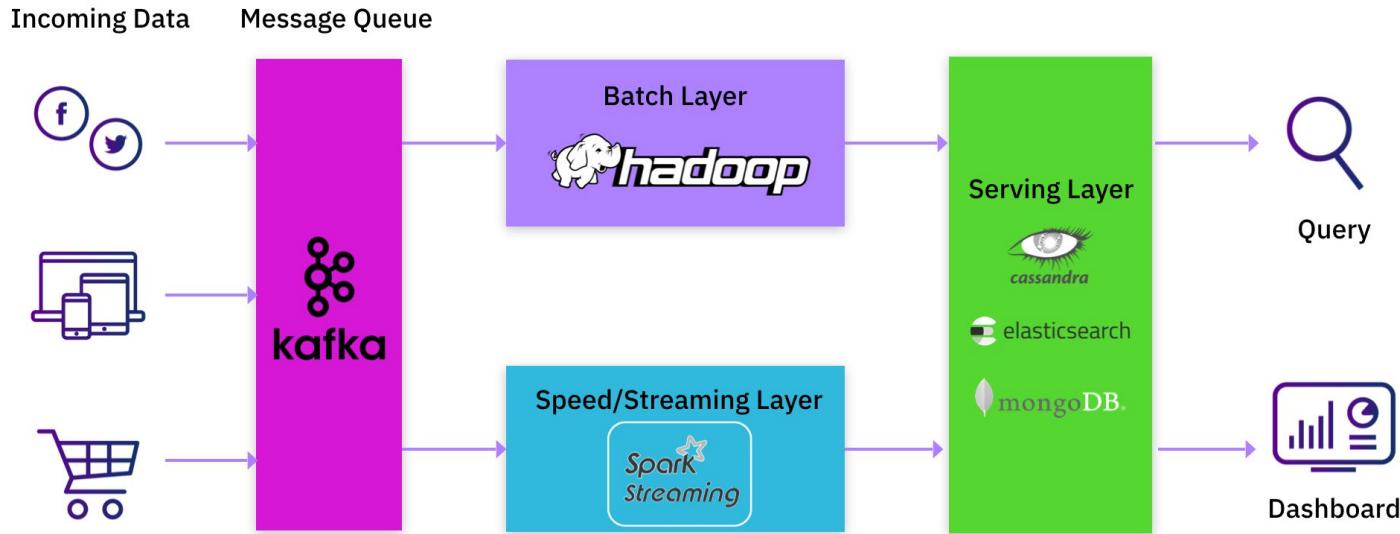
Arquitetura Kappa

Kappa Architecture



Siddharth Mittal

Arquitetura lambda



Lambda vs Kappa

		
Processing paradigm	Batch + Streaming	Streaming
Re-processing paradigm	Every Batch cycle	Only when code changes
Resource consumption	Function = Query (All data)	Incremental algorithms, running on deltas
Reliability	Batch is reliable, Streaming is approximate	Streaming with consistency (exactly once)

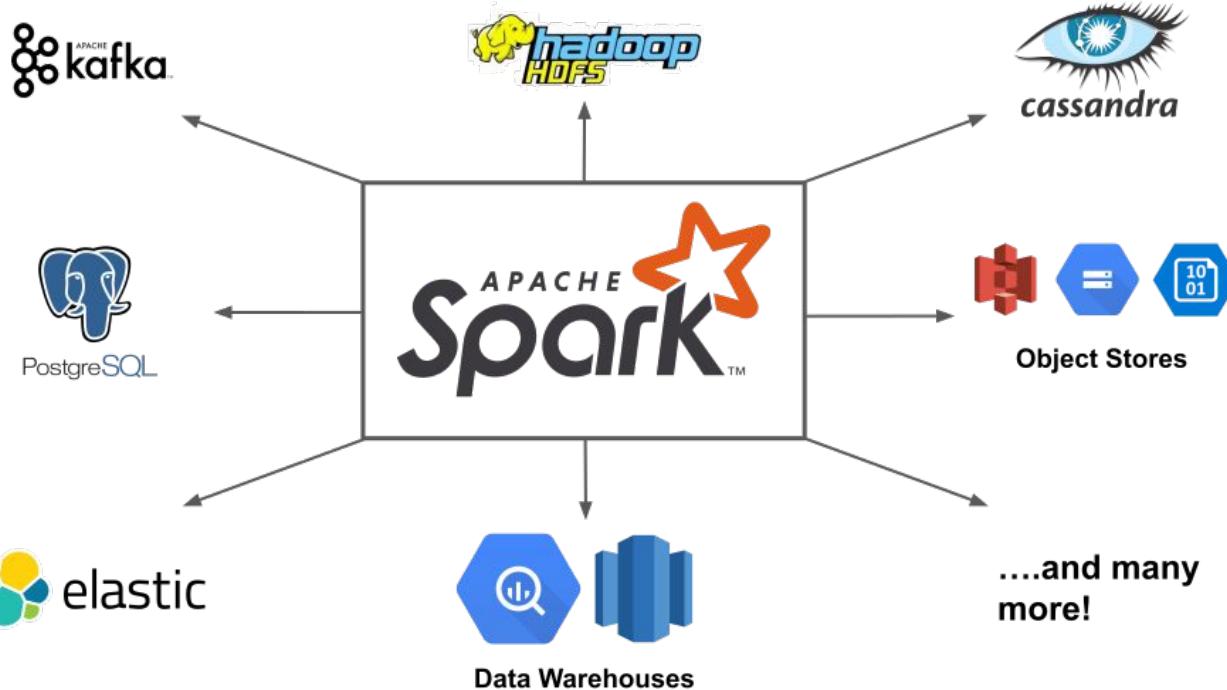
O que é o Apache Spark?



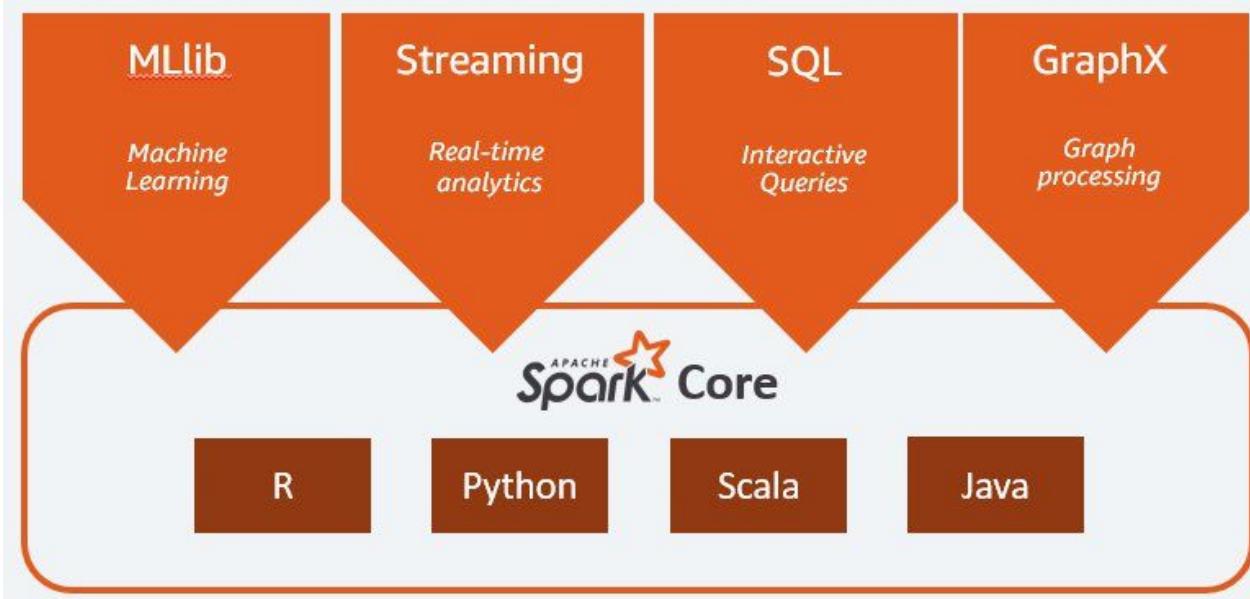
O que é o Cluster?



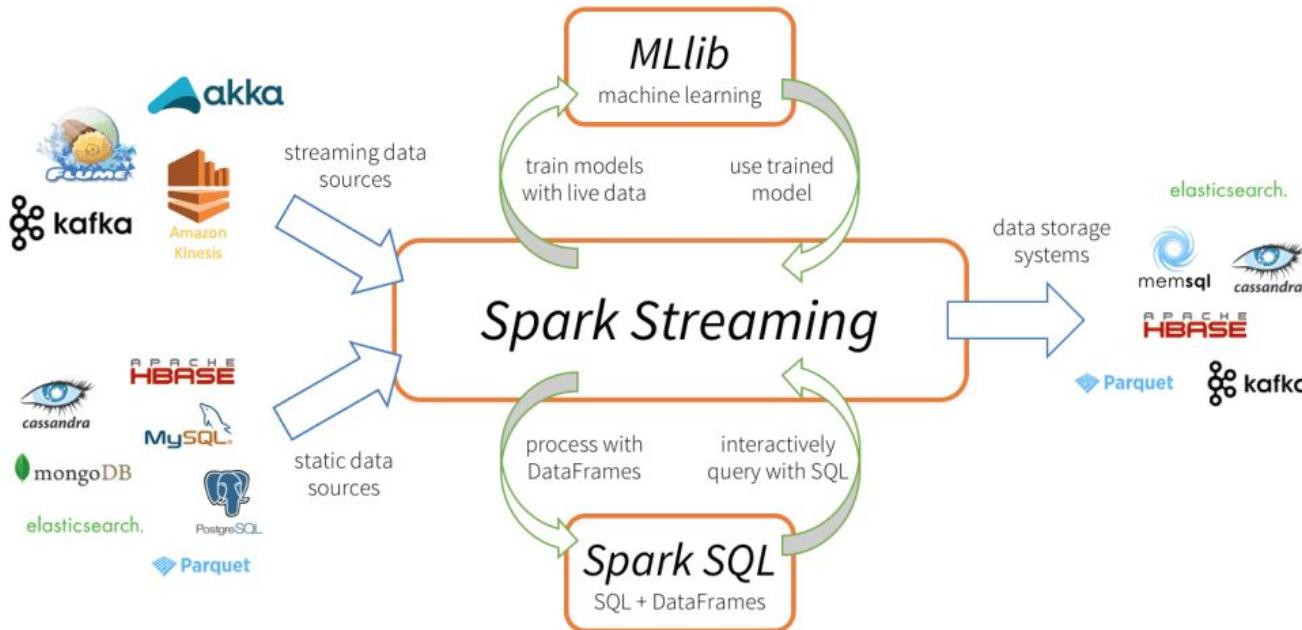
O que é o Apache Spark?



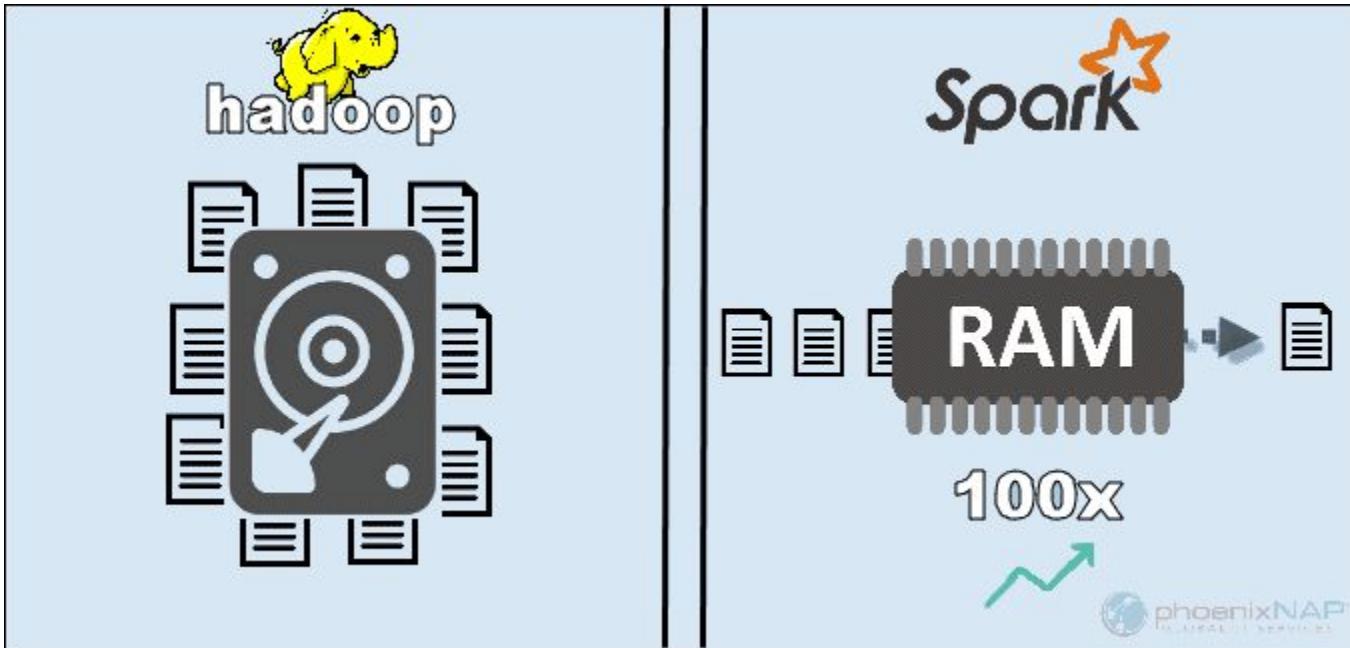
O que é o Apache Spark?



O que é o Apache Spark?



