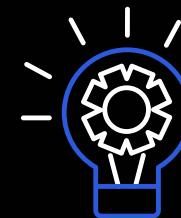




insight into value

BIG DATA ANALYTICS
BUSINESS INTELLIGENCE
INFORMATION MANAGEMENT
PERFORMANCE MANAGEMENT

Conteúdo Programático



DATA ENGINEERING

BIG DATA

CLOUD COMPUTING

APACHE HADOOP

APACHE SPARK

DATABRICKS



The Presenter

ALAN SILVA

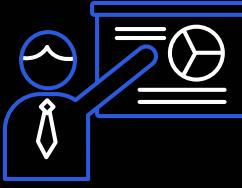
Engenheiro de Dados

Pós Graduado em Ciência de Dados e Big Data

Engenheiro da Computação

Linkedin: www.linkedin.com/in/alancarlossilva



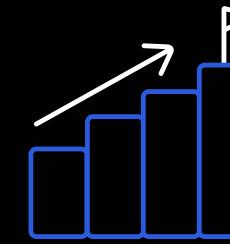


Vagas Engineers vs. Scientists

102.452



Engineers Scientists



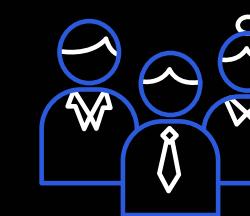
Salaries

- Data Engineer = \$110,000 - \$155,000
- Data Scientist = \$130,000 - \$185,000

Companies



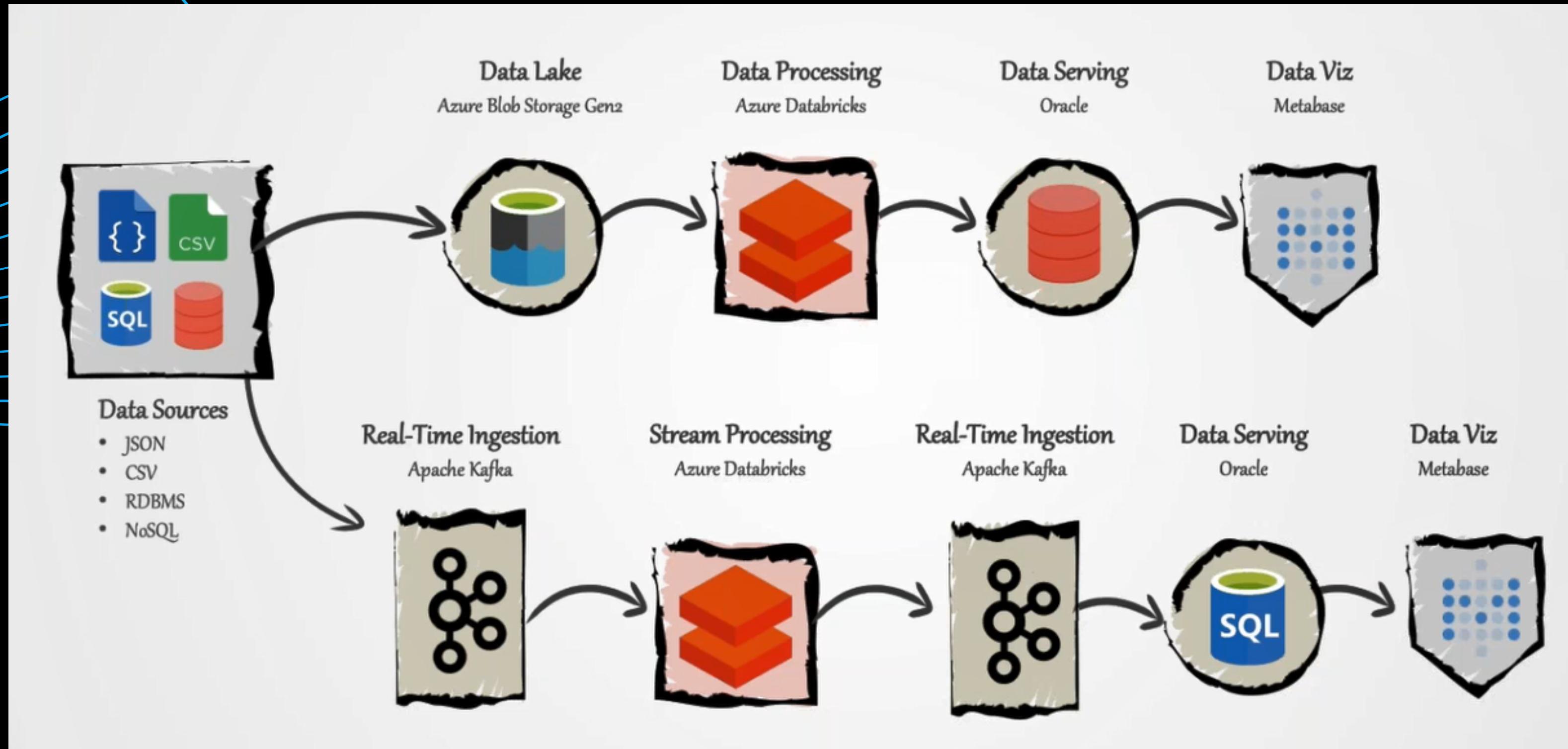
Data Engineer Resume



- Knowledge = Data Warehouse, ETL/ELT, Data Modeling
- Cloud Computing = AWS, GCP, Azure
- Programming Language = Python, Scala, BASH
- Big Data = Apache Hive, Apache Spark, Apache Kafka, Apache Airflow, Apache Druid
- RDBMS = MySQL, PostgreSQL
- NoSQL = Cassandra, MongoDB, ElasticSearch
- OS = Linux

Lambda Architecture

Batch-Layer



Fast-Layer



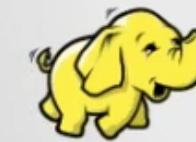
Big Data



Apache Kafka [2014]

speed of data generation & processing large datasets, ability to ingest data as fast as possible

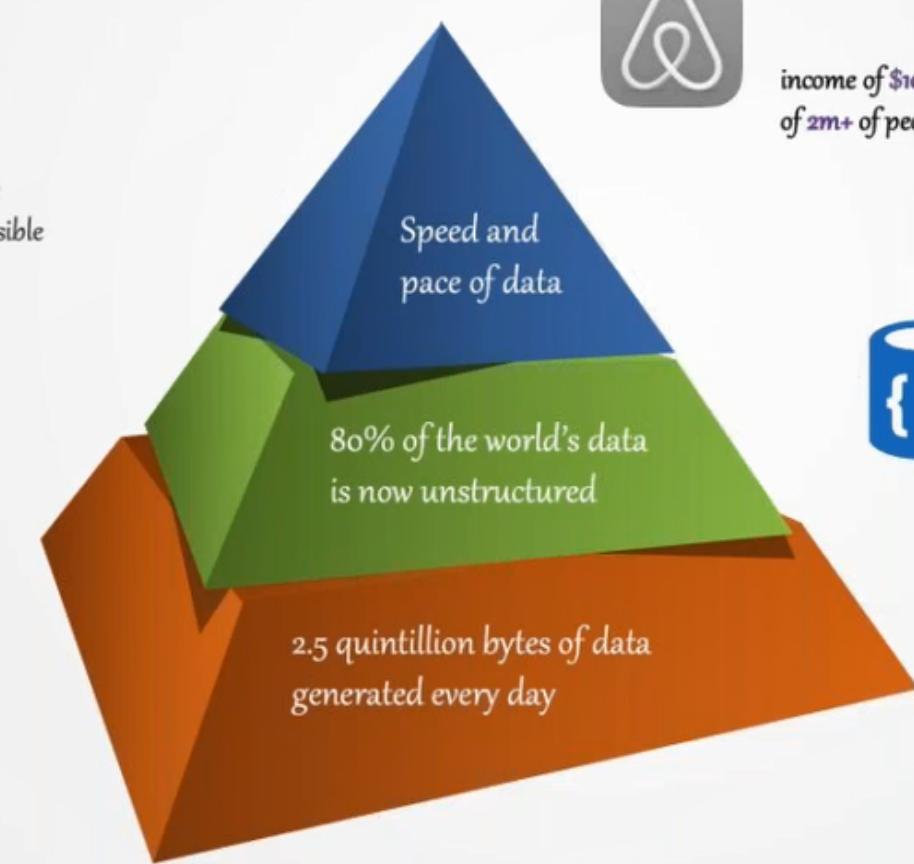
Batch, Near, Real-Time



Apache Hadoop [2006]

quantity of generated & stored data, size of data & value of potential insights

MB | GB | TB | PB



Netflix, Inc.

137 million users worldwide with consumption of 25% of the world's internet bandwidth



Spotify

191 million users worldwide with more than 30 million of songs available



Lyft, Inc.

1m+ million of rides per day and 30M+ of users worldwide



NoSQL [2009]

type & nature of data, different data sources & mappings, easier for developers

Key-Value Pairs, Column-Family, JSON, Graph



Cloud Computing





Google Cloud Platform vs. Microsoft Azure vs. Amazon Web Services

81.7%

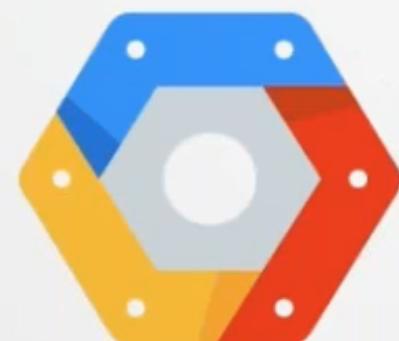
- 2019 = \$9 B
- 2018 = \$3 B

75.9%

- 2019 = \$33 B
- 2018 = \$110 B

46.3%

- 2019 = \$70 B
- 2018 = \$35 B



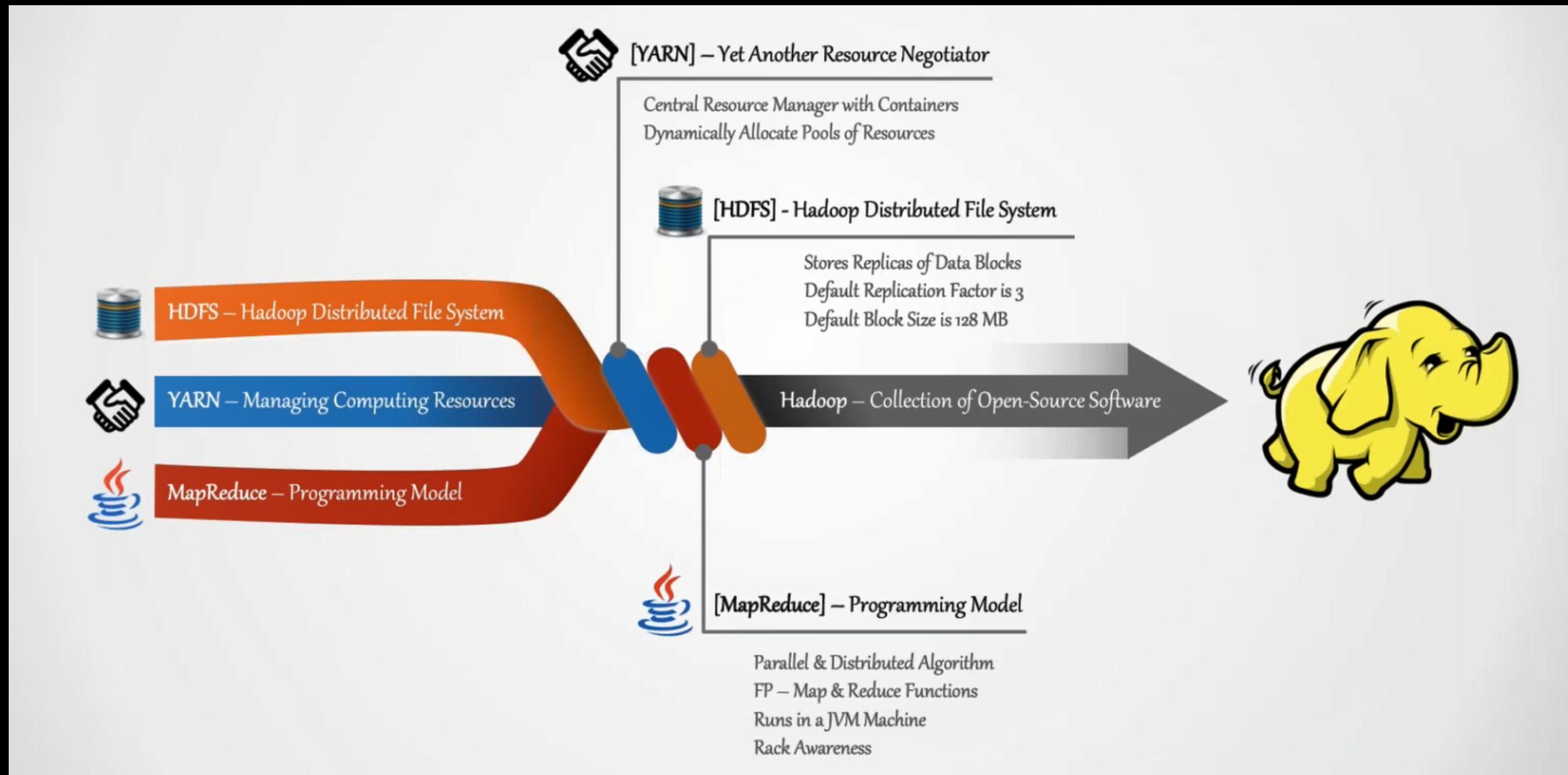


#Spark

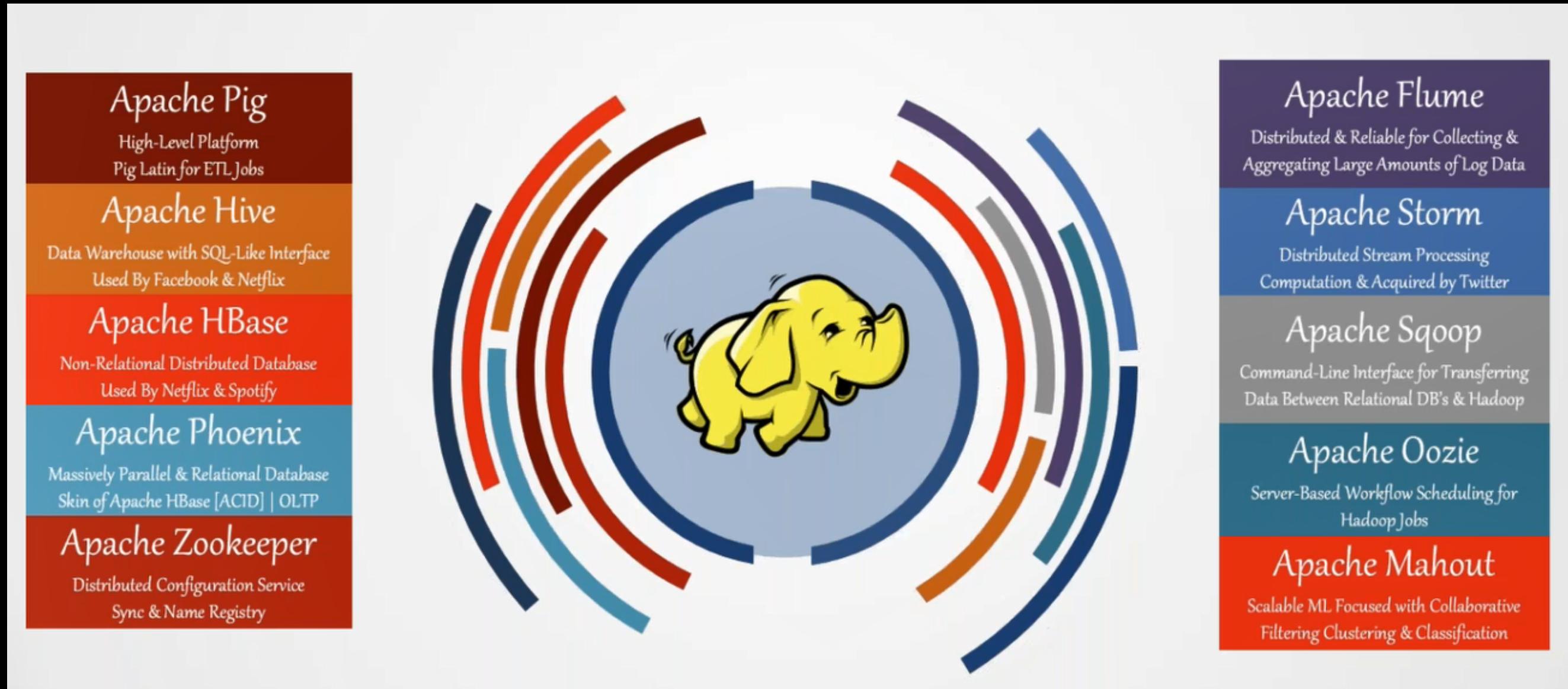
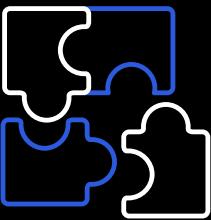
What is the Apache Spark?

Precisamos entender o Hadoop para
entendermos o Spark

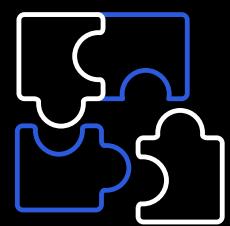
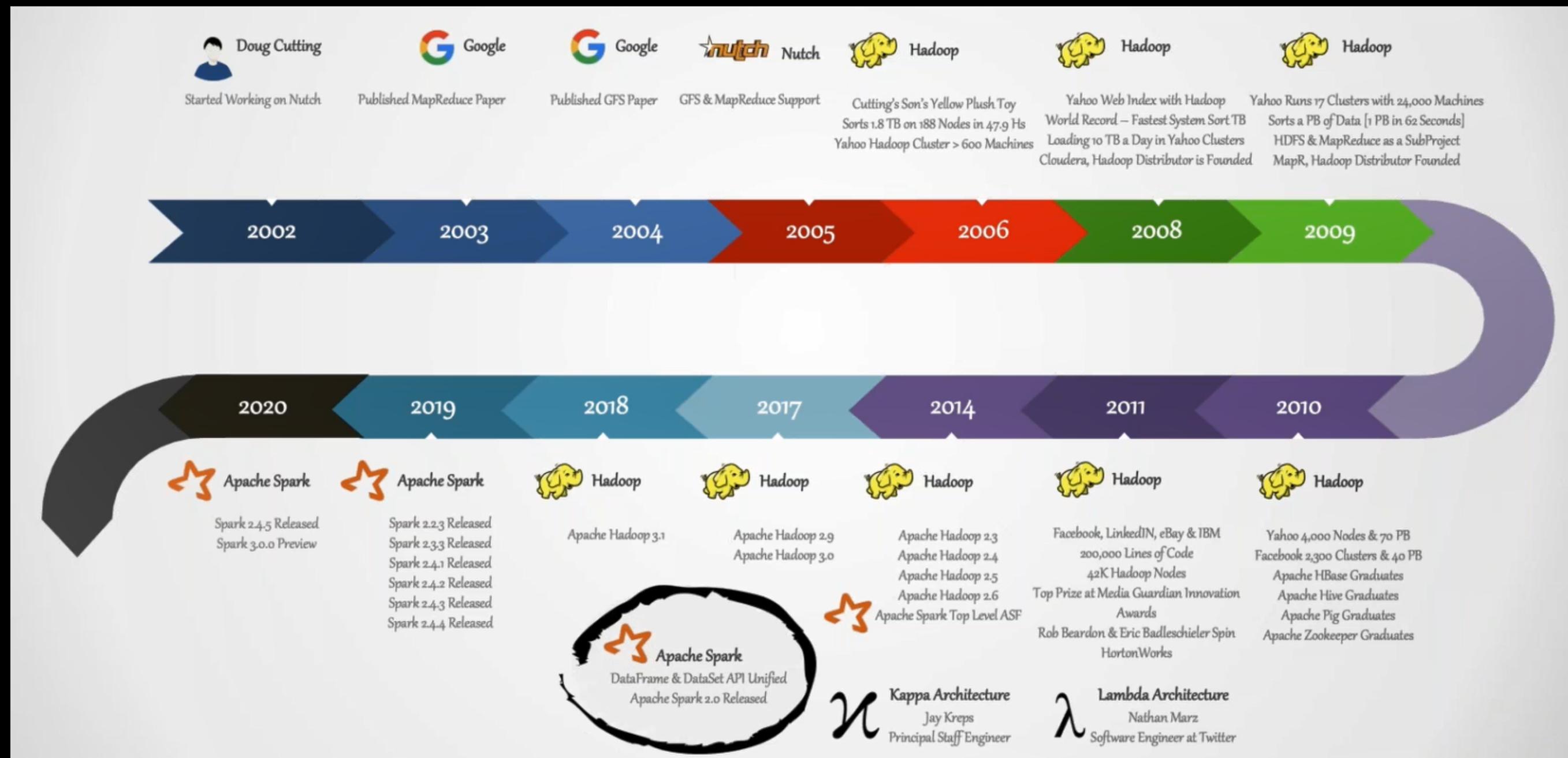
Apache Hadoop [Fundamentals]



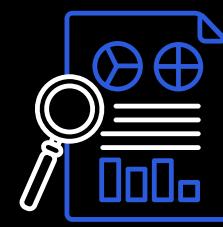
Ecosystem of Hadoop Animal Zoo [2006~2014]



History of Apache Hadoop & Apache Spark

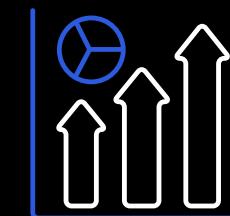


Big Data-as-a-Services [BDaaS]

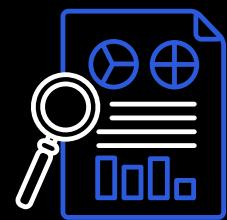


The diagram illustrates the global reach of four Big Data-as-a-Service platforms, each represented by a circular icon and a logo:

- Hadoop-as-a-Service [Haas]**: Collection of Open-Source Software. Fully-Managed Cloud Service. Options: Amazon AWS, Microsoft Azure, Google Cloud Platform.
- Cloud DataProc**: Fully-Managed Cloud Service. Cloud Native Apache Hadoop & Apache Spark. Provisioning Time of 90 Seconds.
- Amazon Elastic MapReduce [EMR]**: Easy Run & Scale Big Data Frameworks. Includes Apache Hadoop, Apache Spark, HBase, Presto & Hive.
- HDInsight**: Easy & Cost-Effective for Open-Source Analytics with Apache Hadoop 3.0. Includes Apache Hadoop, Apache Spark, Apache Kafka, Apache HBase, Apache Hive LLAP, Apache Storm, Machine Learning.