Spark Vs Hadoop



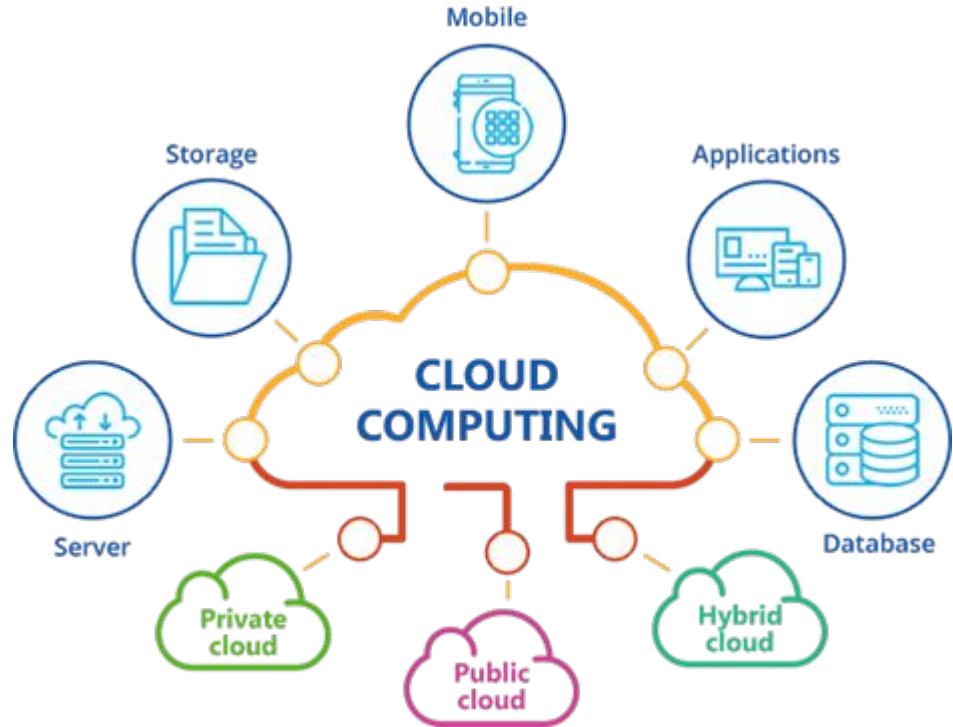
Spark Vs Hadoop

Factors	Spark	Hadoop MapReduce
Speed	100x times than MapReduce	Faster than traditional system
Written In	Scala	Java
Data Processing	Batch / real-time / iterative / interactive / graph	Batch processing
Ease of Use	Compact & easier than Hadoop	Complex & lengthy
Caching	Caches the data in-memory & enhances the system performance	Doesn't support caching of data

O que é Databricks?



Por que Databricks?



Por que Databricks?



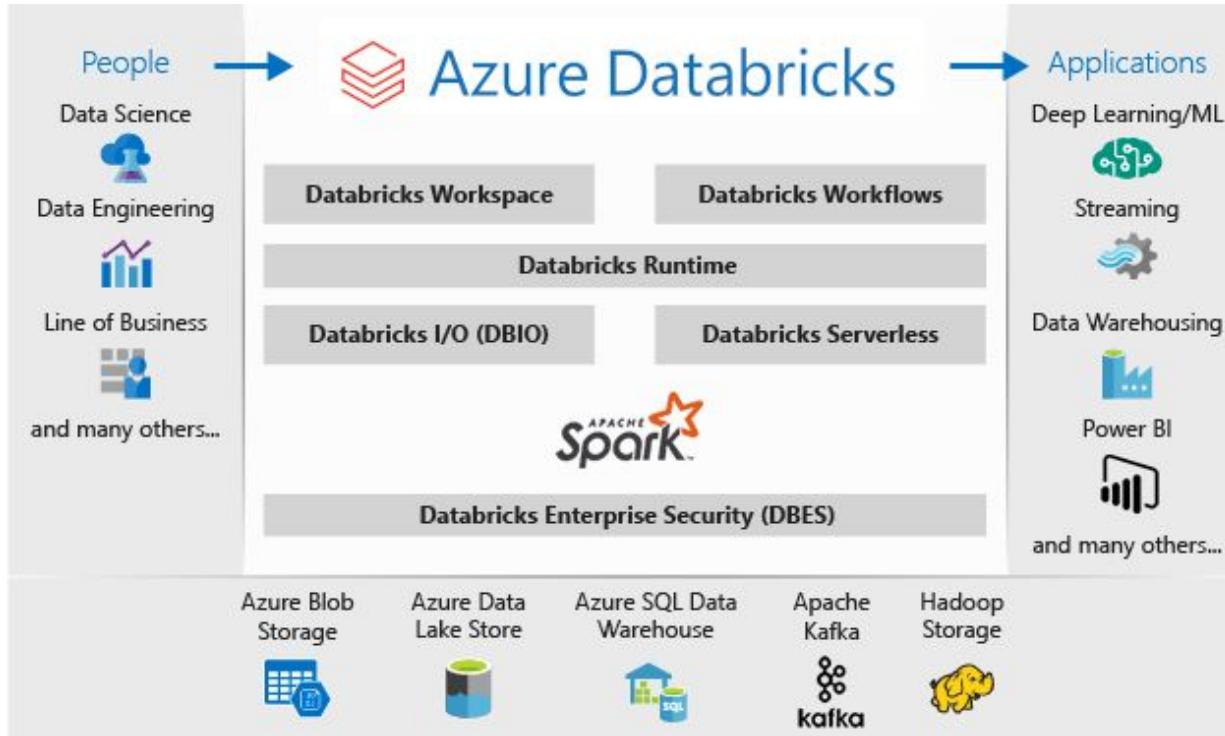
Google Cloud



databricks



Por que Databricks?



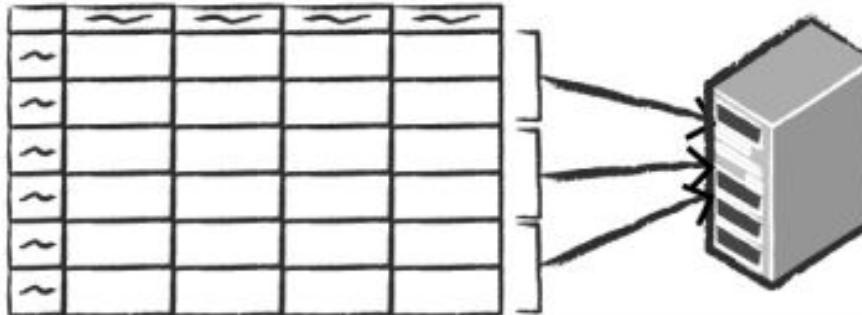
Hands on!

Dataframes e Partições no Spark

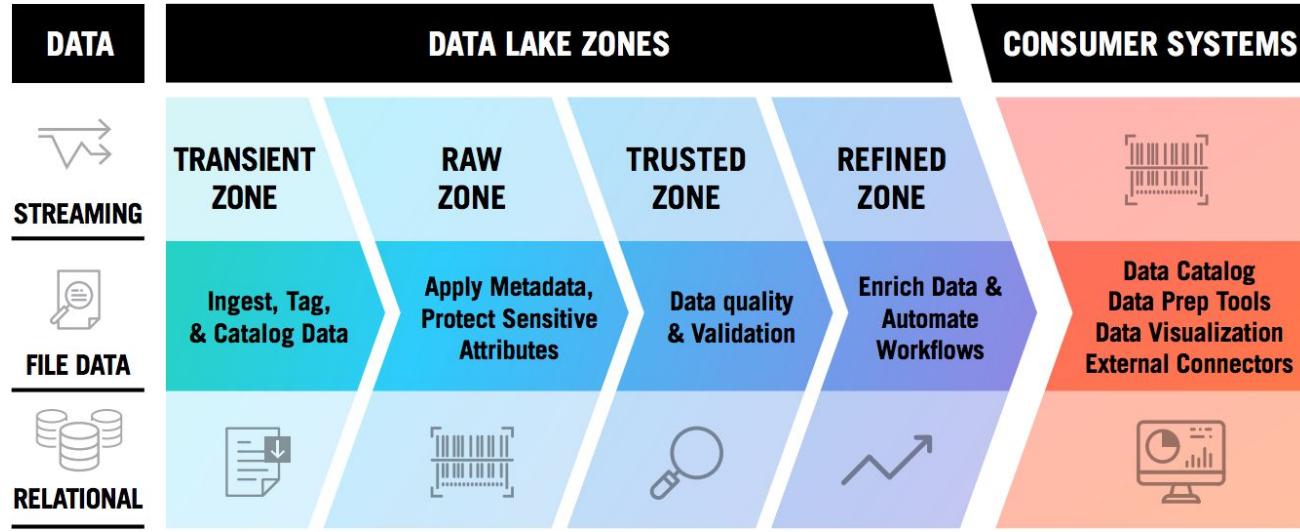
Spreadsheet on
a single machine



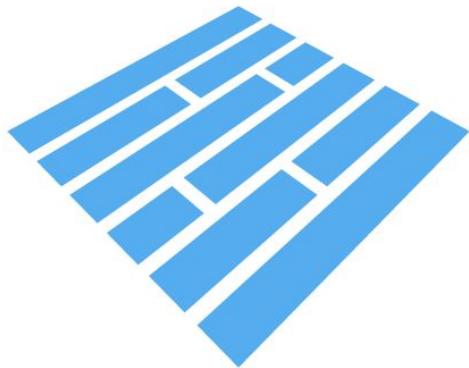
Table or Data Frame
partitioned across servers
in a data center



Dataframes e Partições no Spark



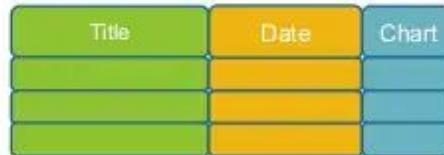
Apache Parquet



Parquet

Apache Parquet

Data In Columns On Disk



Row-Oriented data on disk

Led Zeppelin IV	11/08/1971	1	Houses of the Holy	03/28/1973	1	Physical Graffiti	02/24/1975	1
-----------------	------------	---	--------------------	------------	---	-------------------	------------	---

Column-Oriented data on disk

Led Zeppelin IV	Houses of the Holy	Physical Graffiti	11/08/1971	03/28/1973	02/24/1975	1	1	1
-----------------	--------------------	-------------------	------------	------------	------------	---	---	---

Apache Parquet

	day	location	product	sale
row 1	2017-01-01	l1	p1	300
row 2	2017-01-01	l1	p2	40
row 3	2017-01-01	l2	p1	44
row 4	2017-02-01	l1	p1	200

Traditional Memory Buffer	
row 1	2017-01-01
	l1
	p1
	300
row 2	2017-01-01
	l1
	p2
	40
row 3	2017-01-01
	l2
	p1
	44

Columnar Storage	
day	2017-01-01
day	2017-01-01
day	2017-01-01
day	2017-01-02
location	l1
location	l1
location	l2
location	l1
product	p1
product	p2
product	p1
product	p1



Apache Parquet

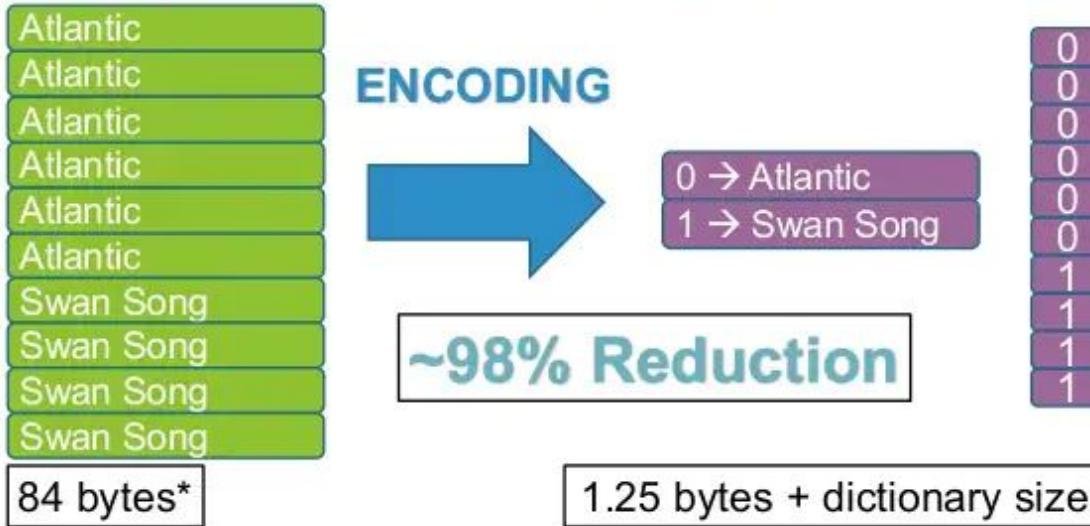
Encoding: Incremental Encoding



*not counting delimiters

Apache Parquet

Encoding: Dictionary Encoding



Apache Parquet

Partitioning

```
dataFrame  
  .write  
  .partitionBy("Whatever", "Columns", "You", "Want")  
  .parquet(outputFile)  
  
// For a common example  
dataFrame  
  .write  
  .partitionBy("Year", "Month", "Day", "Hour")  
  .parquet(outputFile)
```



Apache Parquet



VS



- Não existe estatísticas dos dados no arquivo.
- É preciso ler todo o arquivo para definição dos data types.
- Não permite compressão por coluna.
- Ocupa mais espaço e mais processamento.

Apache Parquet

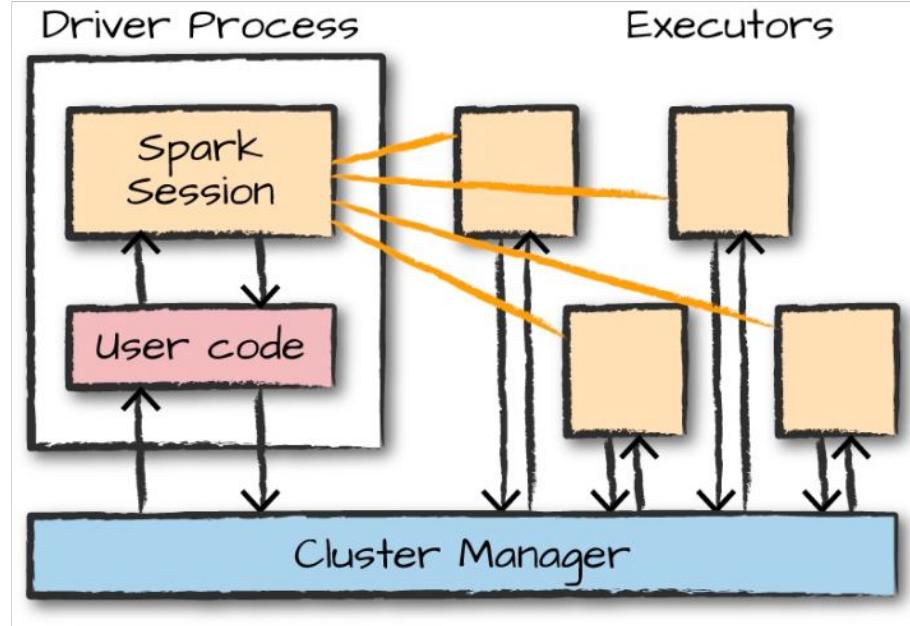
The following table compares the savings as well as the speedup obtained by converting data into Parquet from CSV.

Dataset	Size on Amazon S3	Query Run Time	Data Scanned	Cost
Data stored as CSV files	1TB	236 seconds	1.15 TB	\$5.75
Data stored in Apache Parquet Format	130 GB	6.78 seconds	2.51 GB	\$0.01
Savings	87% less when using Parquet	34x faster	99% less data scanned	99.7% savings

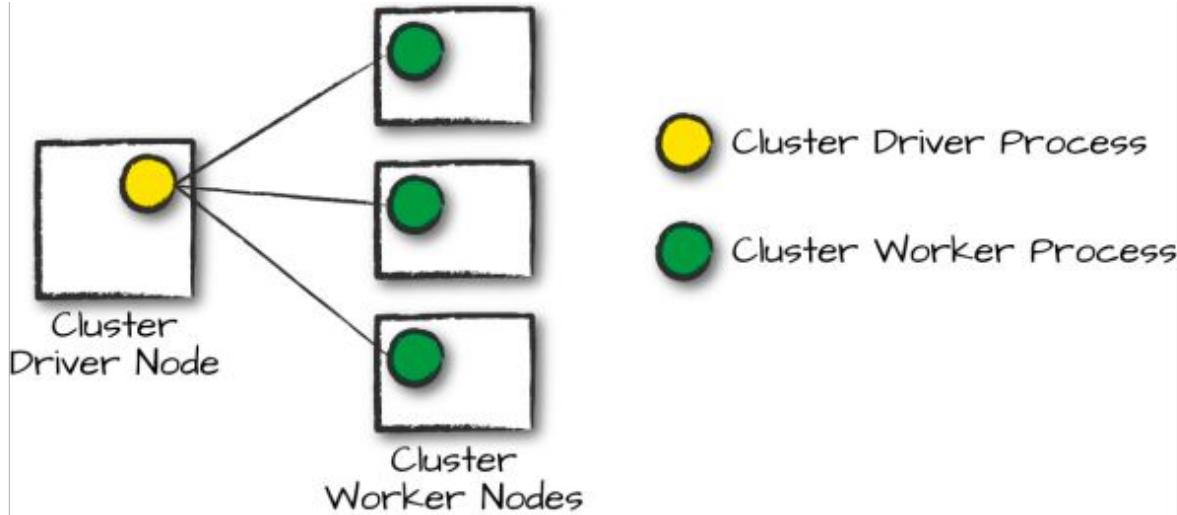
Hands on!

Arquitetura Spark

- Master: Sessões, Planos de Execução
- Slaves: Executores, máquinas escravas para processamento de jobs.



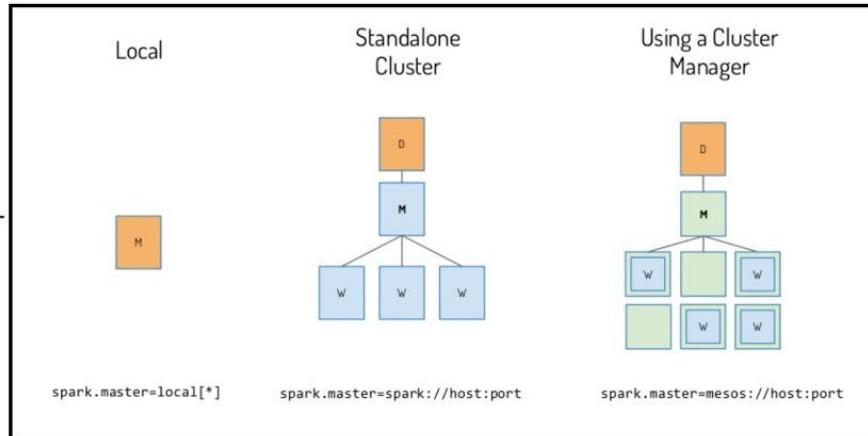
Cluster Communication



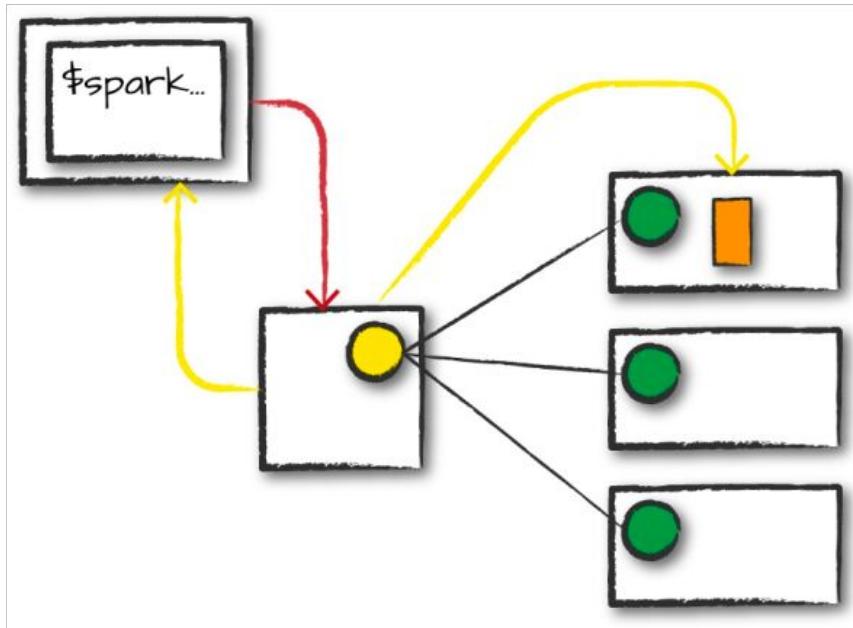
Deployment modes



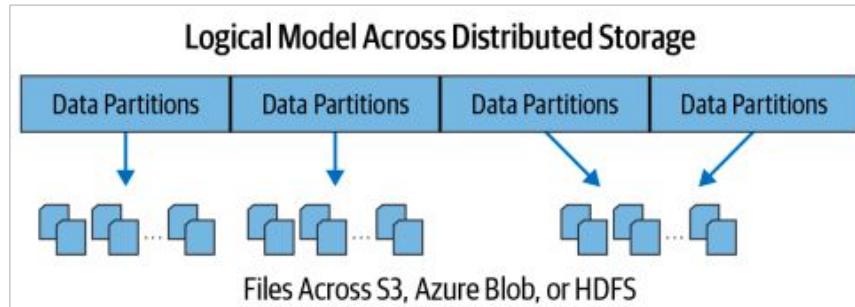
Deployment Options



Executando uma aplicação



Leitura de dados



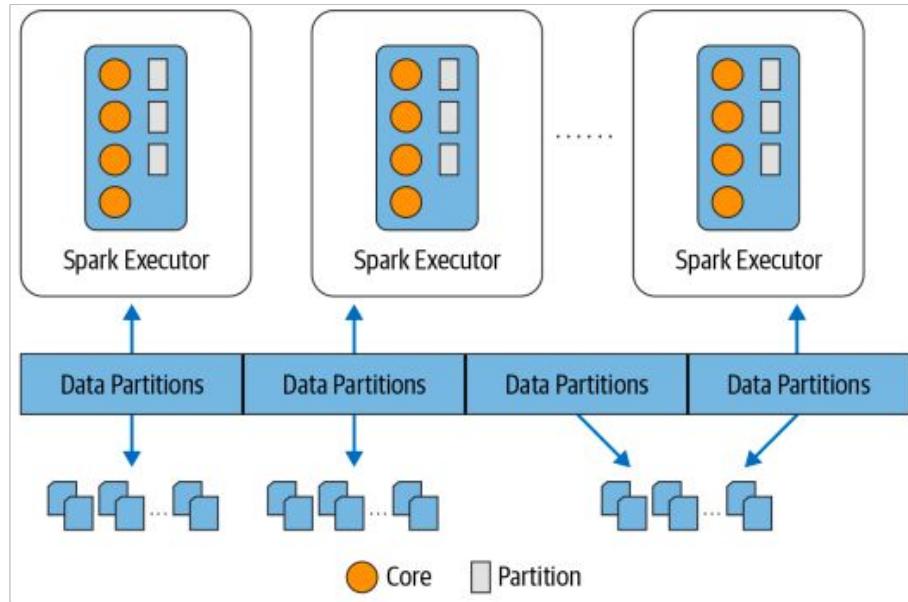
Storage



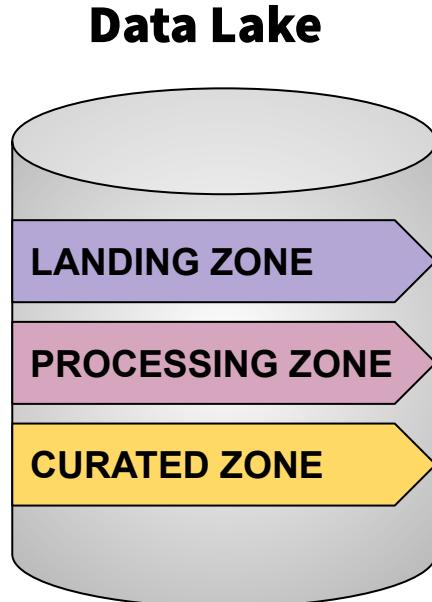
Database (jdbc)



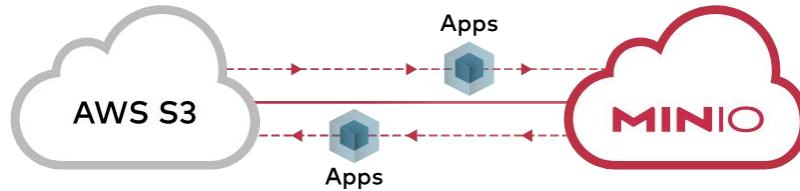
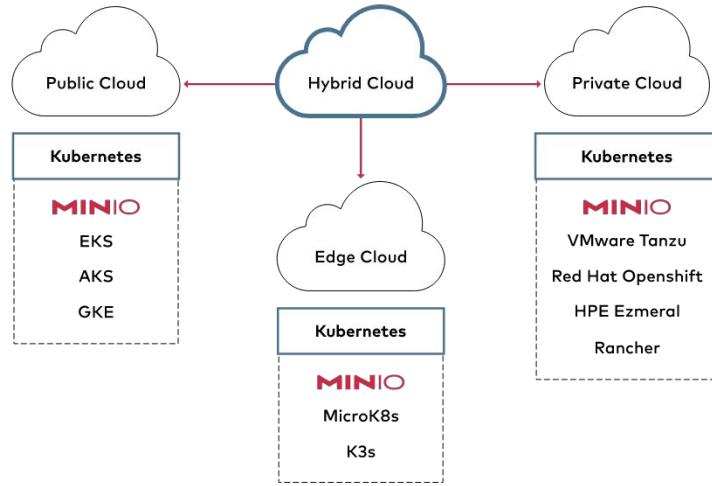
Leitura e processamento de dados



Criando um ambiente local

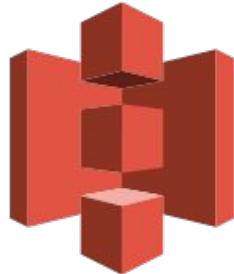


Criando um ambiente local



Hands on!