Develop Java MapReduce Program for Apache Hadoop on [HDInsight]



Apache Spark [Fundamentals]



Apache Spark

Open-Source Distributed Cluster-Computing Framework
Implicit Data Parallelism & Fault Tolerance
Optimized for Memory Computation
Written In - Scala
100x - MapReduce Jobs & 10x - Disk-Based Operations



Performance

Daytona Gray

- 100 TB in 23 Minutes with 206 EC2 VMs 
- 100 TB in 72 Minutes with 2,100 Nodes 

Use-Cases

AirBnB – ML & Streaming on Mobile App

Uber – Continuous ETL Pipeline by Processing Apache Kafka Data and Storing on HDFS

Netflix – Data Analytics Computation & ML at Scale

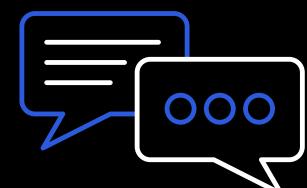
Core APIs

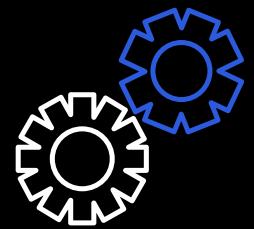
SQL
Java
Scala
Python
R



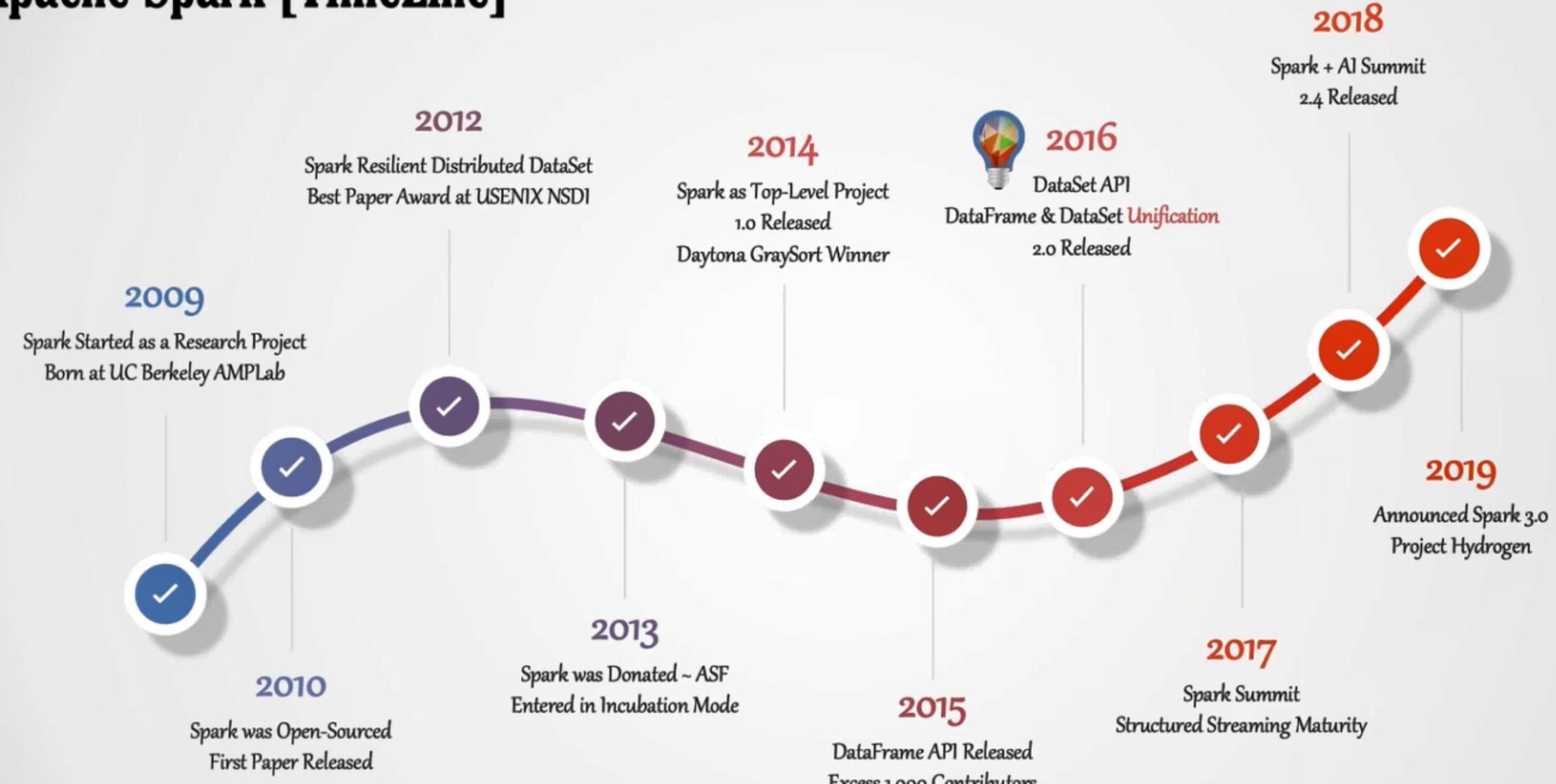
Processing Structures

RDD – Resilient Distributed DataSet
Spark Streaming – Processing Data Streams using DStreams
Spark-SQL, DataSets & DataFrames – Processing Structured Data
Structured Streaming – Processing Structured Data Streams

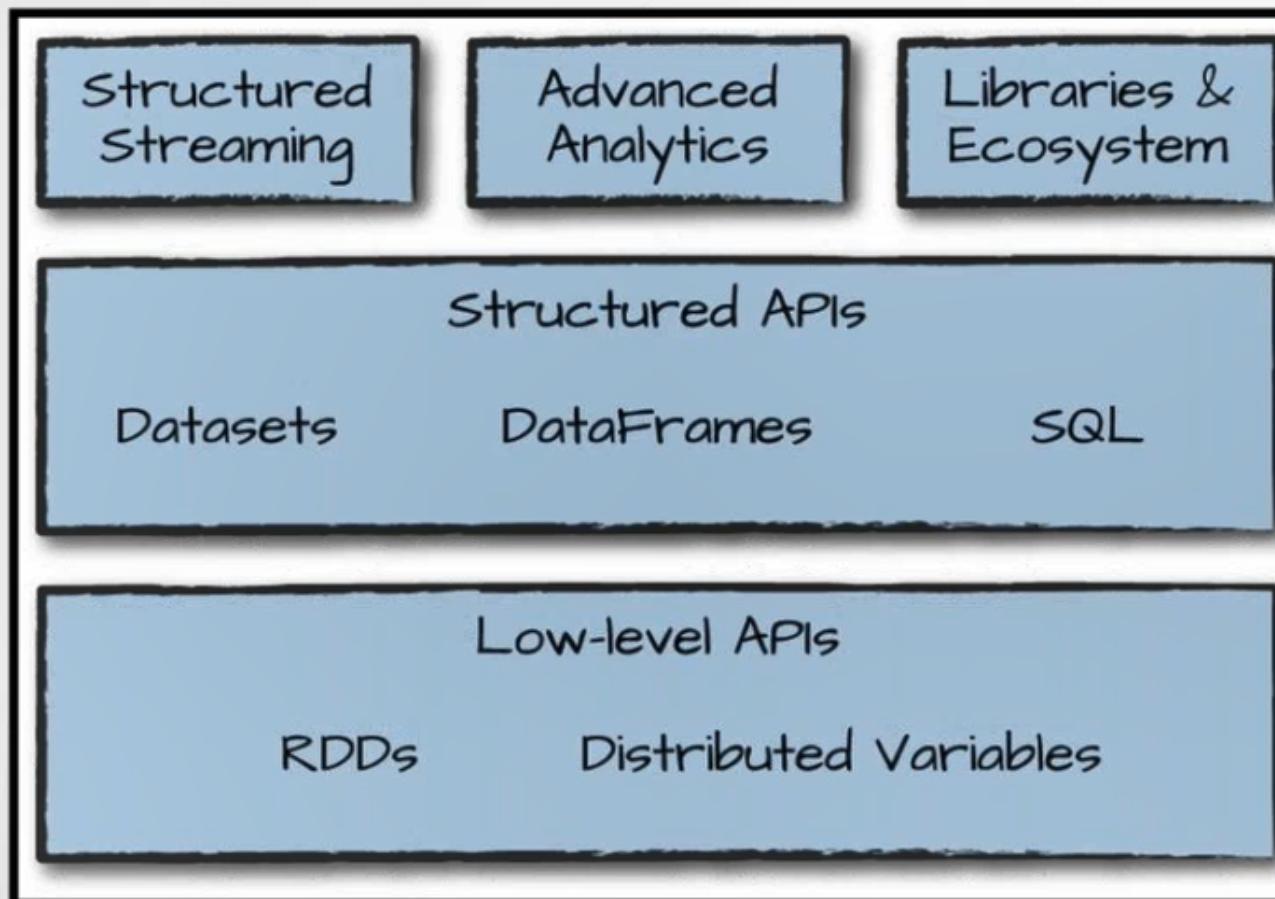




Apache Spark [TimeLine]



Apache Spark [Components]



developers stopped making individual processors faster and switched toward adding more parallel cpu cores all running at the same speed. this change meant that suddenly applications needed to be modified to add parallelism in order to run faster, which set the stage for new programming models such as apache spark.



High-Level APIs

the fundamental abstraction that you will use to write most of your data flows.

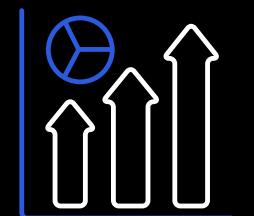
- fundamentally, spark is a distributed programming model in which the user specifies transformations, these transforms build a directed acyclic graph of instructions
- actions, begins the process of executing these dags as a single job by breaking down into stages and tasks to execute across the cluster

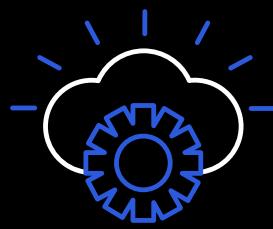


Low-Level APIs

virtually everything in spark is built on top of RDDs.

- you need some functionality that you cannot find in higher level apis for example a very tight control over physical data placement across the cluster
- you need to maintain some legacy codebase written in RDDs
- add some custom shared variable manipulation





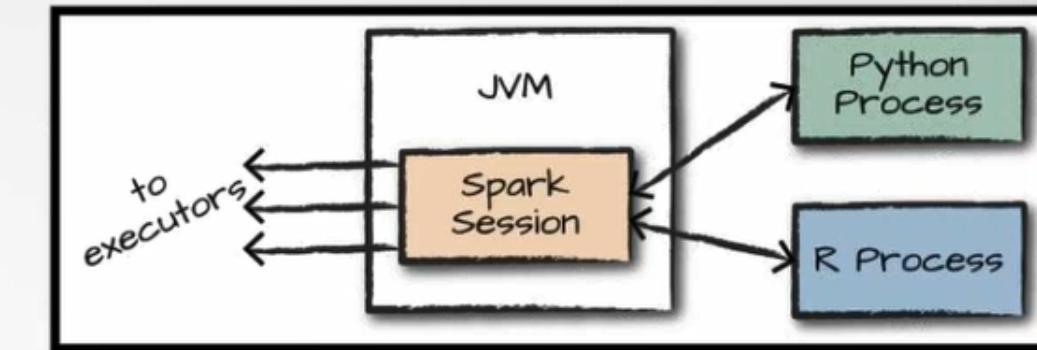
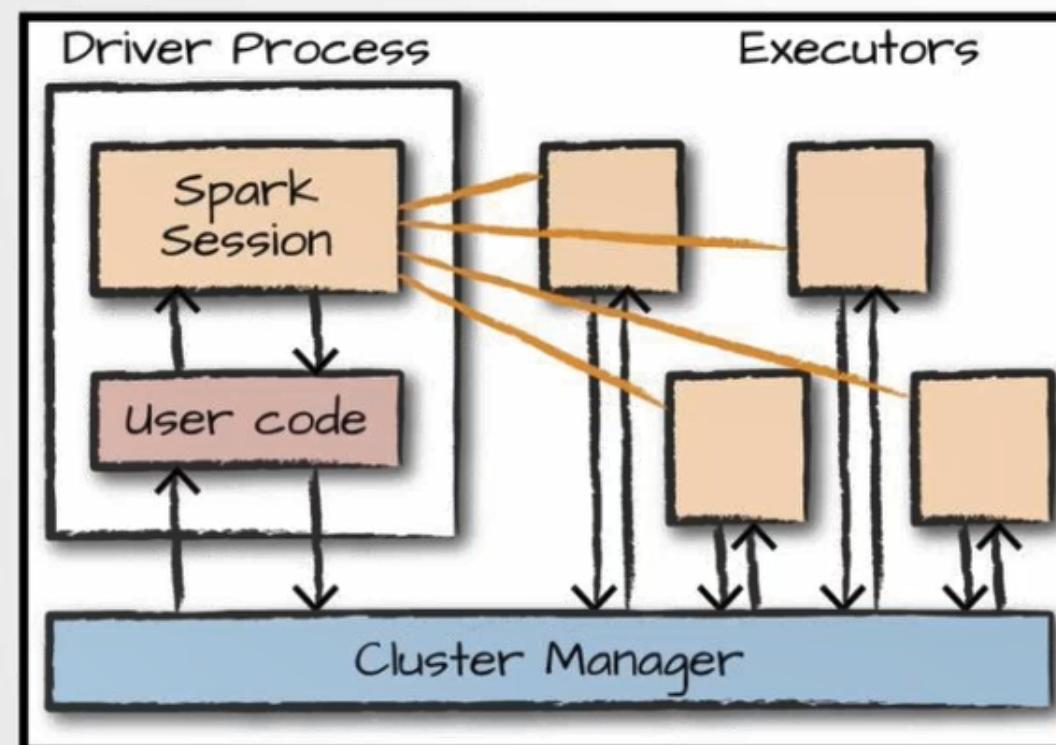
Apache Spark [Architecture]



Driver Process

the driver runs the main function, sits on a node in the cluster, and is responsible for

- maintaining information about the spark application
- responding to a user's program input
- analyze, distribute and schedule work across the executors



each language API maintains the same core concepts that we described earlier. there is a spark session object available to the user, which is the entrance point to running spark code. when using spark from python or r, you don't write explicit JVM instructions; instead, you write python and r code that spark translates into code that it then can run on the executor JVMs.



Executors

responsible for carrying out the work that the driver assigns them

- executing code assigned by the driver
- reporting the state of the computation

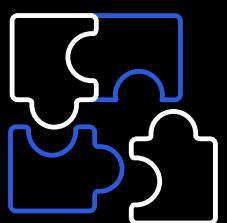


RDD | DataFrame & DataSet [Comparison]

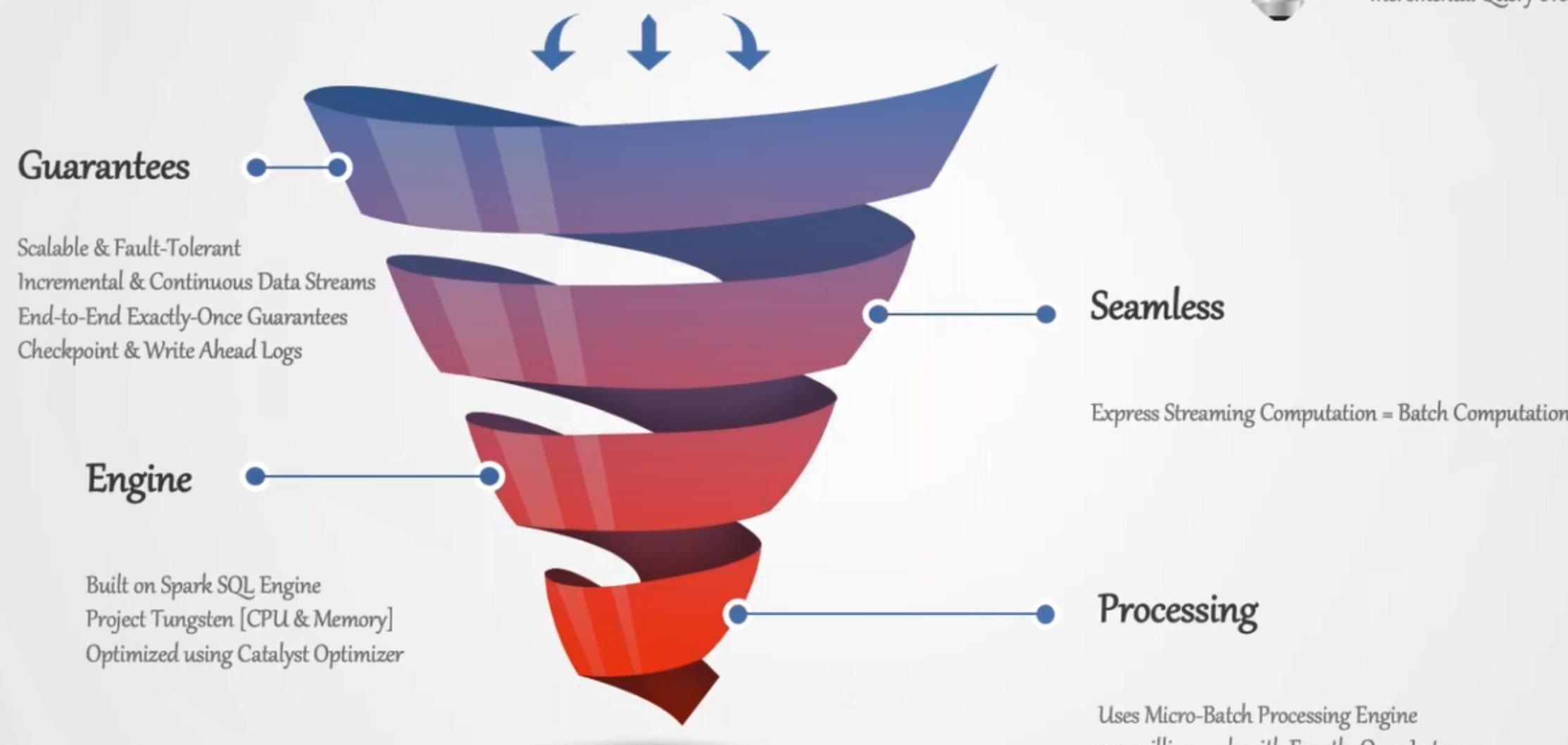


Apache Spark
DataFrame & DataSet API Unified
Apache Spark 2.0

	RDD	DataFrame	DataSet
› Structured & Unstructured	✓	✓	✓
› Java & Scala	✓	✓	✓
› Python & R	✓	✓	✗
› Any Data Source	✓	✗	✓
› Schema Infer	✗	✓	✓
› Optimization Engine	✗	✓	✓
› Fast Aggregation	✗	✓	✓
› In-Memory Serialization	✗	✓	✓

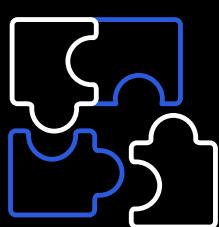


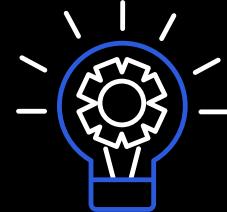
Apache Spark [Structured Streaming]



Philosophy [Moto]

Treat Data Streams = Unbounded Tables
Incremental Query over Streams



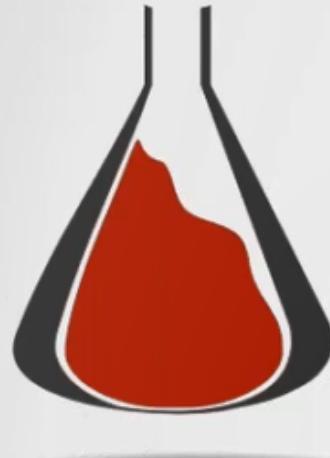


Spark ML [MLlib]



MLlib – Apache Spark's Scalable Machine Learning Library

RDD-Based API vs. DataFrame-Based API



MLlib RDD-Based API [`spark.mllib`] Package [Maintenance Mode]
DataFrame-Based API [`spark.ml`] Package

DataFrame-Based API [Spark-ML]

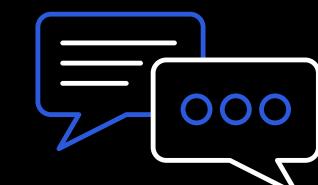
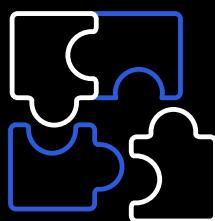
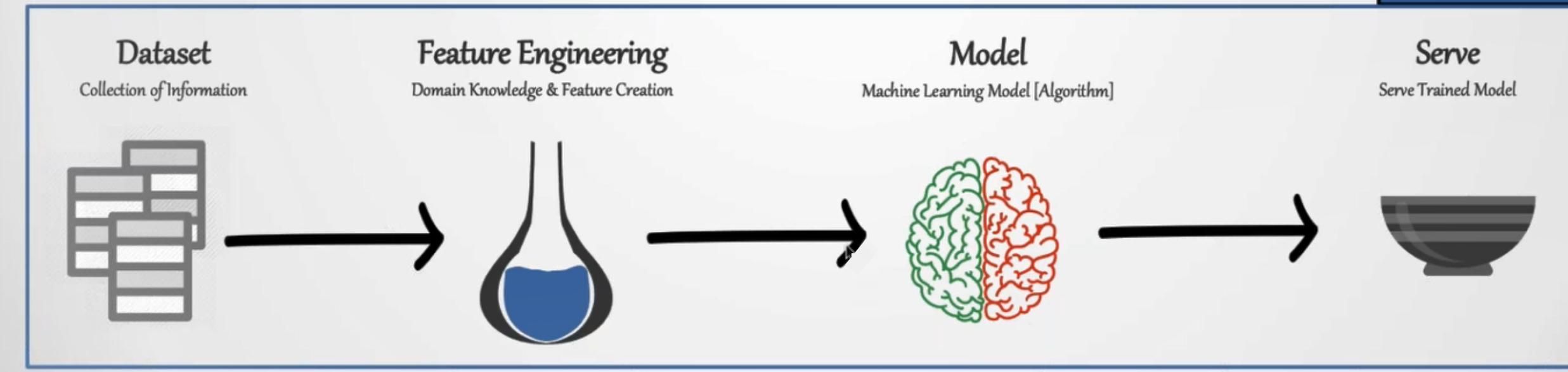
- User-Friendly API
- Data Sources
- SQL & DataFrame Queries
- Tungsten & Catalyst Optimizations
- Facilitate Practical ML Pipelines [Feature Transformation]



ML Library [MLlib]

ML Algorithms – Classification | Regression | Clustering & Collaborative Filtering
Featurization – Extraction | Transformation | Dimensionality Reduction & Selection
Pipelines – Constructing | Evaluating & Tuning ML Pipelines
Persistence – Saving & Loading Algorithms, Models & Pipelines
Utilities – Linear Algebra | Statistics & Data Handling

ML Pipeline

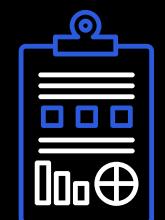
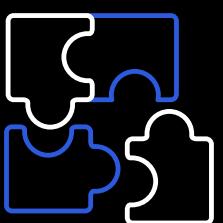


RDD | DataFrame & DataSet [Comparison]



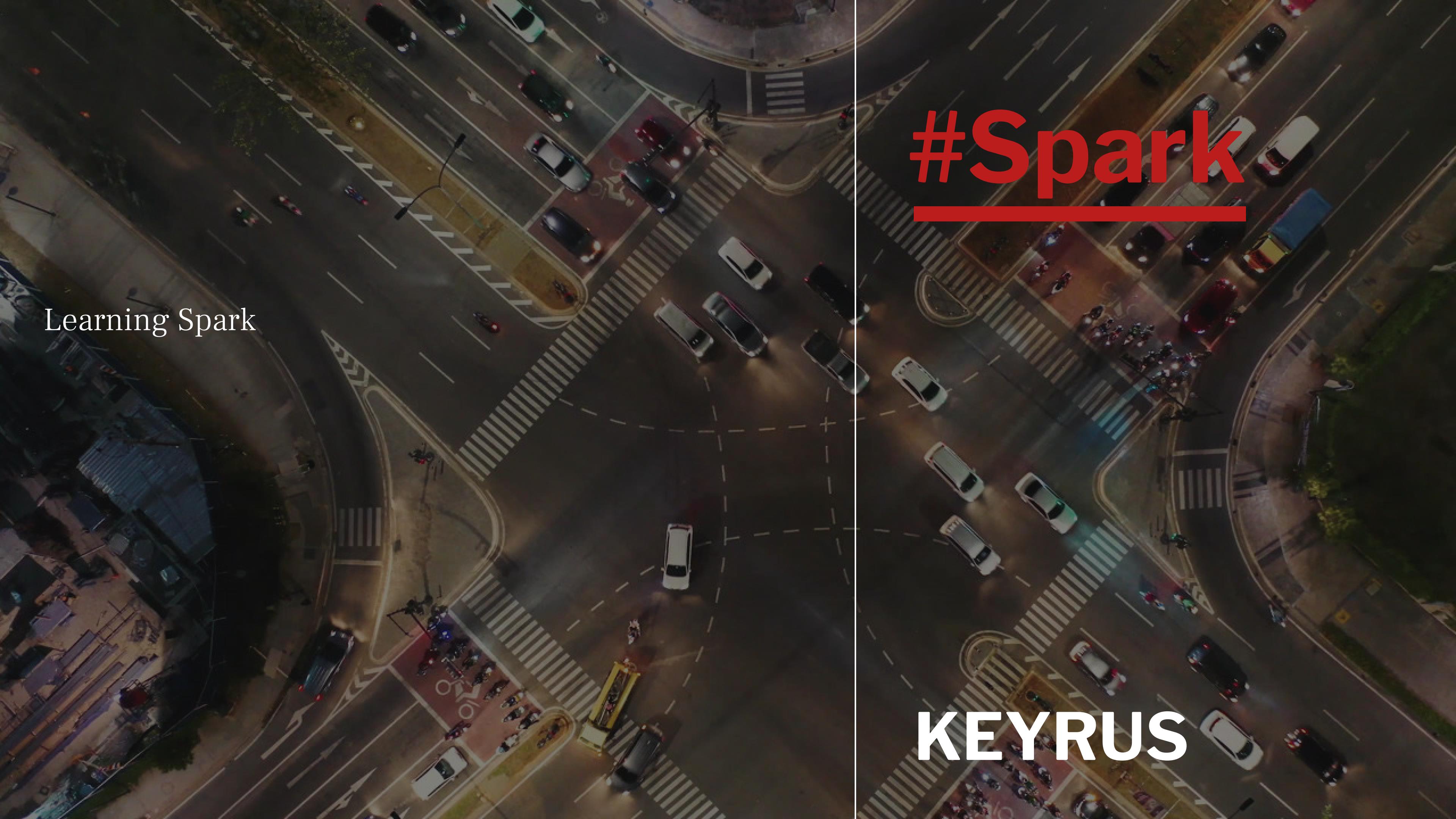
Apache Spark
DataFrame & DataSet API Unified
Apache Spark 2.0