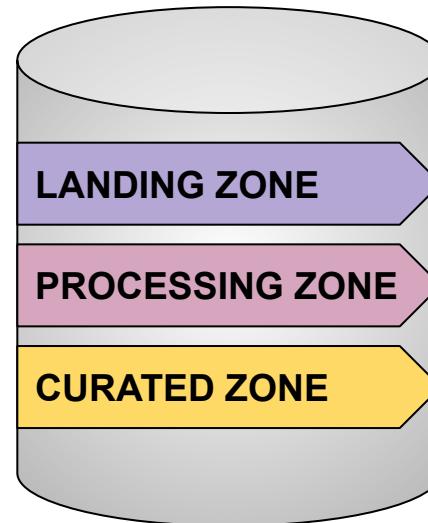
Criando um ambiente em Cloud



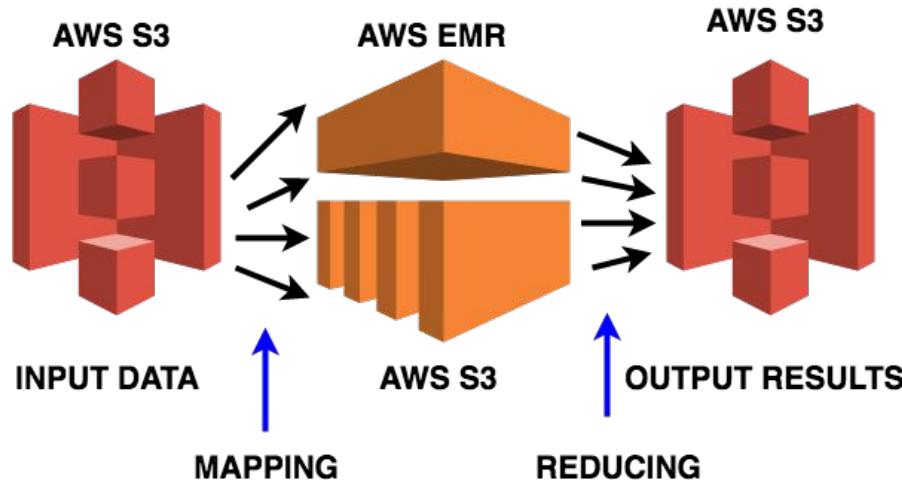
amazon
S3



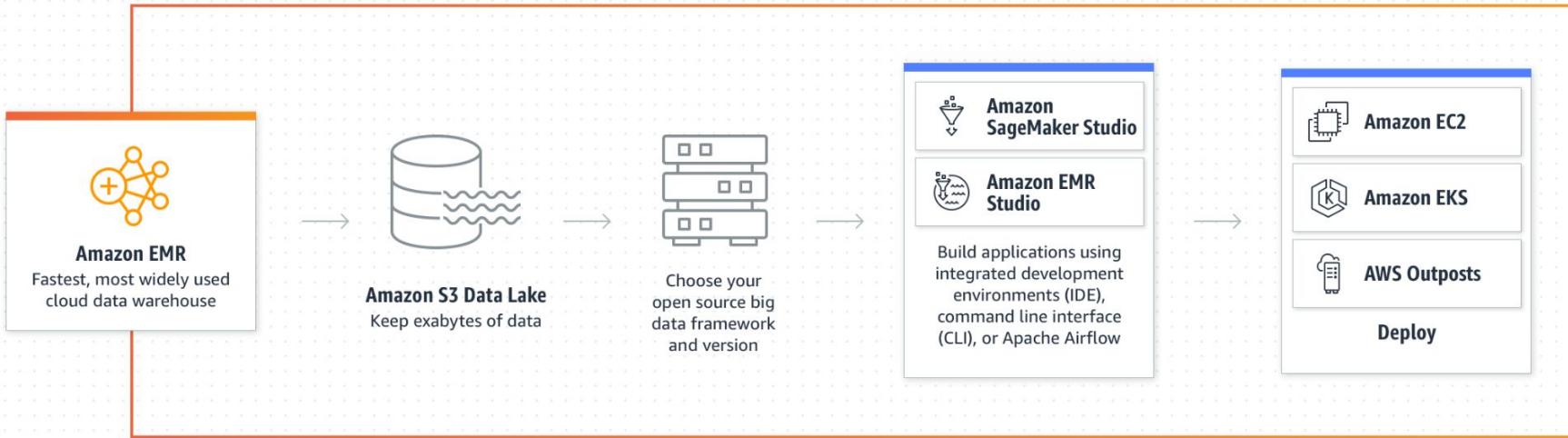
Data Lake



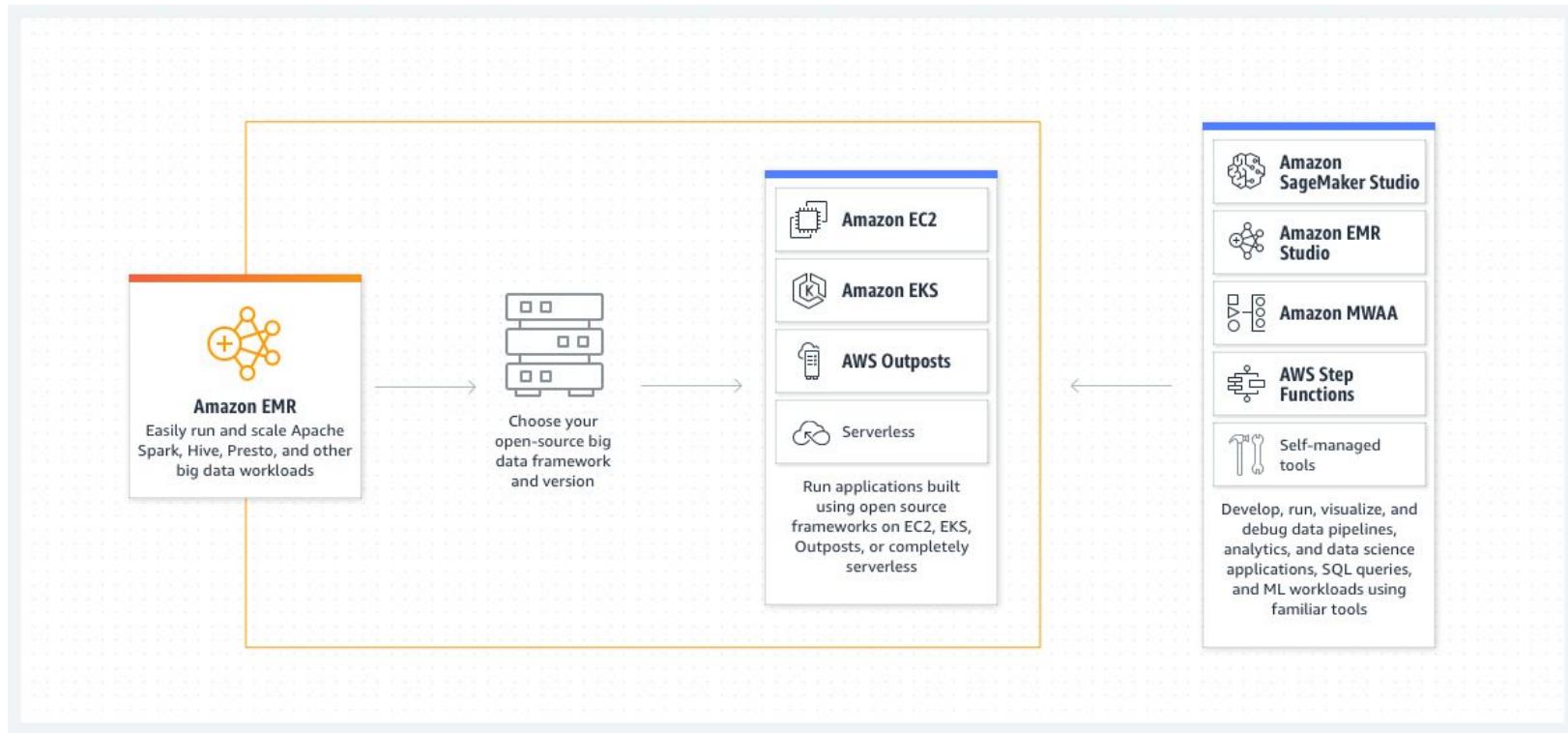
Criando um ambiente em Cloud



Amazon EMR

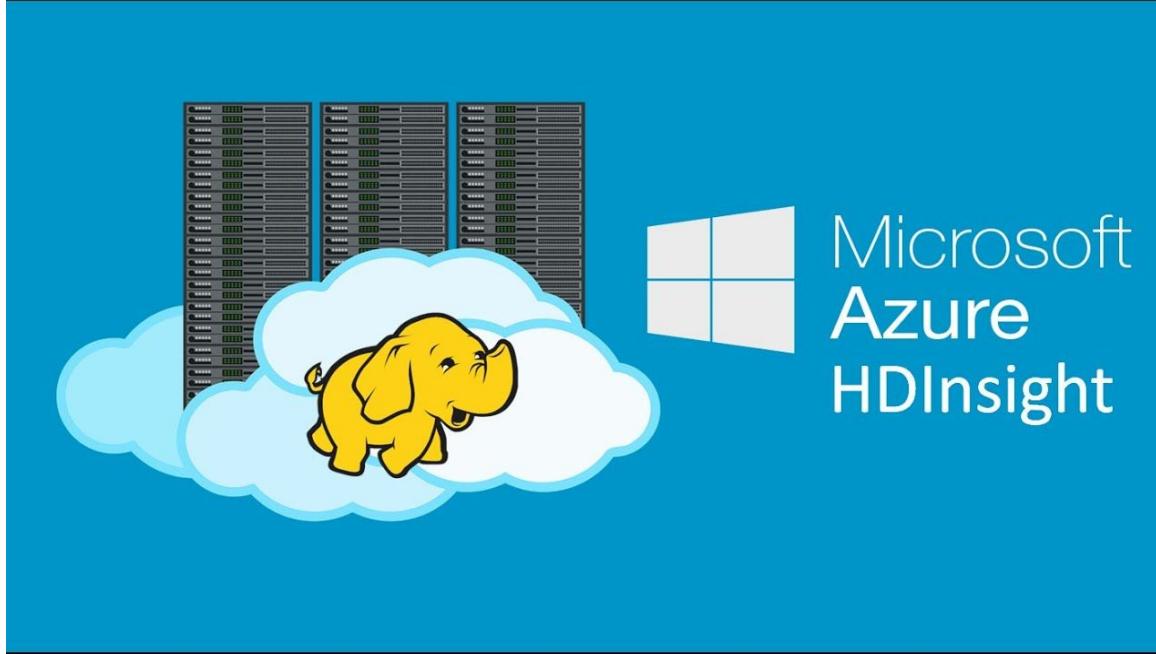


Amazon EMR

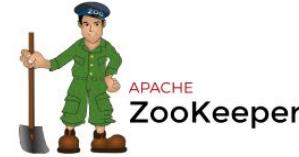


Hands on!