	RDD	DataFrame	DataSet
› Structured & Unstructured	✓	✓	✓
› Java & Scala	✓	✓	✓
› Python & R	✓	✓	✗
› Any Data Source	✓	✗	✓
› Schema Infer	✗	✓	✓
› Optimization Engine	✗	✓	✓
› Fast Aggregation	✗	✓	✓
› In-Memory Serialization	✗	✓	✓



Develop PySpark Program for Apache Spark on [HDInsight]



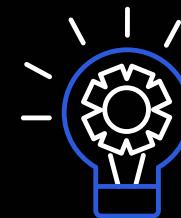
An aerial photograph of a multi-lane road intersection at night. The scene is illuminated by streetlights and vehicle headlights, showing several cars and a few people on the sidewalks. Buildings are visible along the left side of the street.

Learning Spark

#Spark

KEYRUS

Conteúdo Programático



DATA ENGINEERING

DATALAKE

SPARKSQL

DATAFRAME

EXECUTANDO SPARK EM UM CLUSTER

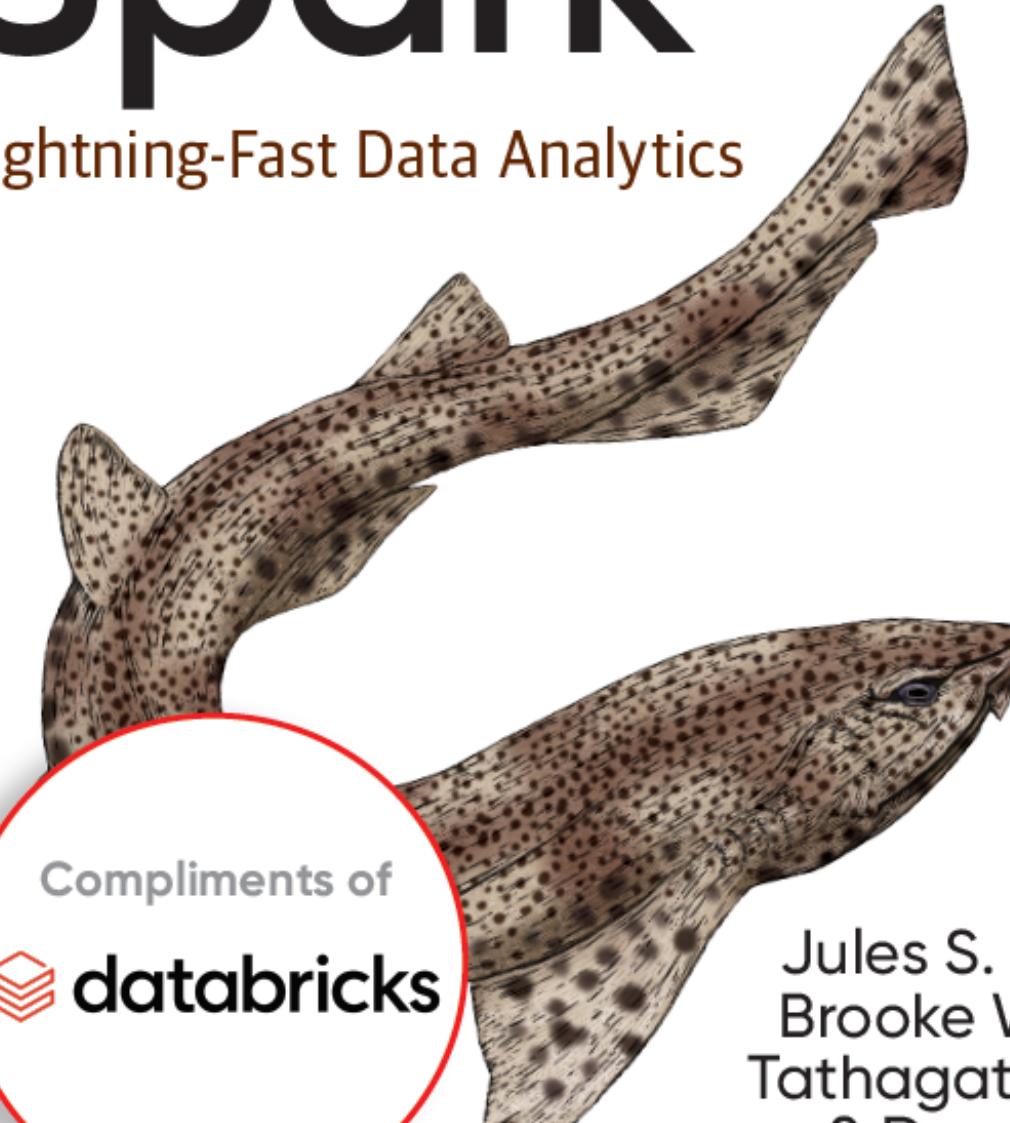
DATABRICKS

O'REILLY®

2nd Edition
Covers
Apache Spark 3.0

Learning Spark

Lightning-Fast Data Analytics



Compliments of



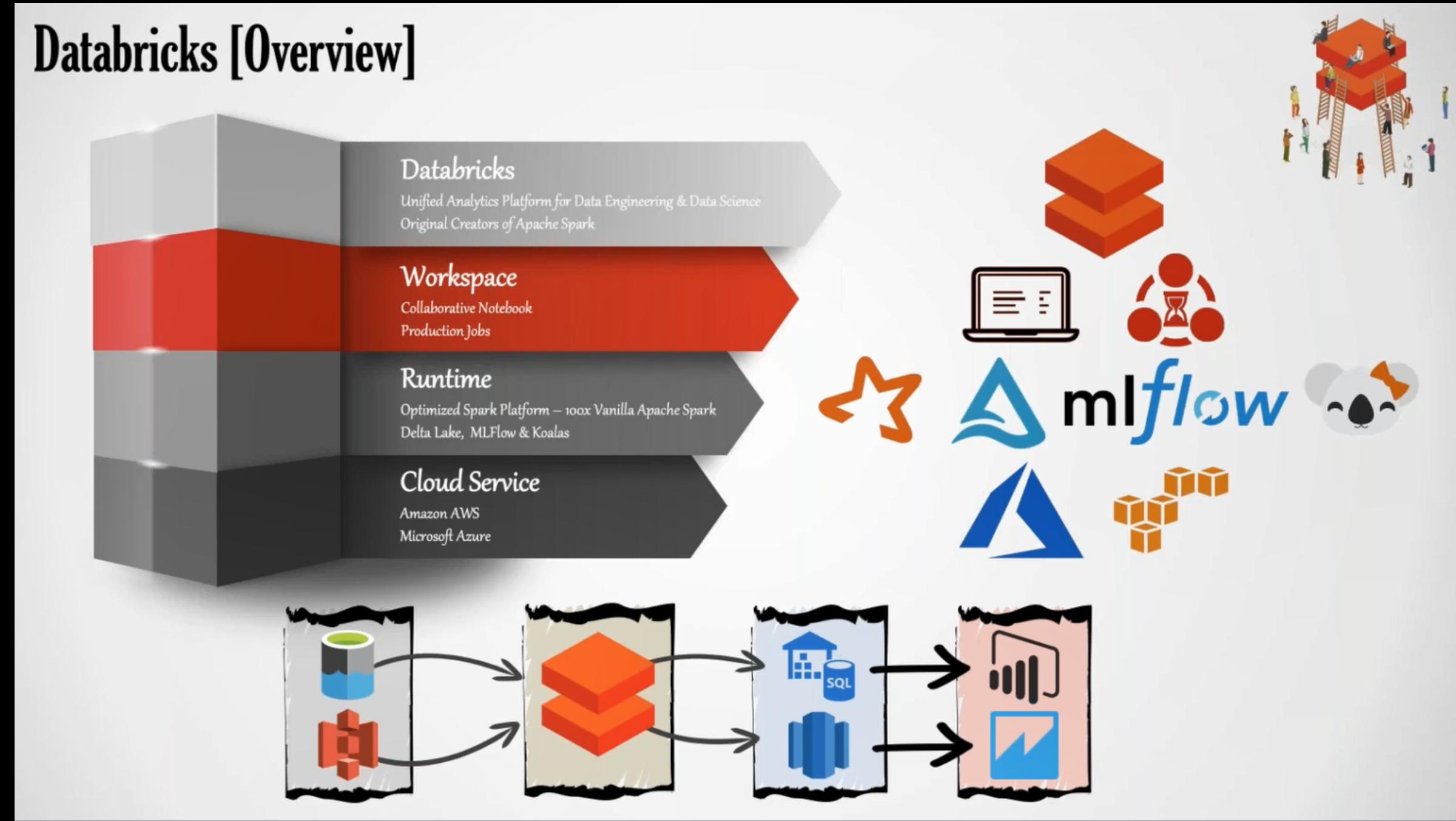
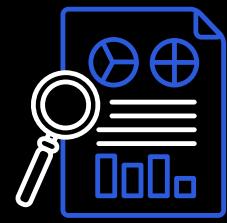
Jules S. Damji,
Brooke Wenig,
Tathagata Das
& Denny Lee

Foreword by Matei Zaharia

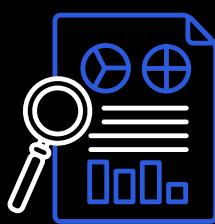
Using the Scala or PySpark Shell on [HDInsight]



Databricks [Overview]



Databricks [TimeLine]



Databricks Founded
Creators of Apache Spark
Raise \$14M in Series A Funding

2013

Databricks Cloud Publicly Available
Amazon AWS

2014

2015

2016

2017

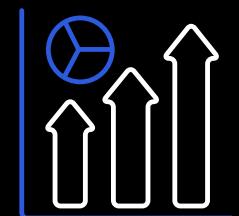
2018

2019

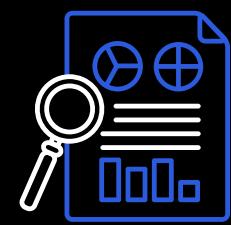
Databricks Cloud Announced
Amazon AWS Only
Raise \$33M in Series B Funding