HDinsight



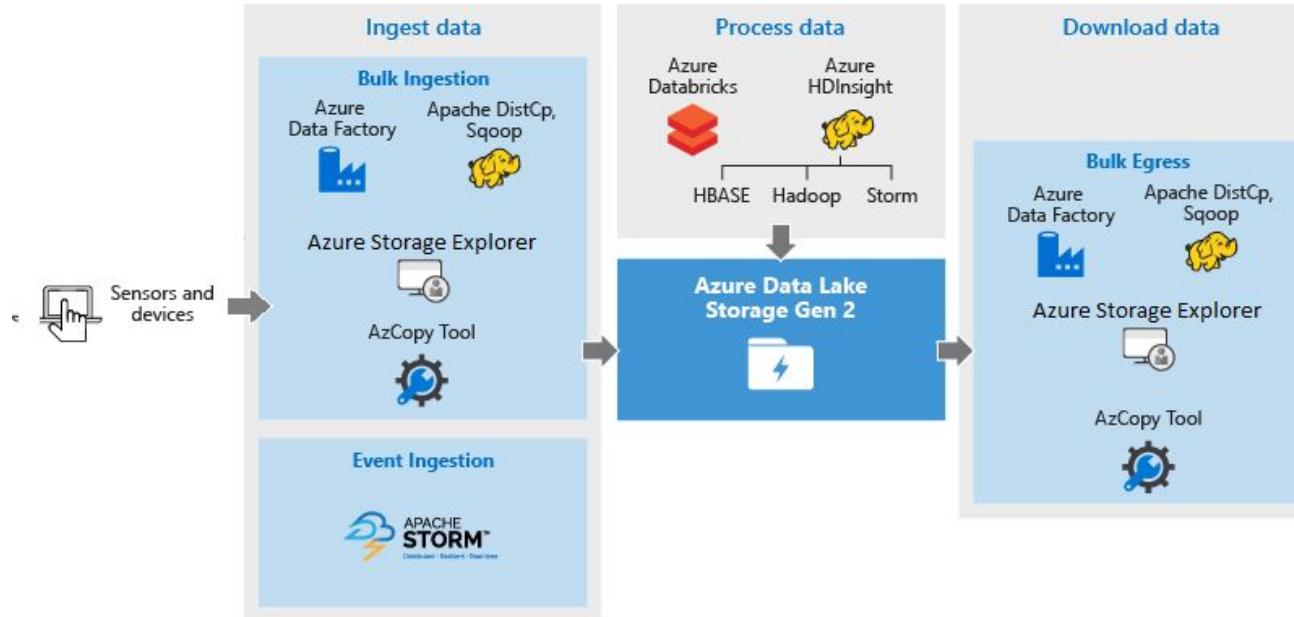
HDinsight



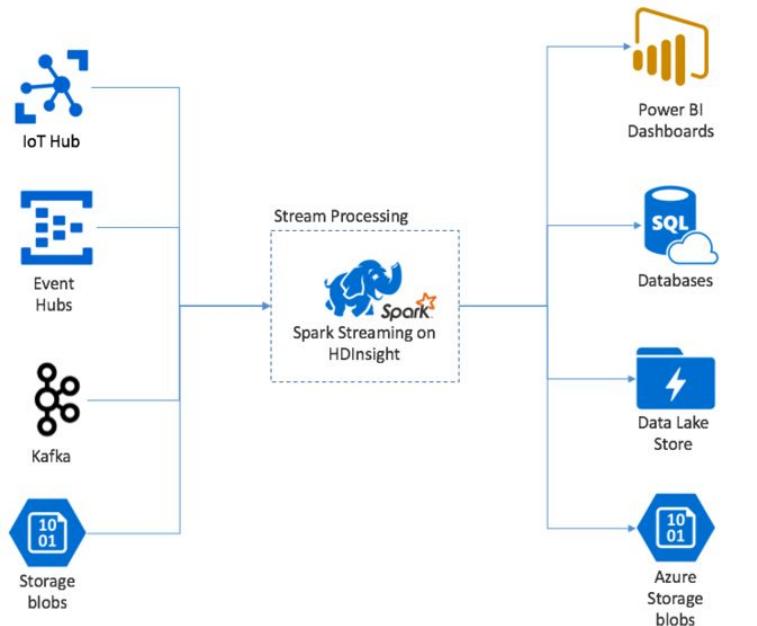
Apache Ambari

Stack

Data Lake



HDinsight



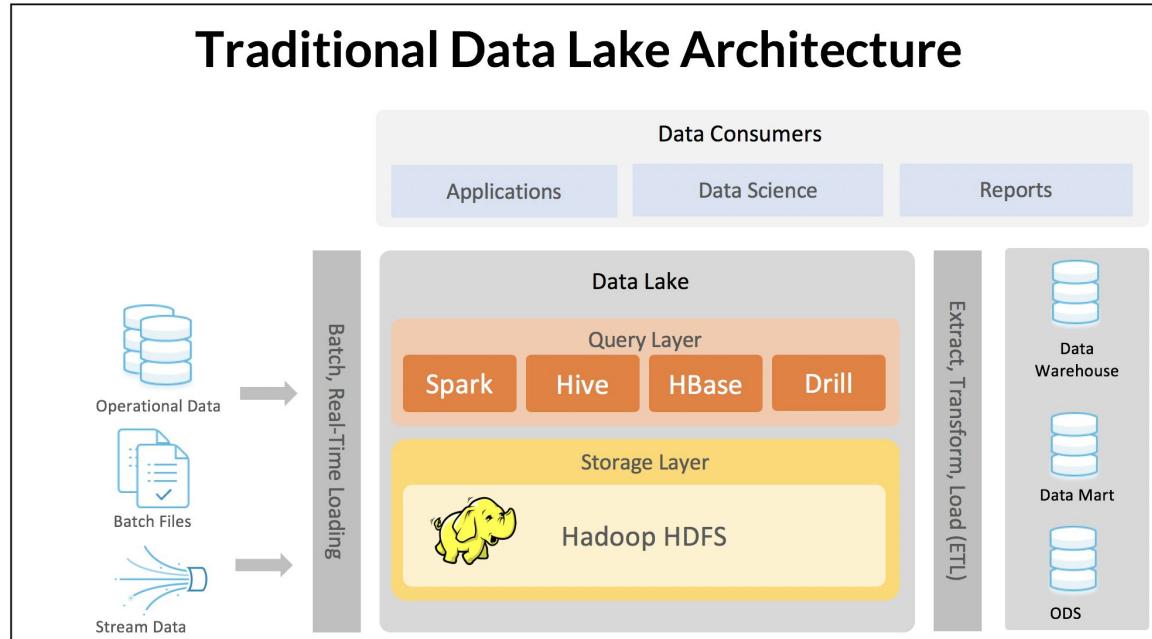
Hands on!

Hadoop

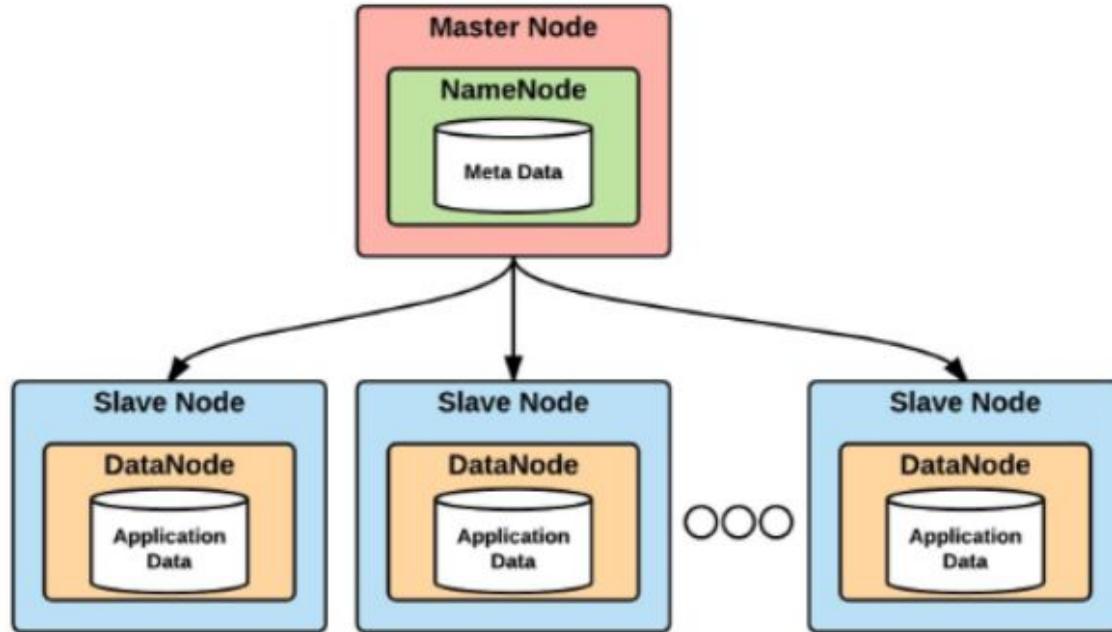


Hadoop

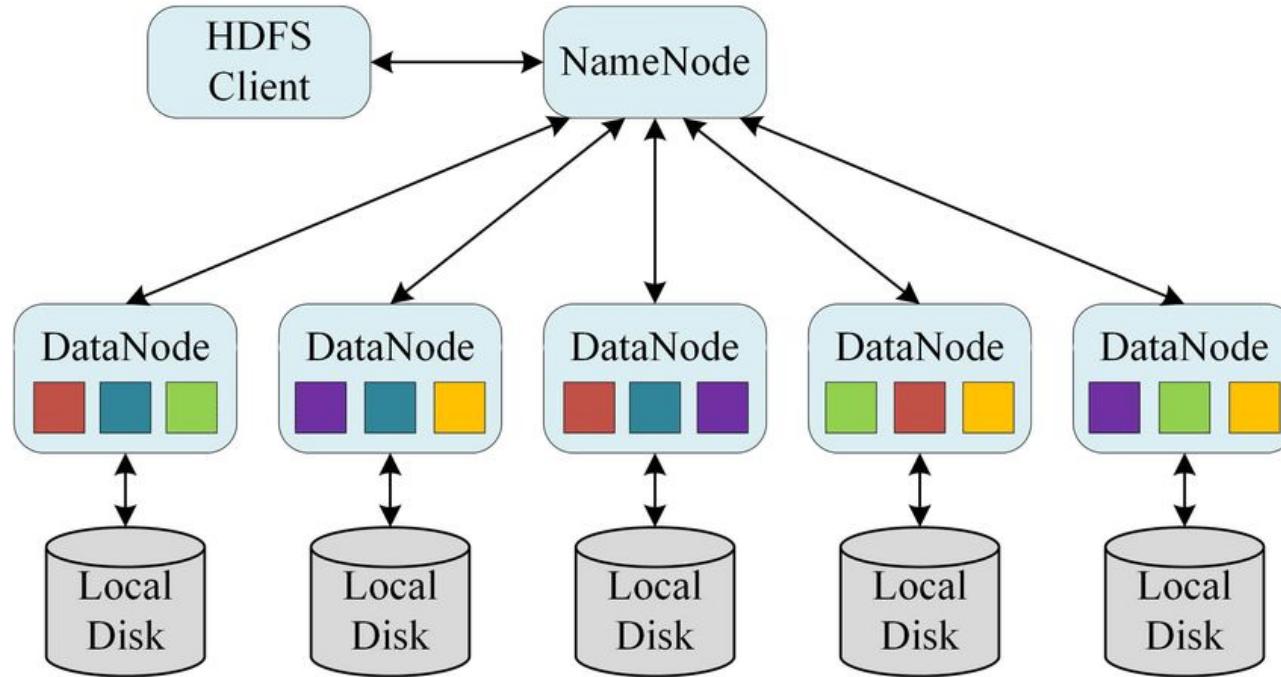
Traditional Data Lake Architecture



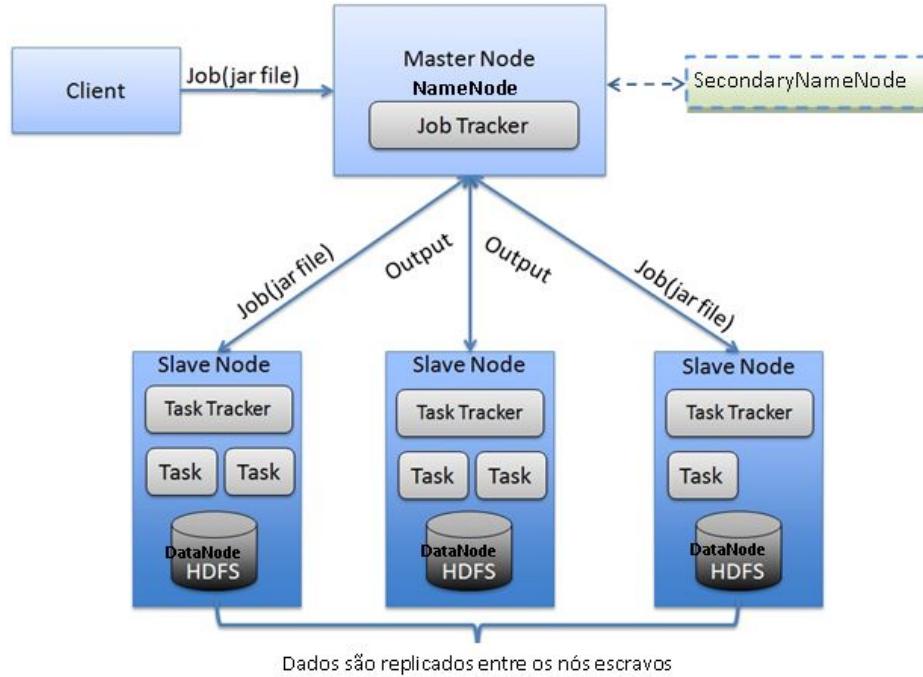
Hadoop



Hadoop



Hadoop



Hadoop

Factors	Spark	Hadoop MapReduce
Speed	100x times than MapReduce	Faster than traditional system
Written In	Scala	Java
Data Processing	Batch / real-time / iterative / interactive / graph	Batch processing
Ease of Use	Compact & easier than Hadoop	Complex & lengthy
Caching	Caches the data in-memory & enhances the system performance	Doesn't support caching of data

Hands on!

Hive



Hive

Challenge...

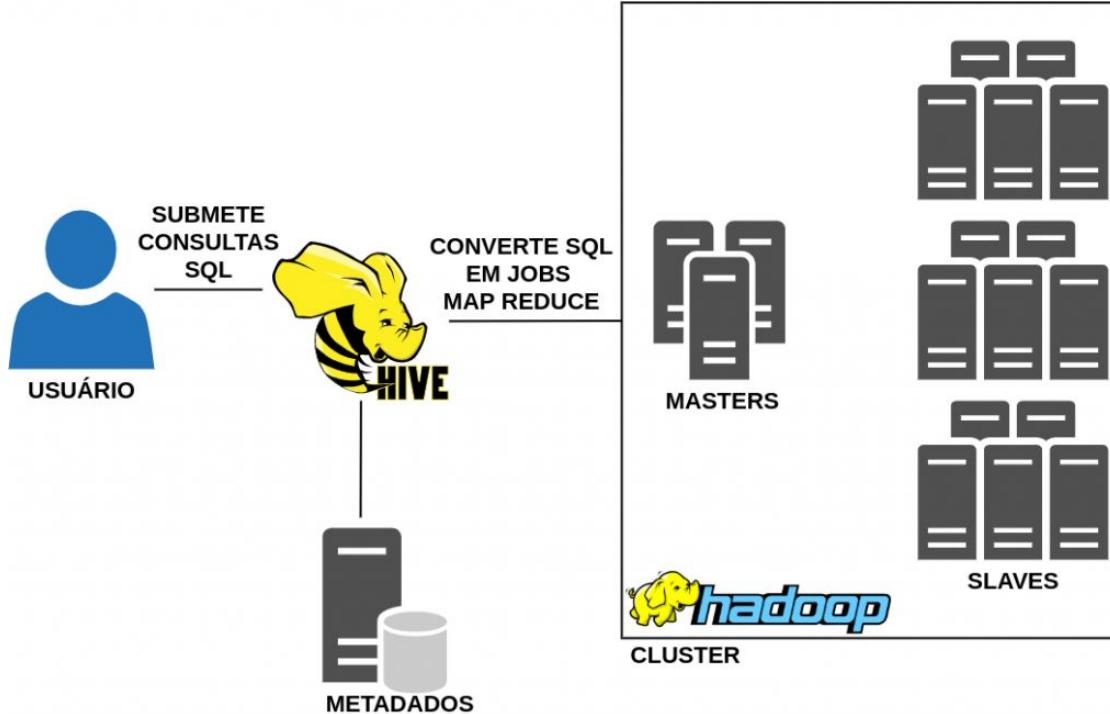


Traditional RDBMS... X

Solution...