Databricks Debuts Free Community Edition
Raise \$60M in Series C Funding

Microsoft Announce GA of
Azure Databricks



Getting Started with Databricks

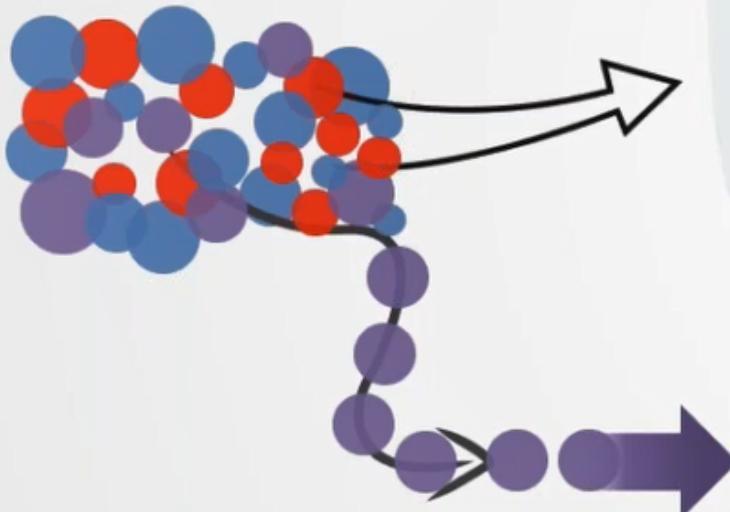


Lambda Architecture



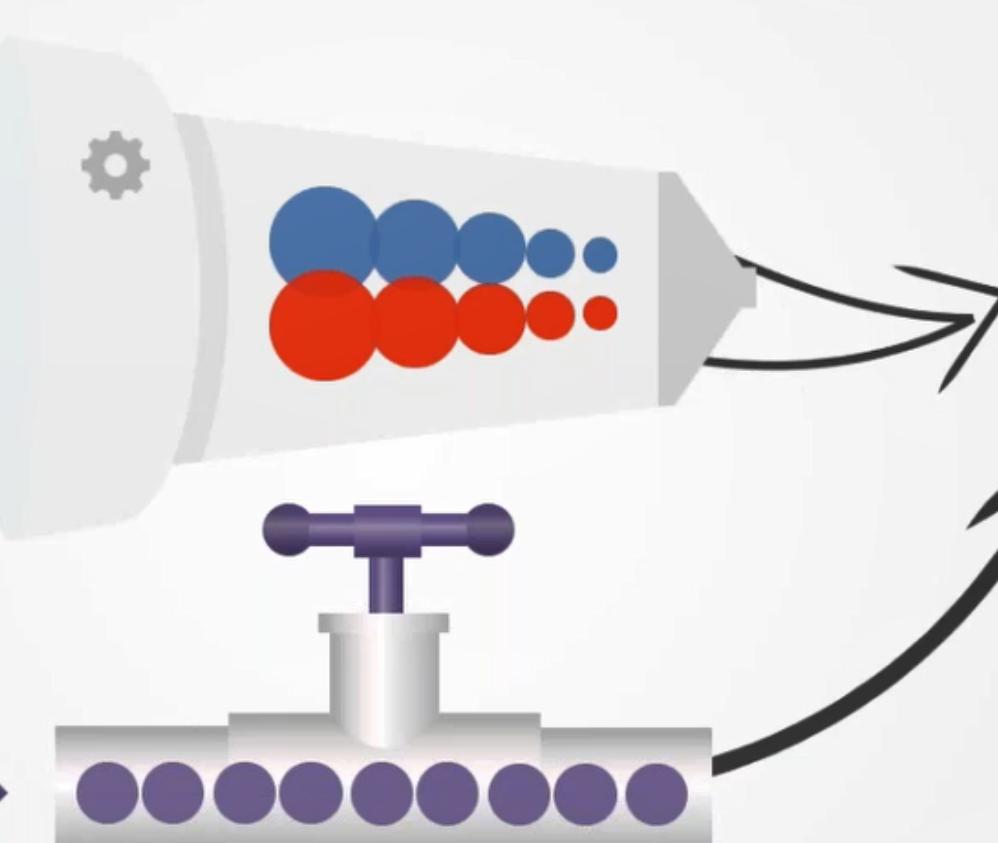
Data Source

Relational & NoSQL DBs
Web & Mobile Apps
CSV, Logs, XML, JSON
Emails, Documents, Audio, Video



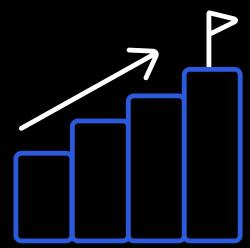
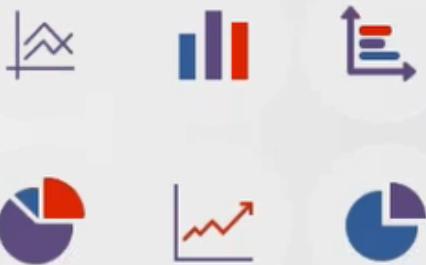
Speed-Layer

Real-Time Ingestion - Feed & Store Data in Real-Time
Stream Processing - Compute & Transform Inbound Data



Serving-Layer

Output of Batch & Speed Layer
Processed & Computed Data
Ad-Hoc Queries & Data Analysis



Batch-Layer

Data Storage - Data Stored in a Data Lake
Batch-Processing - New Data is Computed & Processed



Nathan Marz

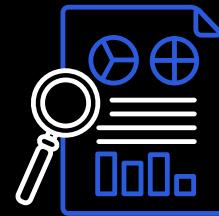
Coined in 2011
Founder of Red Planet Labs
Software Engineer at Twitter

Lambda Architecture – Cloud Agnostic & Simplified



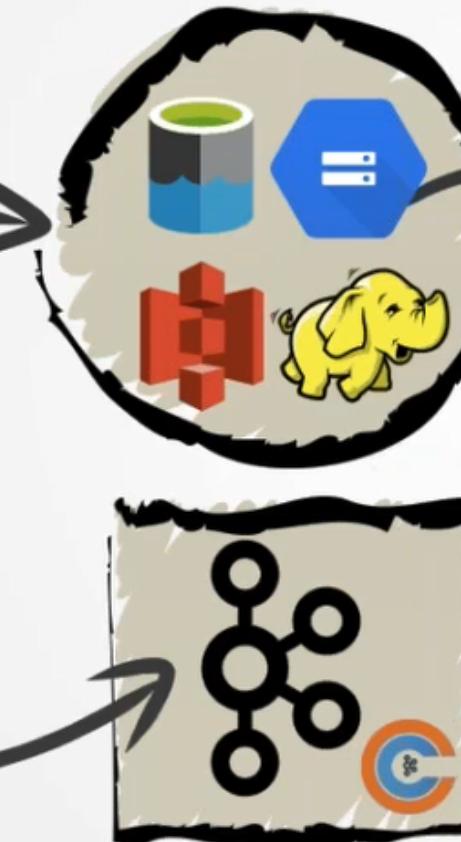
Data Source

JSON & CSV
SQL Server & Redis
Internet – Siri | Spotify | YouTube



Batch-Layer

Data Storage - Data Lake Storage Gen2 | GCS | S3 | HDFS
Batch-Processing - Apache Spark | Databricks



Speed-Layer

Real-Time Ingestion - Apache Kafka [Confluent]
Stream Processing - Apache Kafka [Confluent] | Apache Spark [Databricks]



Serving-Layer

Azure SQL Data Warehouse
Amazon Redshift
Google BigQuery
Apache Hive
CosmosDB



Kappa Architecture

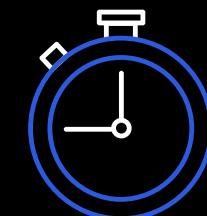
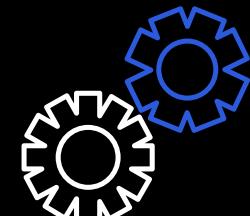
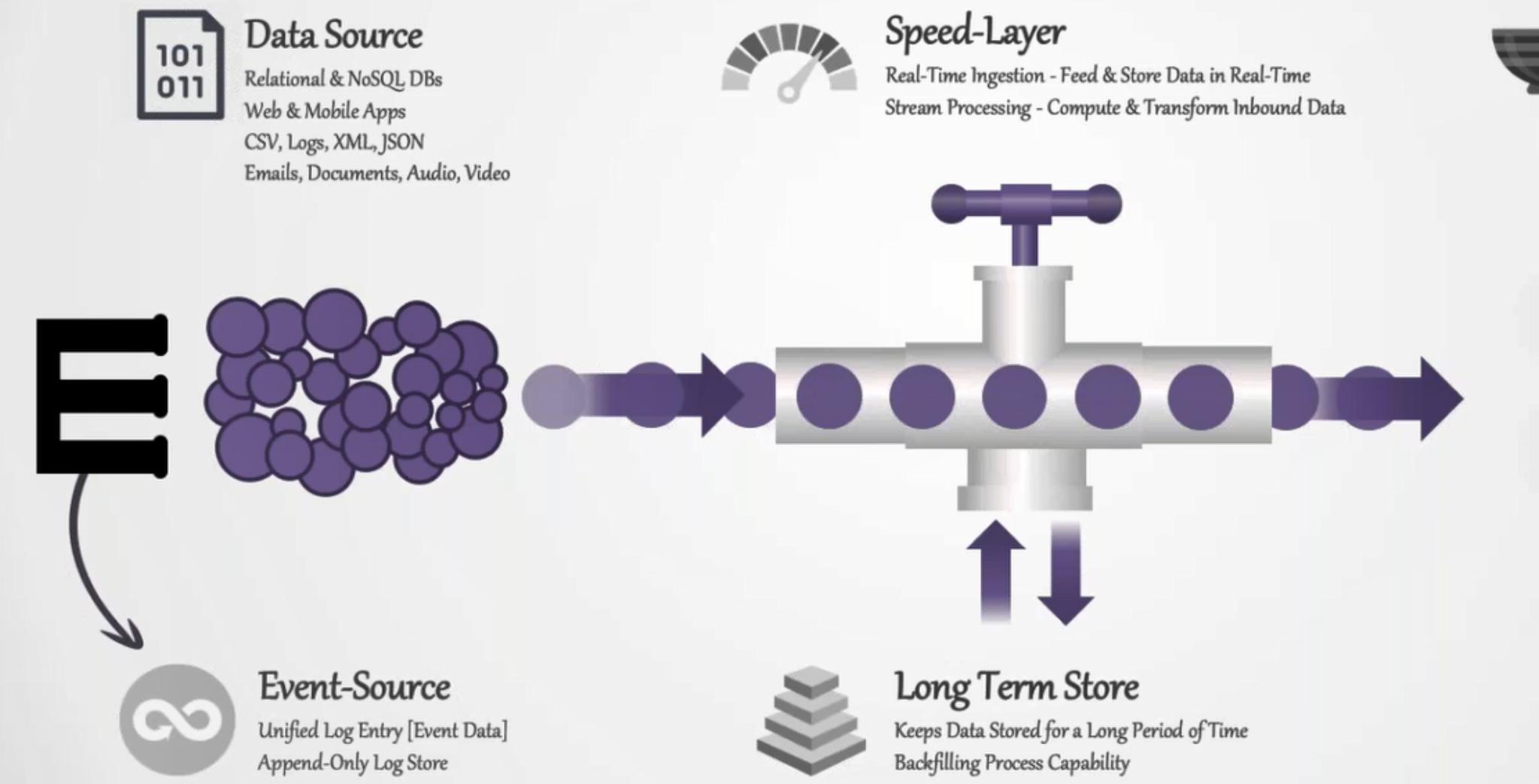