Hive



AWS Athena

Overview of AWS Athena



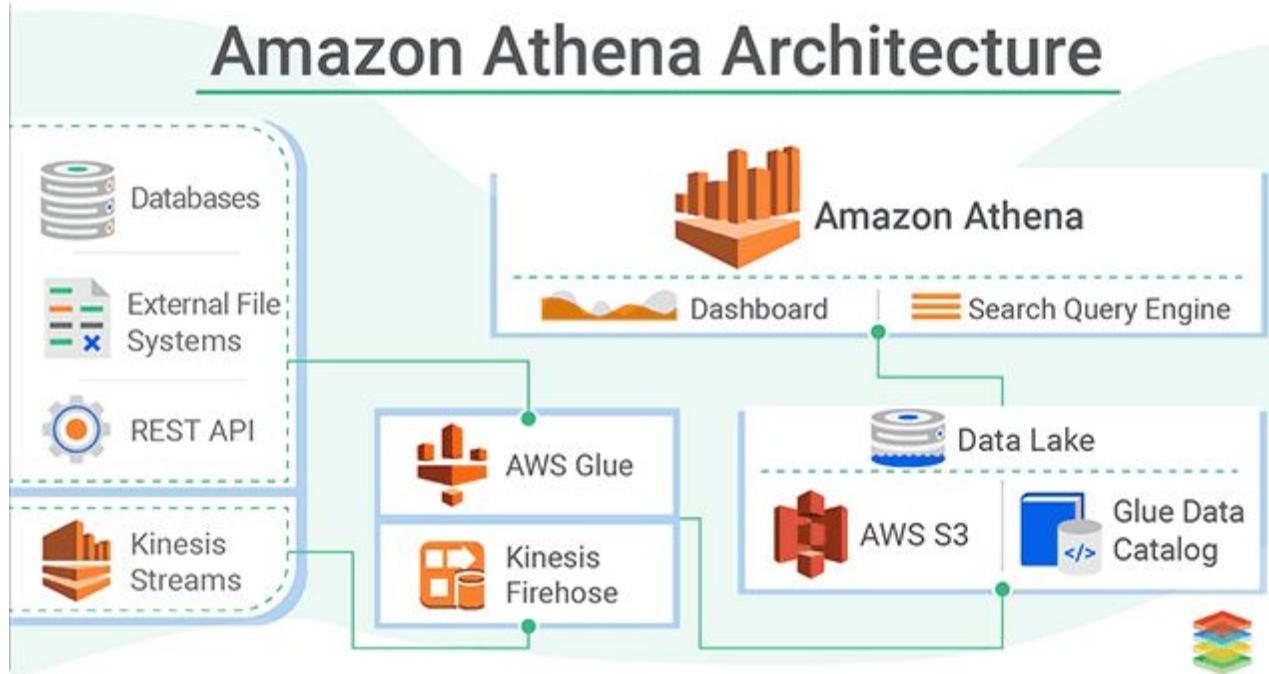
Serverless
Interactive
Query

Amazon Athena

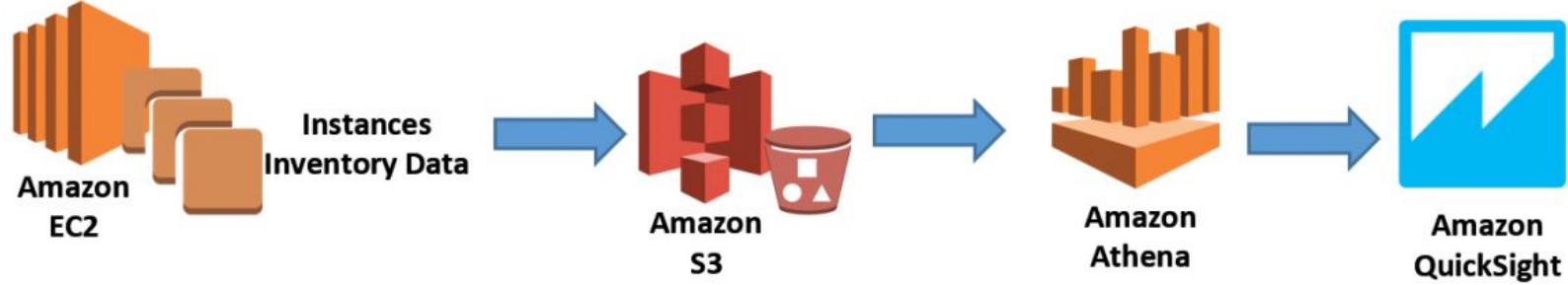
AWS Athena



Arquitetura



Arquitetura



AWS Athena



Athena



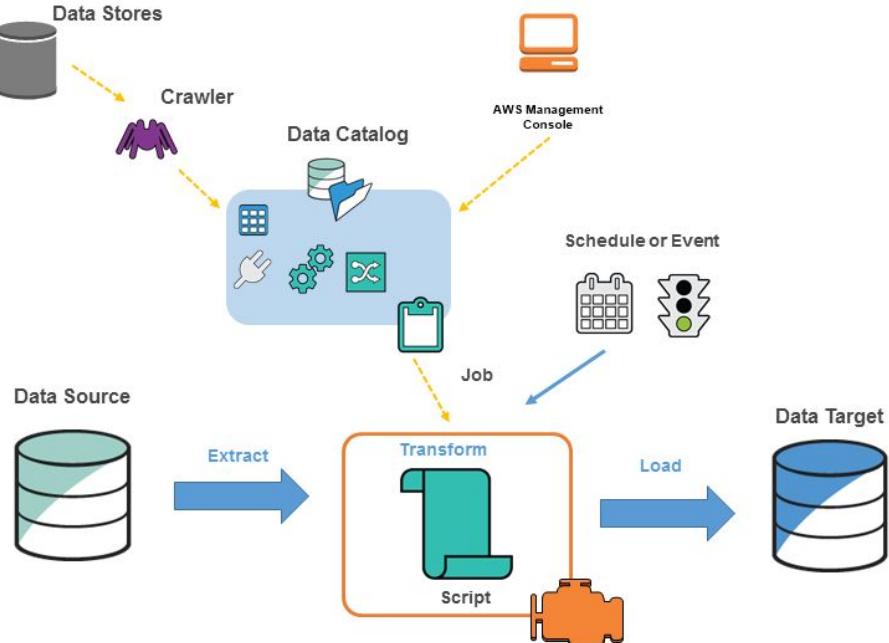
EMR

- Menos utilizado para processamento de dados em larga escala.
- Mais utilizado como uma ferramenta para executar queries ad-hoc.
- Serverless - não precisa gerenciar o cluster.

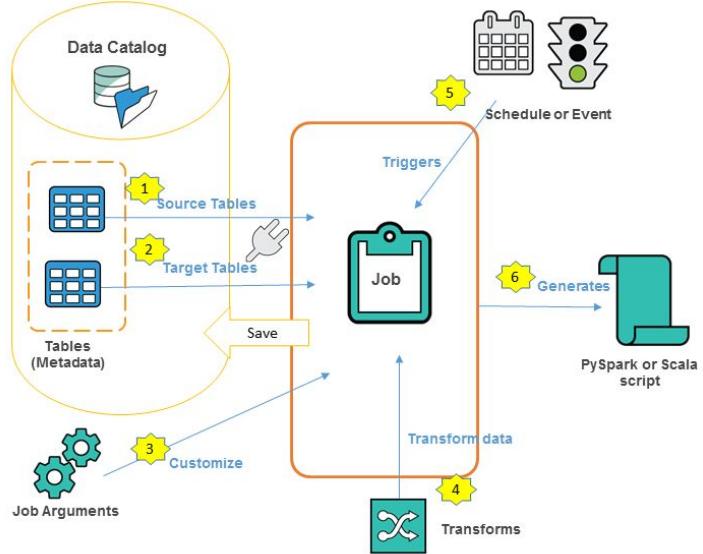
- É muito mais do que execuções de queries.
- Utiliza scripts/códigos customizados para processar e analisar grande volume de dados.
- Necessário gerenciar o cluster do EMR.

Hands on!

AWS Glue + Athena



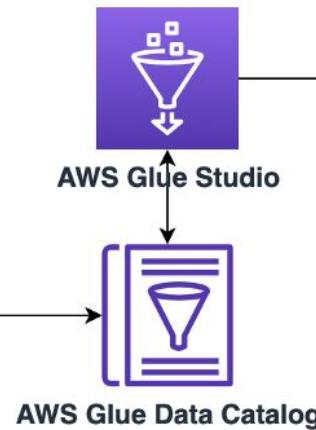
AWS Glue Studio



Data Sources



Data Pipeline



Data Interactive Query



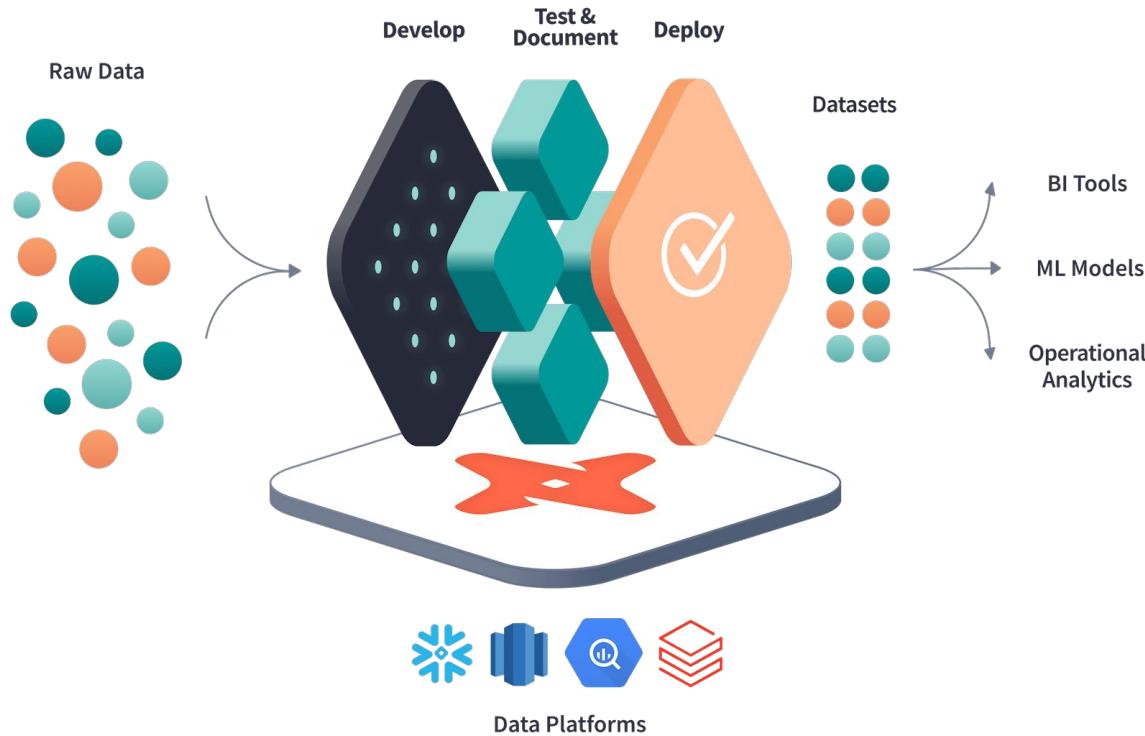
Data Visualization



DBT - Data Build Tool



DBT - Data Build Tool



Data build tool

Version Control and CI/CD

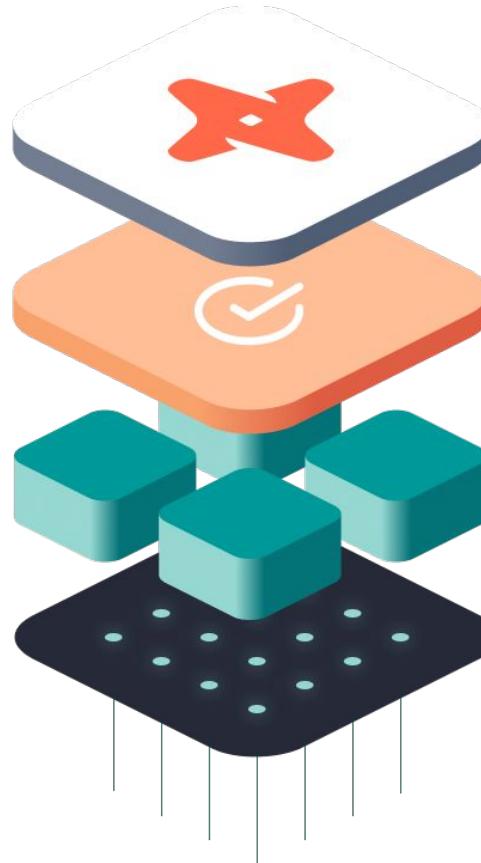
Deploy safely using dev environments.
Git-enabled version control enables
collaboration and a return to previous states.

Test and Document

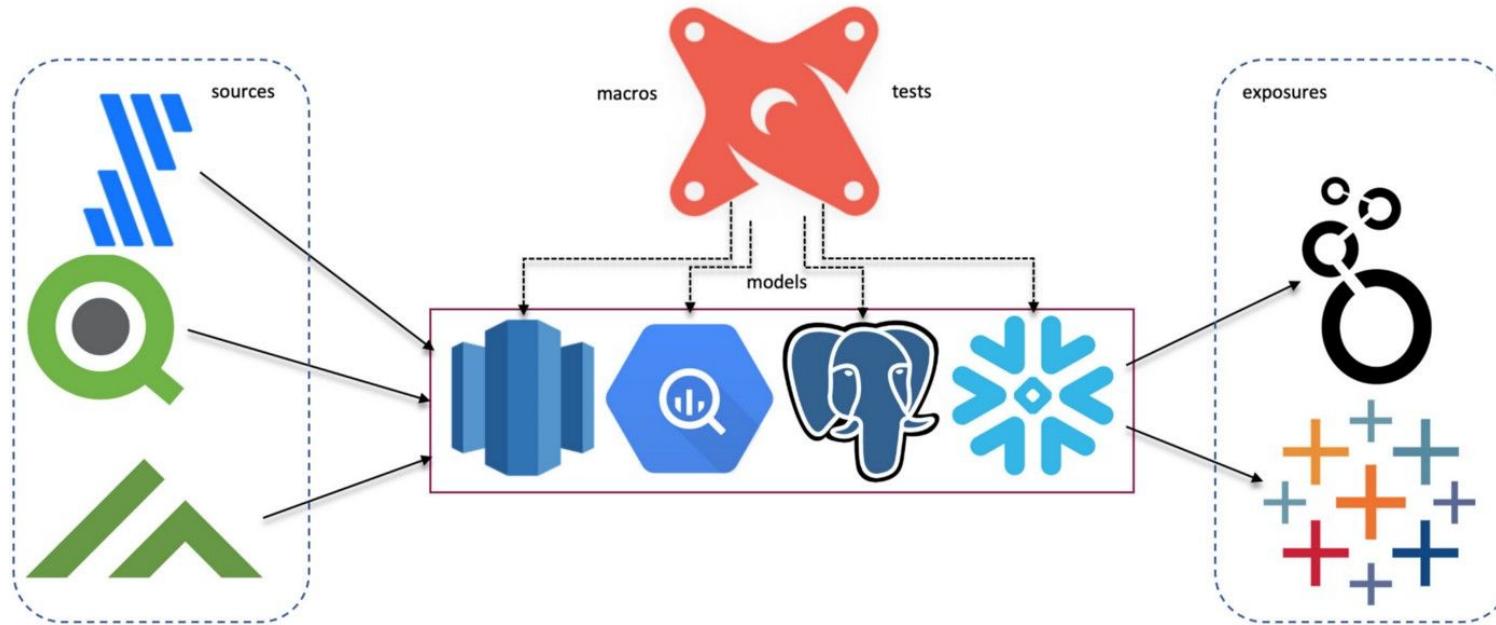
Test every model prior to production, and share
dynamically generated documentation with all
data stakeholders.

Develop

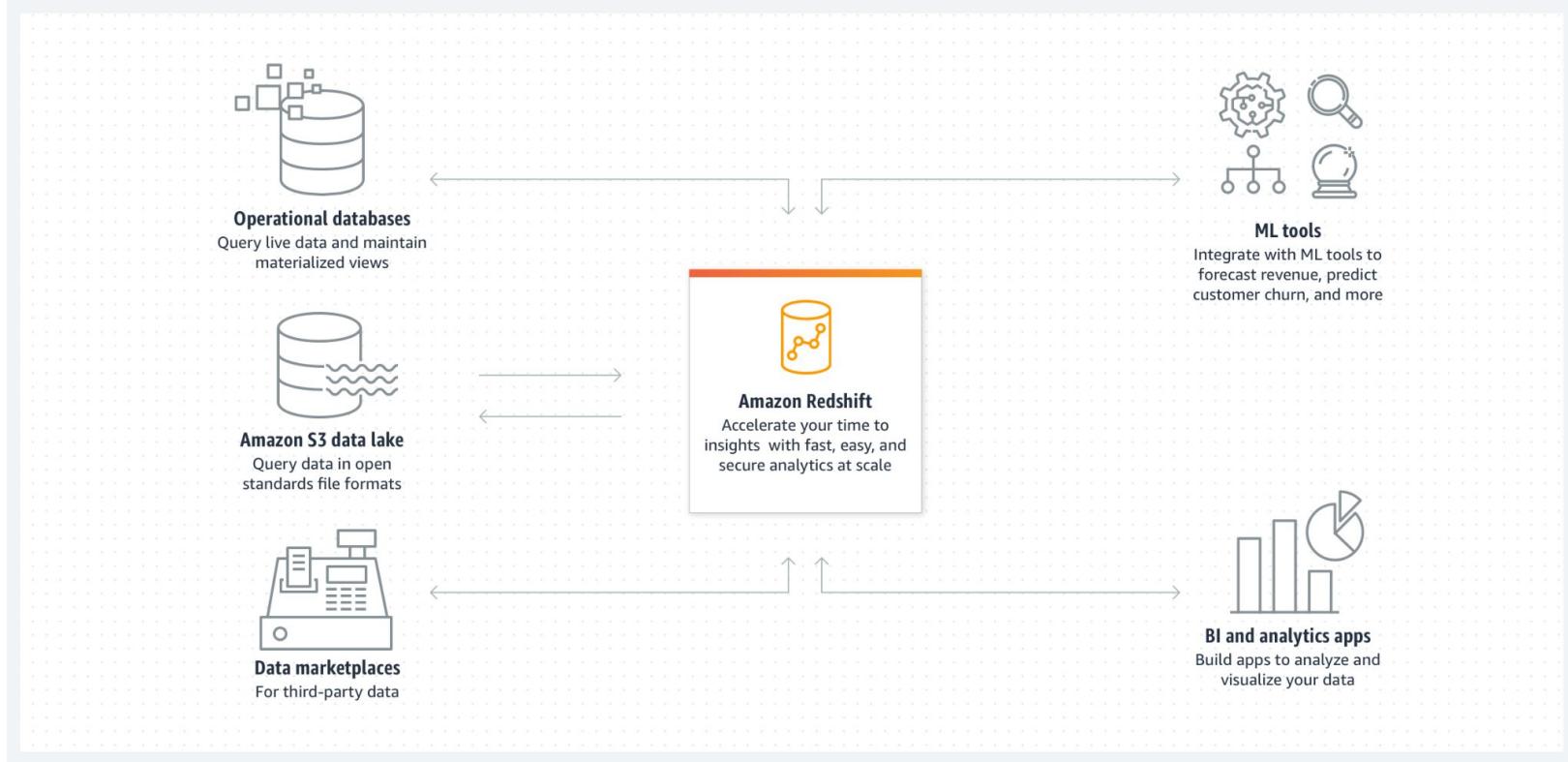
Write modular SQL models with SELECT
statements and the ref() function—dbt handles
the chore of dependency management.



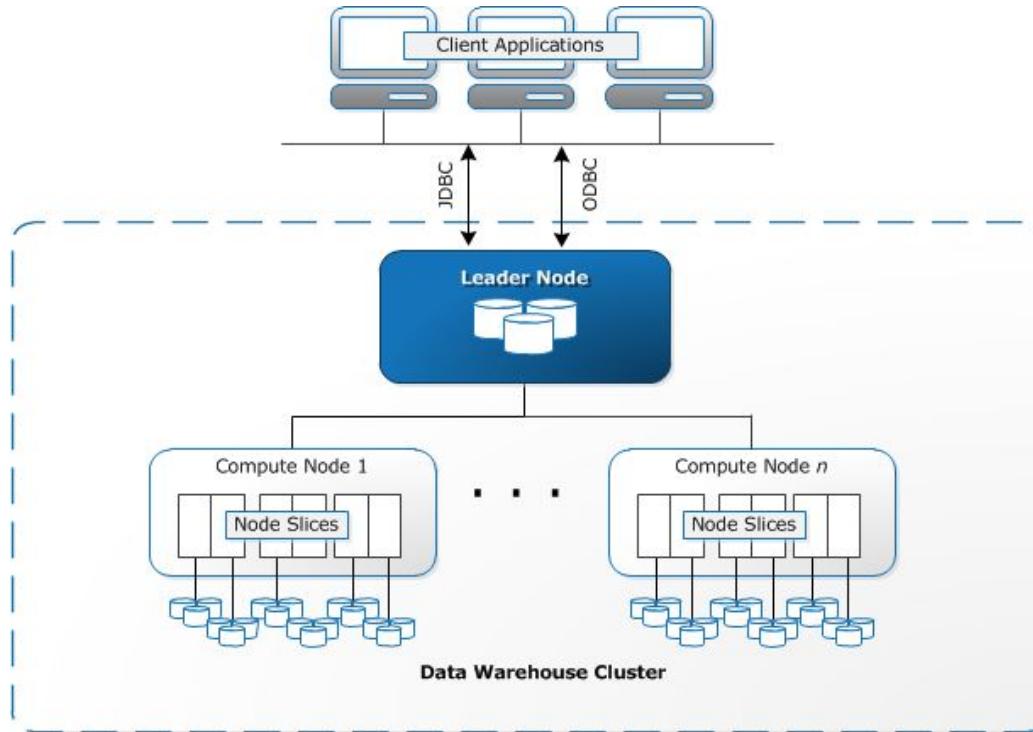
Data build tool



Redshift



Redshift



Redshift

	client_id	client_name	client_email	client_state
+	1	ABC	a@abc.com	MH
row below		DEF	d@def.com	RJ
	3	GHI	g@ghi.com	MP

1 ABC a@abc.com MH	2 DEF d@def.com RJ	3 GHI g@ghi.com MP
--------------------	--------------------	--------------------

Block 1

Block 2

Block 3

Columnar Storage

1 2 3 4 5 6 7 8 9 10	11 12 13 14 15 16 17 18 19	20 21 22 23 24 25 26 27
----------------------	----------------------------	-------------------------

Block 1

Block 2

Block 3

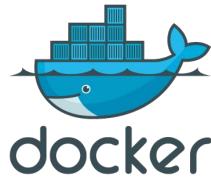
Redshift



**amazon
REDSHIFT**

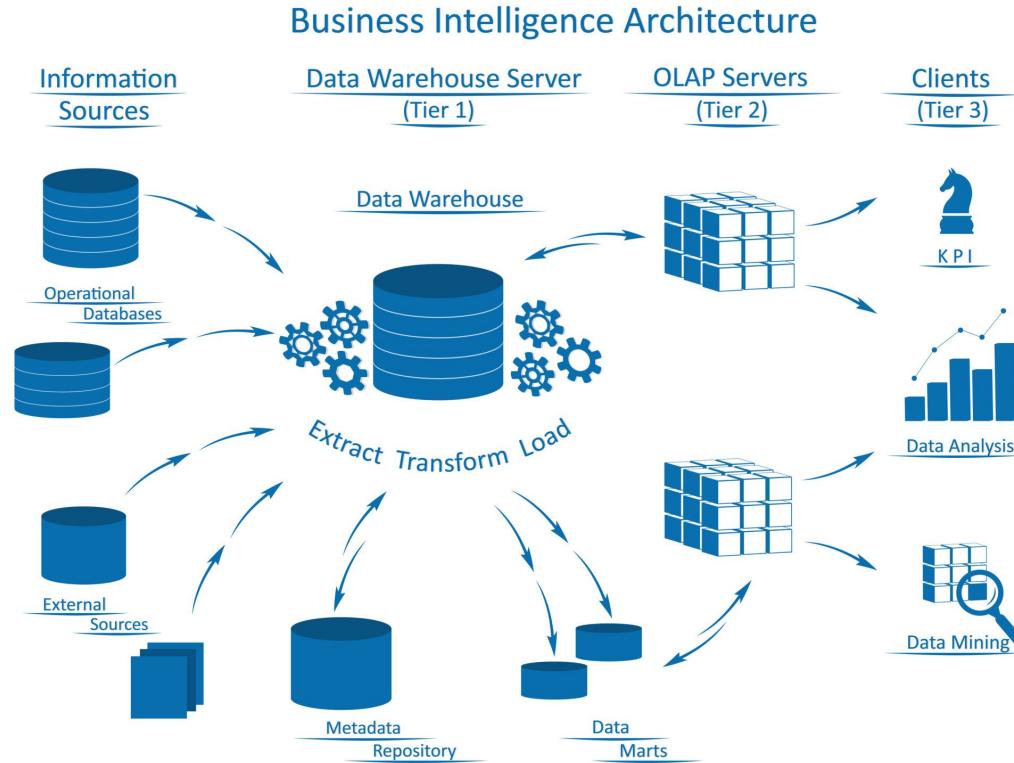
- Compression
- Massively Parallel Processing
- Machine Learning Optimization
- AWS SageMaker

Setup local



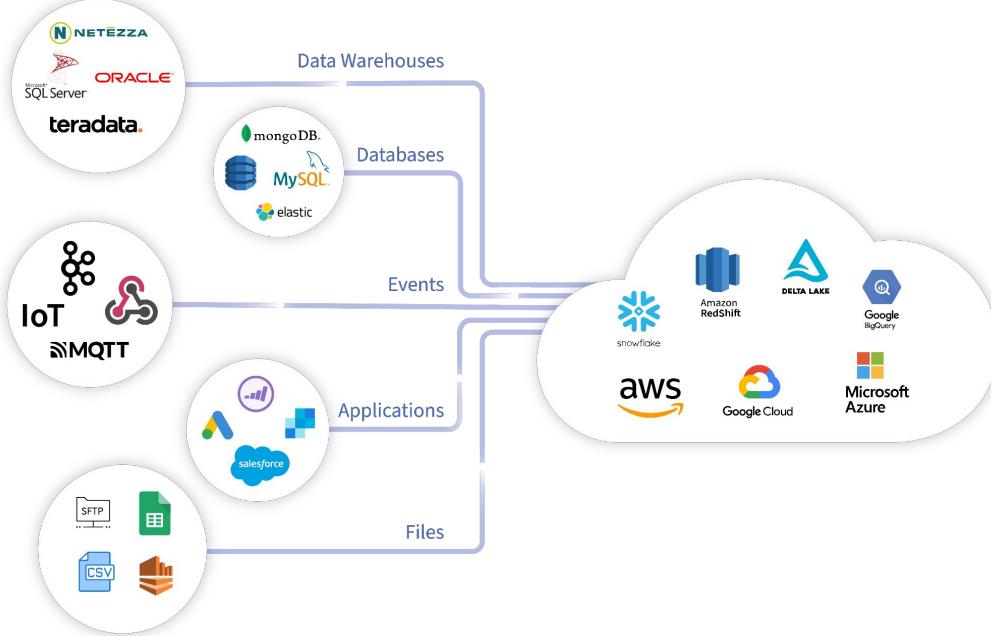
Hands on!