Jay Kreps

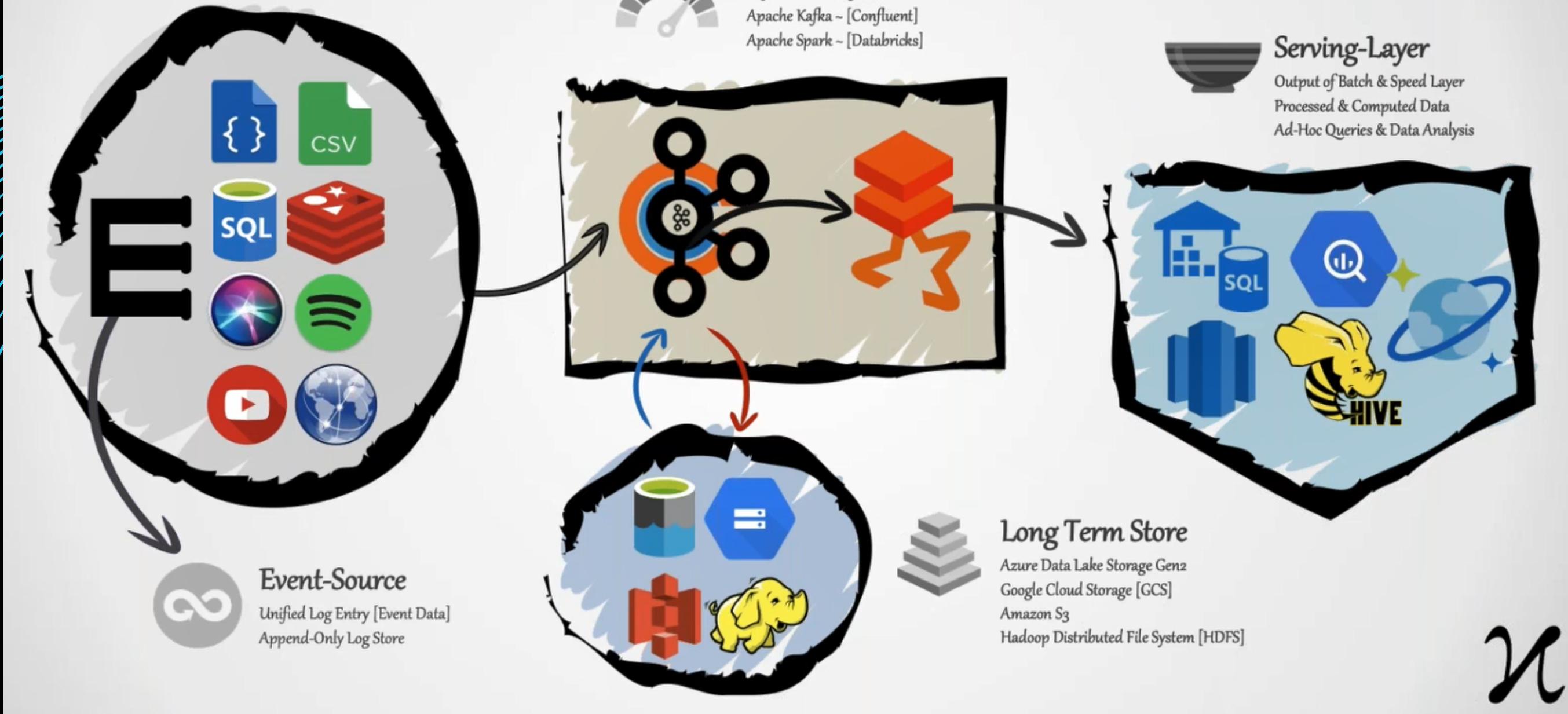
Coined in 2014

Co-Founder & CEO at Confluent

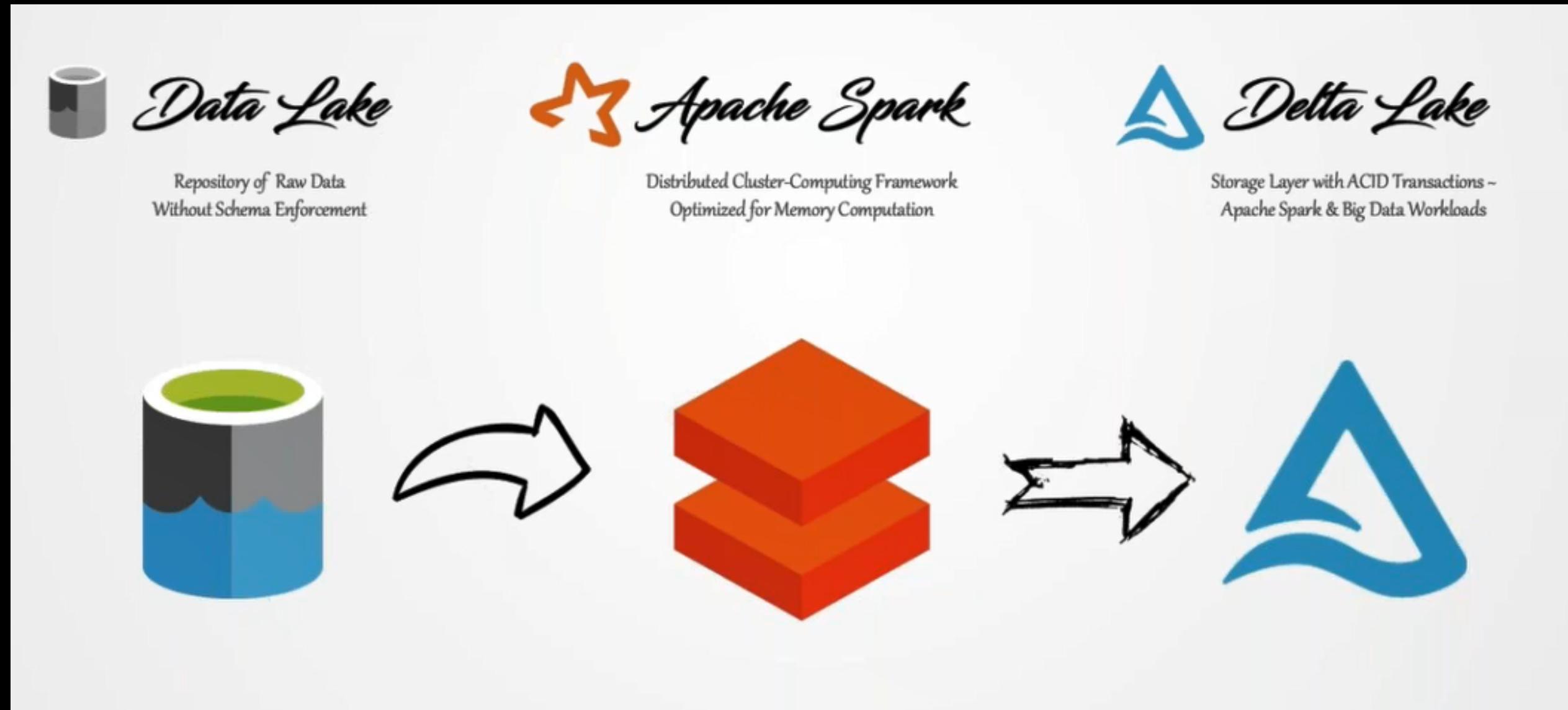
Principal Staff Engineer

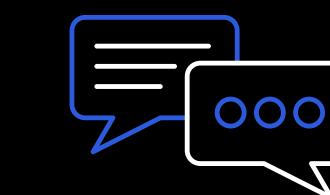
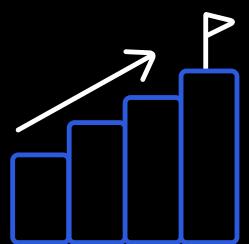
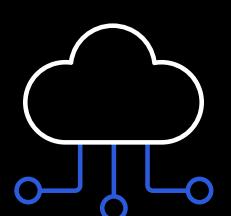


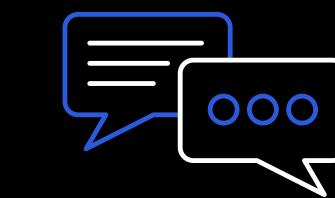
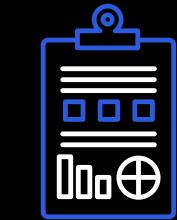
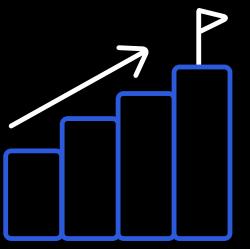
Kappa Architecture – Cloud Agnostic & Simplified



Use Case ~ Batch-ETL

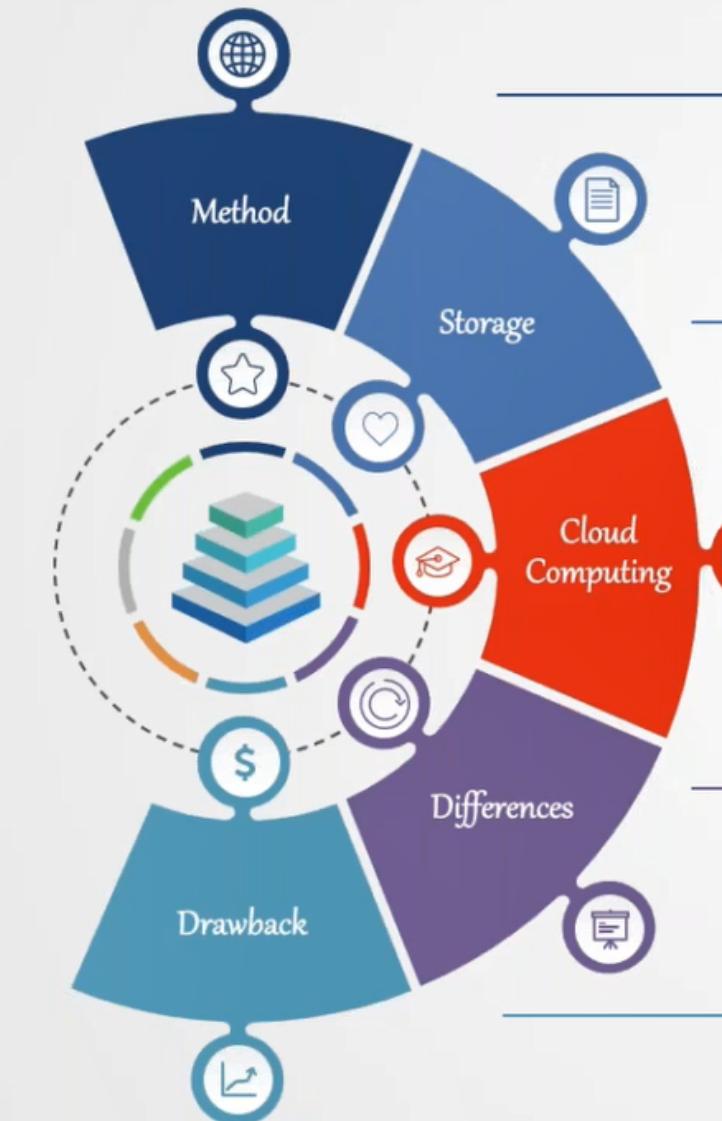








Data Lake [Concepts]



Repository of Raw Data [Data Democratization]
Data Lake – James Dixon in 2010
Democratization of Data – [Unsiloeed Data]



Structured – Relational Databases
Semi-Unstructured – CSV, Logs, XML & JSON
Unstructured – Emails, Documents, Binaries, Audio & Video



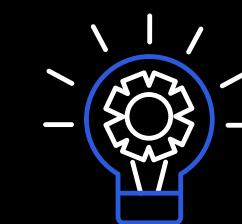
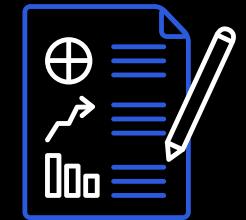
Azure – Azure Blob Storage & Azure Data Lake Storage Gen2
AWS – Amazon S3 & AWS Lake Formation
GCP – Google Cloud Storage



Data Mart – Subset of Decision Support for Departments
Data Warehouse – Decision Making Support for Enterprise-Level
Data Lake – Raw Data [Source of Truth]



Data Swamp – Highly Disorganized Data
Data Governance – Management & Control of Data
Data Quality – Accuracy and Data Veracity
Data Security – Protecting Digital Data



Data Lake [Use-Cases]



Mobile Application
Sending Data Periodically
Storing Data in TXT Format



MongoDB
Sending Data in Near Real-Time
Storing Data in JSON Format



SQL Server
Sending Data in Batches [2 Hs]
Storing Data in Parquet Format



SFTP
Sending Data in Batches [12 Hs]
Storing Data in CSV Format



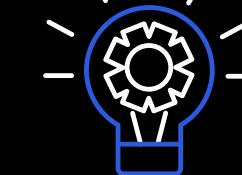
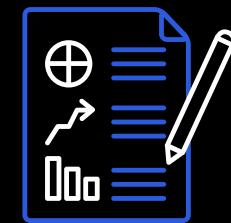
Data Processing [DP]