Ingestion process



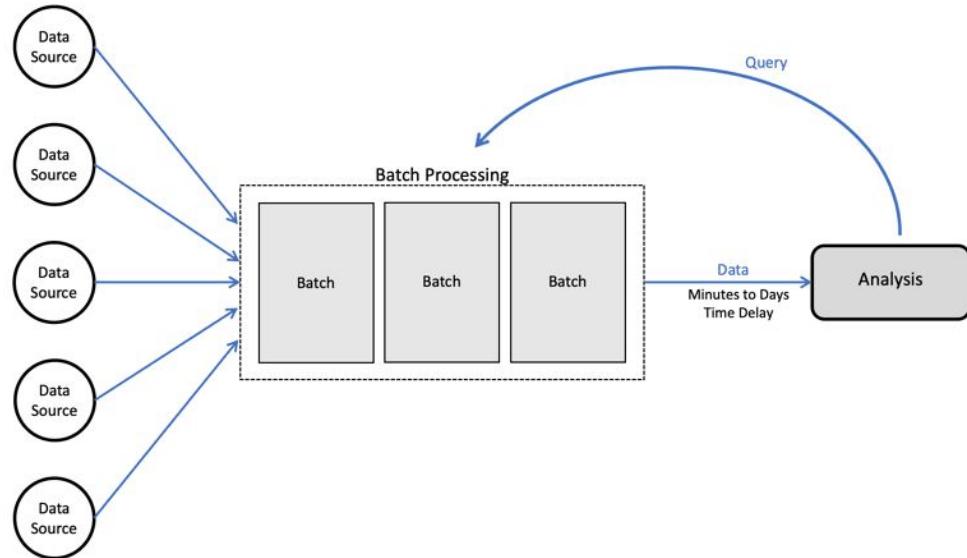
Ingestion types

- Batch
- Real-time
- Lambda architecture



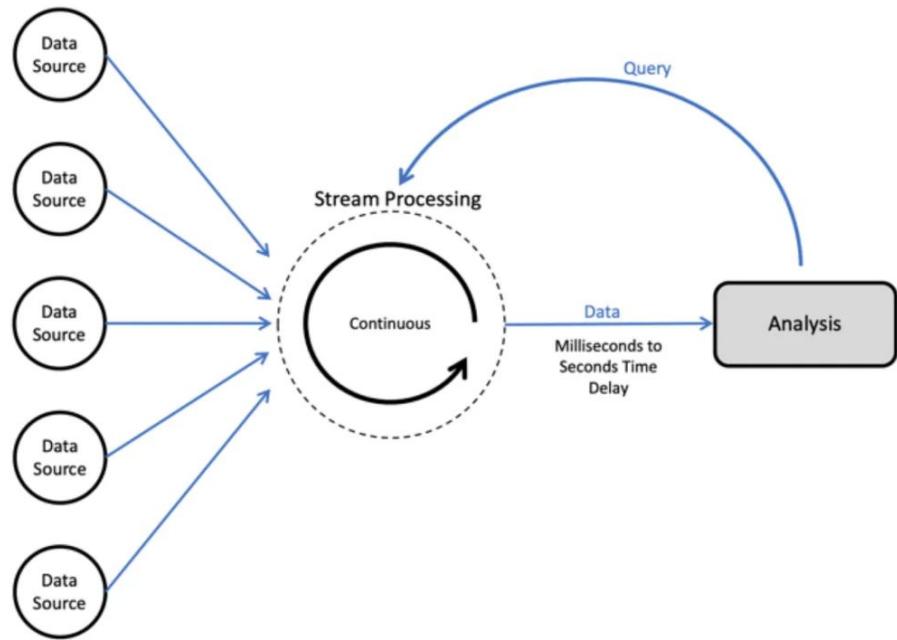
Ingestion types

- **Batch**
- Real-time
- Lambda architecture



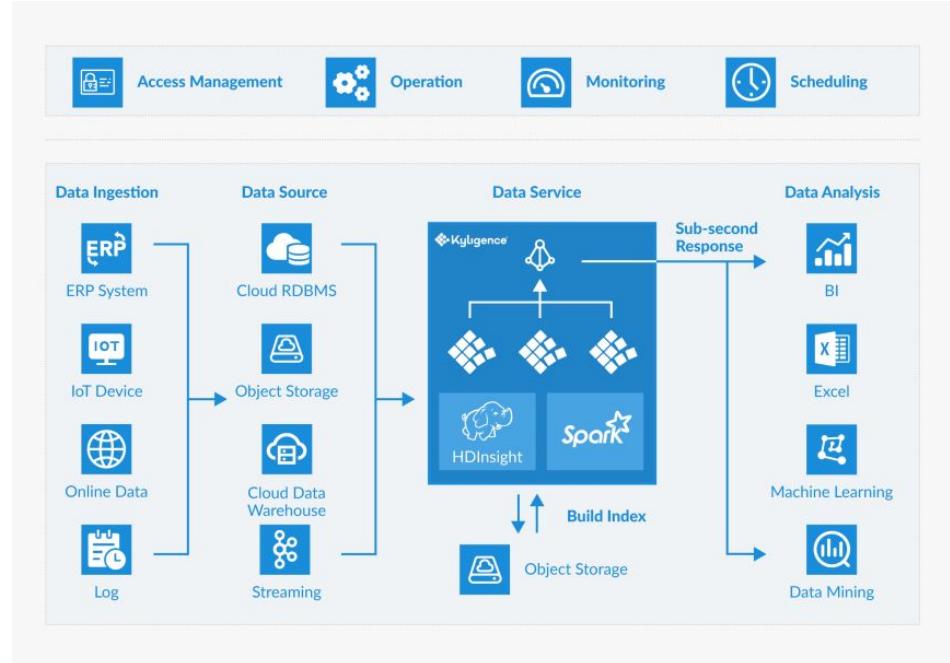
Ingestion types

- Batch
- Real-time
- Lambda architecture



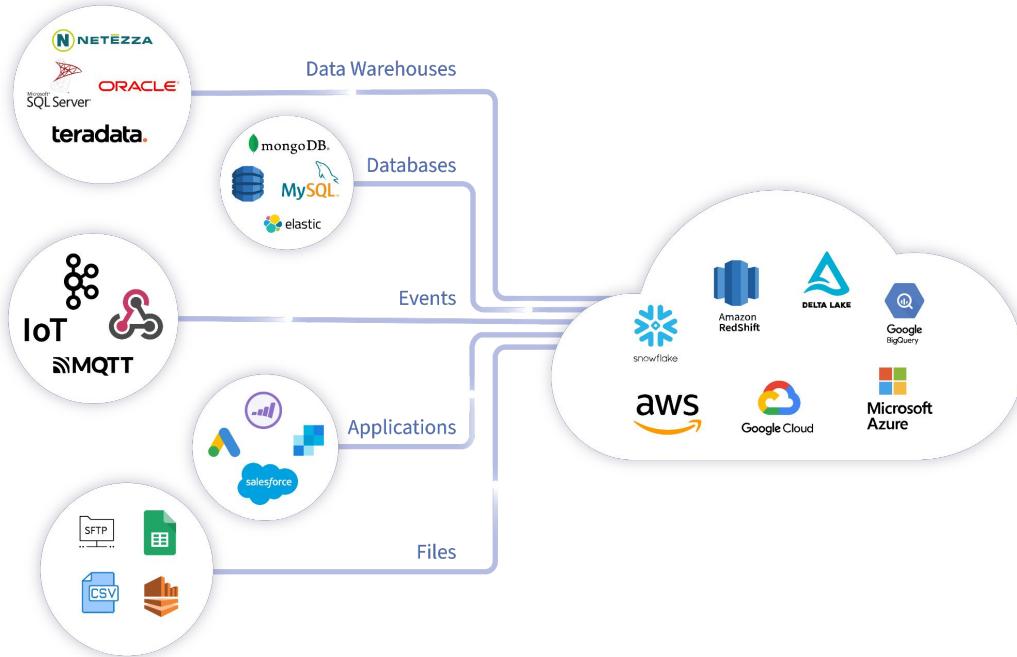
Benefits of Data Ingestion

- Data is readily available
- Data is less complex
- Teams save time and money
- Companies make better decisions
- Teams create better apps and software tools

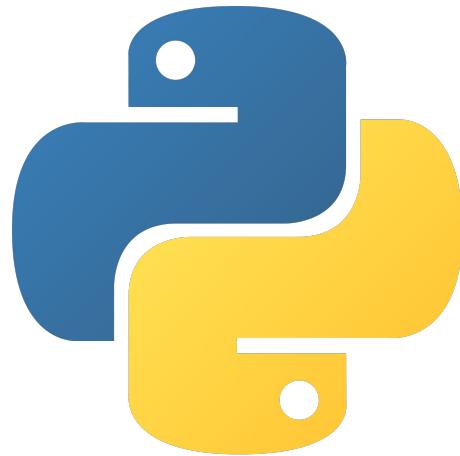
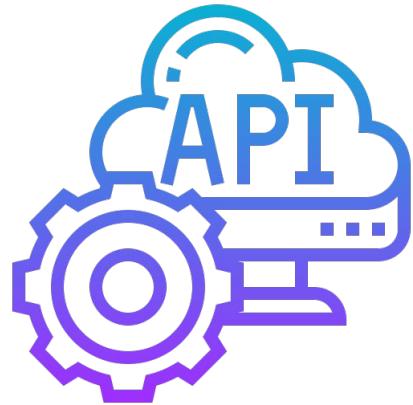


Considerations

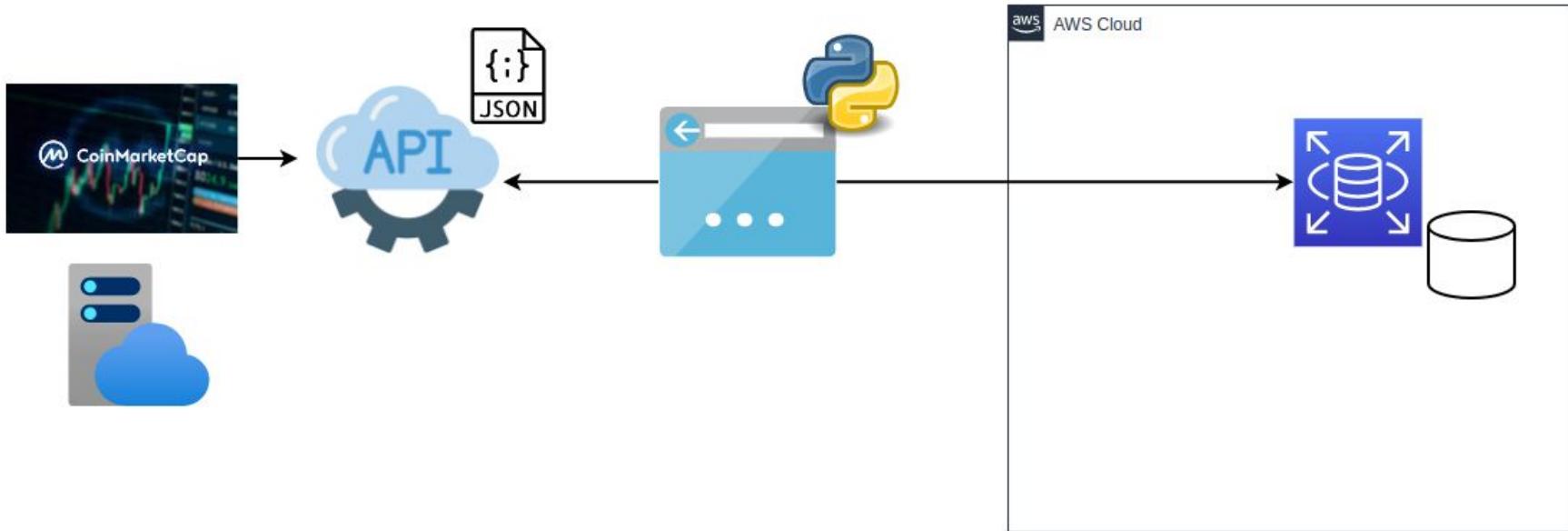
- Data Velocity
- Data Size
- Data Frequency (Change Data Capture?)
- Data Format



Data Pipeline

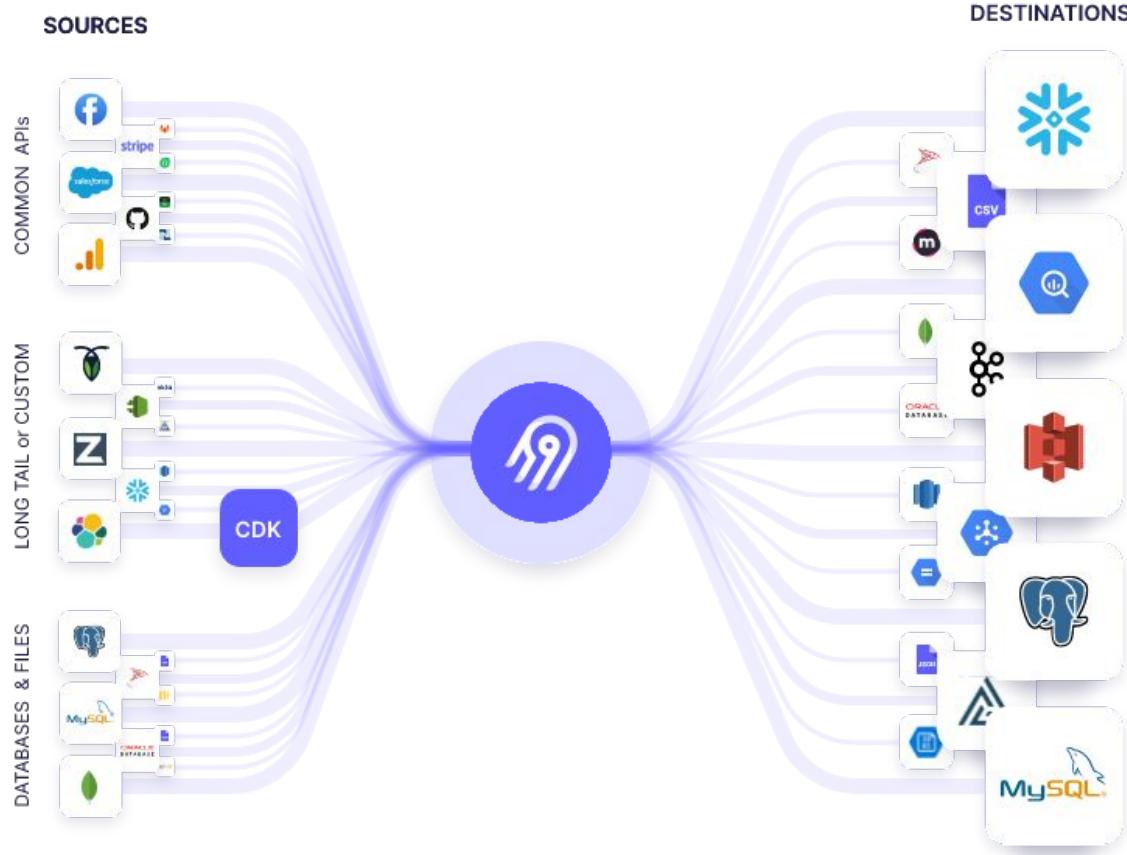


Pipeline



Hands on!

Airbyte



Hands on!