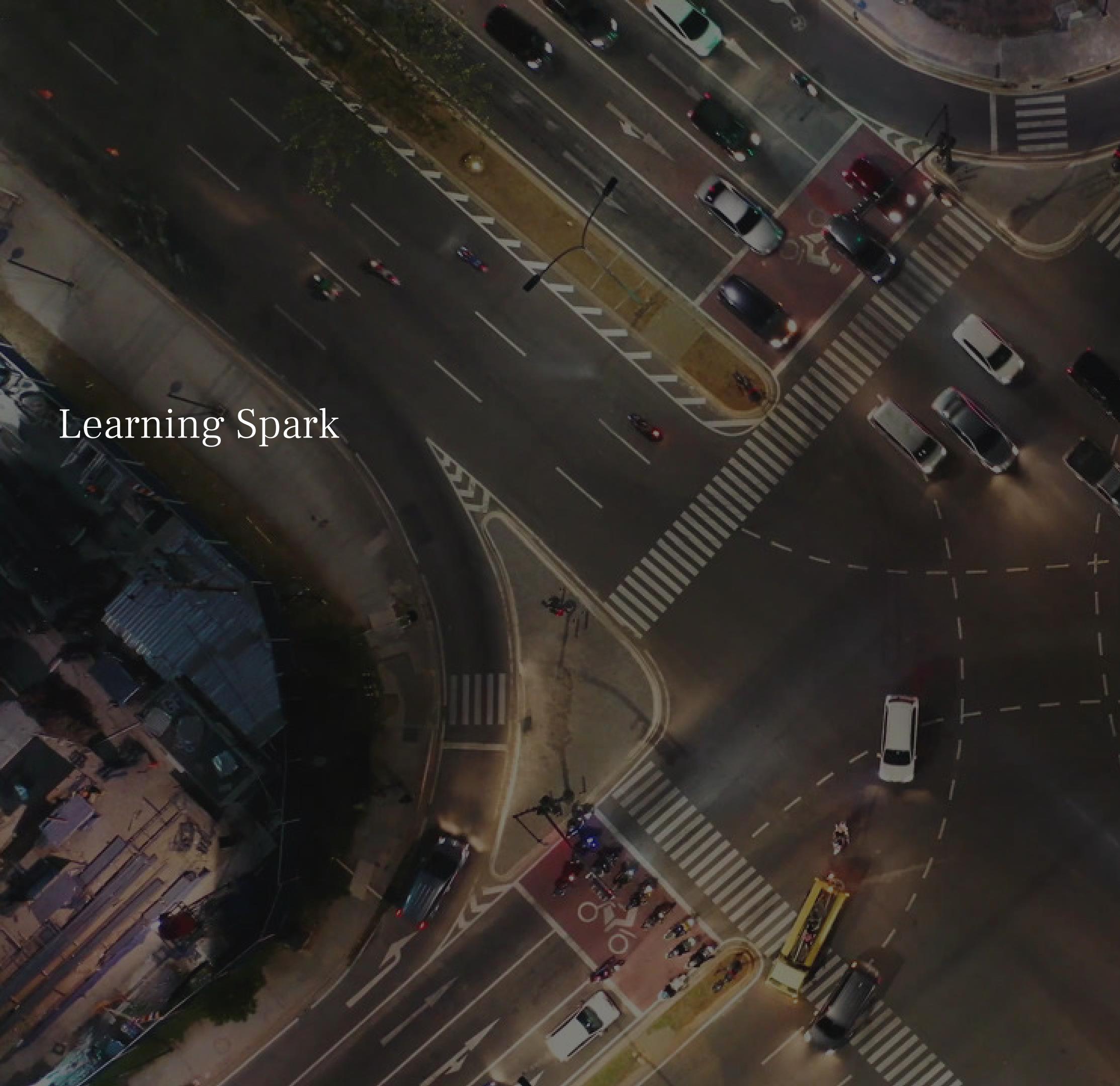
Data Filtering & [ELT]

Machine Learning [ML]

Data Warehouse [Dw]

Data Visualization [DataViz]





Learning Spark

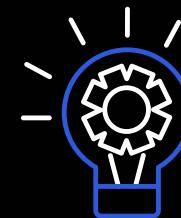


#Spark



KEYRUS

Conteúdo Programático



DATA ENGINEERING

DATA LAKE

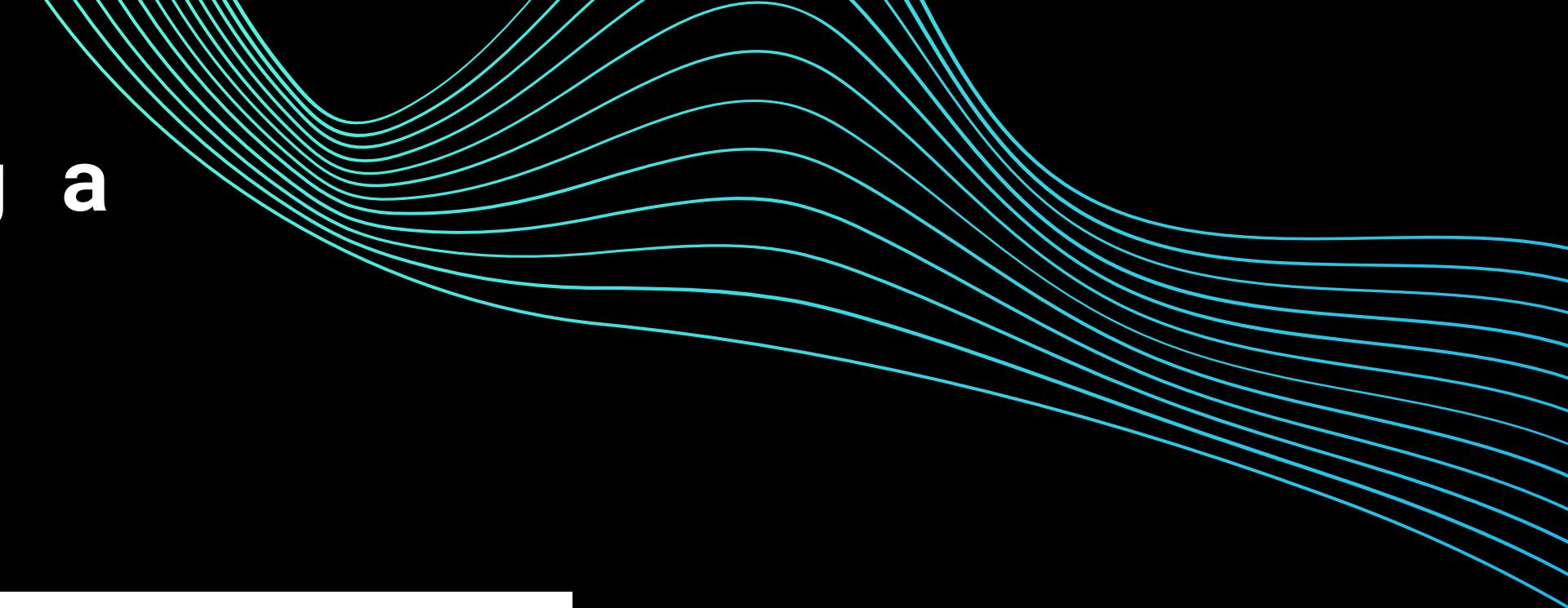
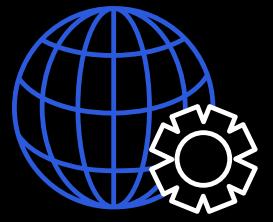
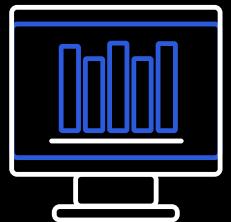
DELTA LAKE

DESIGN PATTERNS

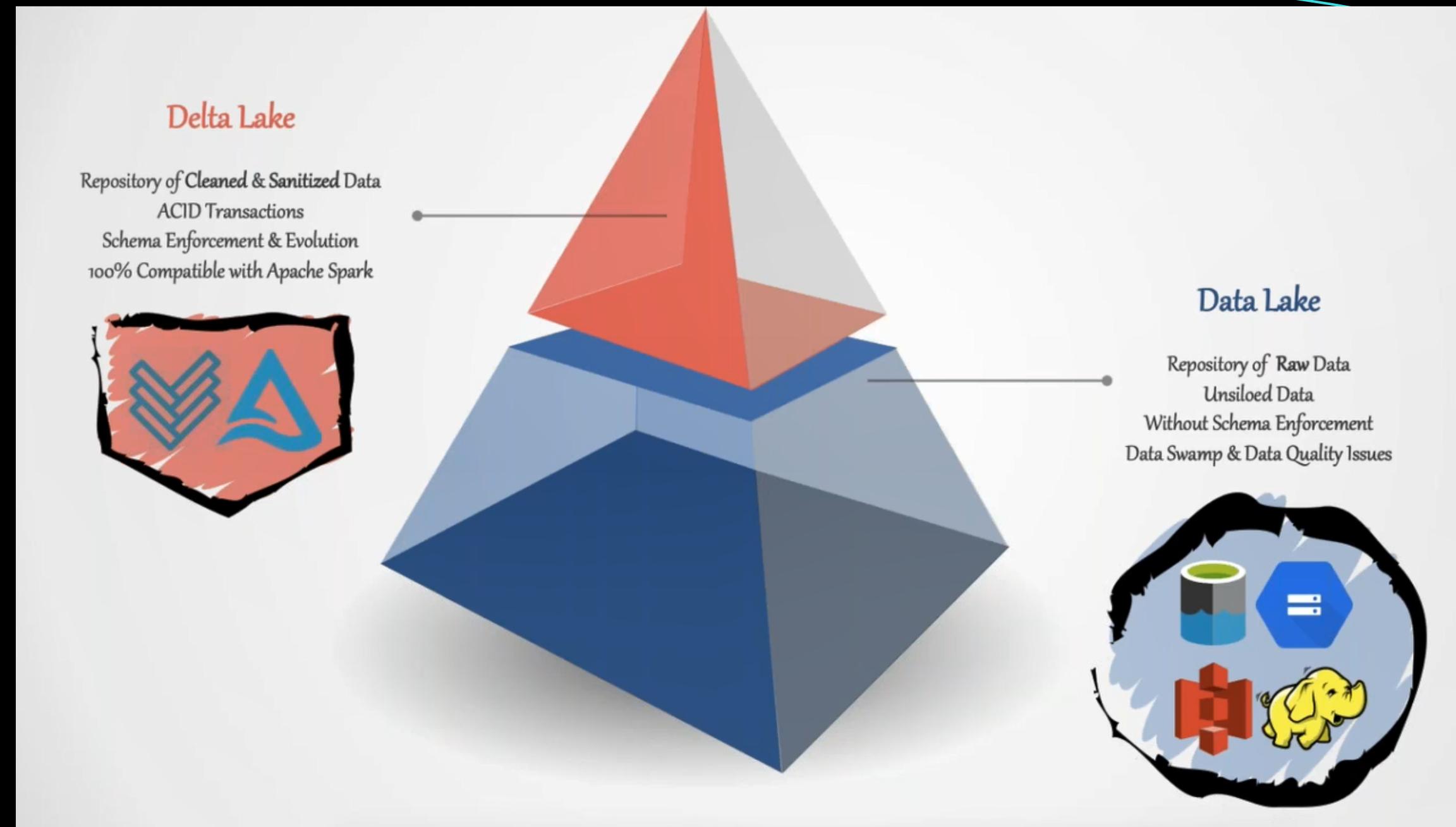
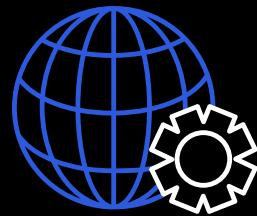
EXECUTANDO SPARK EM UM CLUSTER

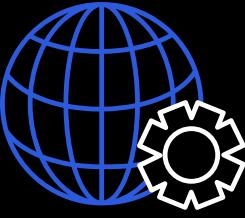
DEMO ON DATABRICKS

Best Practices Organizing a Data Lake Repository



Data Lake vs. Delta Lake





Delta Lake [Features]



Delta Lake

Is an Open-Source Storage Layer ~ ACID Transactions to Apache Spark & Big Data Workloads



unified analytics engine for big data

storage layer [format] with acid capabilities



most common data lakes

ACID Transactions

data lakes typically have multiple data pipelines reading and writing data concurrently, and data engineers must go through a tedious process to ensure data integrity, due to the lack of transactions. it provides serializability, the strongest level of isolation level.

Scalable Metadata Handling

delta lake treats metadata just like data, leveraging spark's distributed processing power to handle all its metadata. as a result, delta lake can handle petabyte-scale tables with billions of partitions and files at ease.

Time Travel [Data Versioning]

delta lake provides snapshots of data enabling developers to access and revert to earlier versions of data for audits, rollbacks or to reproduce experiments.

Open Format

all data in delta lake is stored in apache parquet format enabling delta lake to leverage the efficient compression and encoding schemes that are native to parquet.

Unified Batch & Streaming

a table in delta lake is both a batch table, as well as a streaming source and sink. streaming data ingest, batch historic backfill, and interactive queries all just work out of the box.

Schema Enforcement & Evolution

delta lake provides the ability to specify your schema and enforce it. this helps ensure that the data types are correct and required columns are present, preventing bad data from causing data corruption. delta lake enables you to make changes to a table schema that can be applied automatically, without the need for cumbersome dll.

Audit History

delta lake transaction log records details about every change made to data providing a full audit trail of the changes.

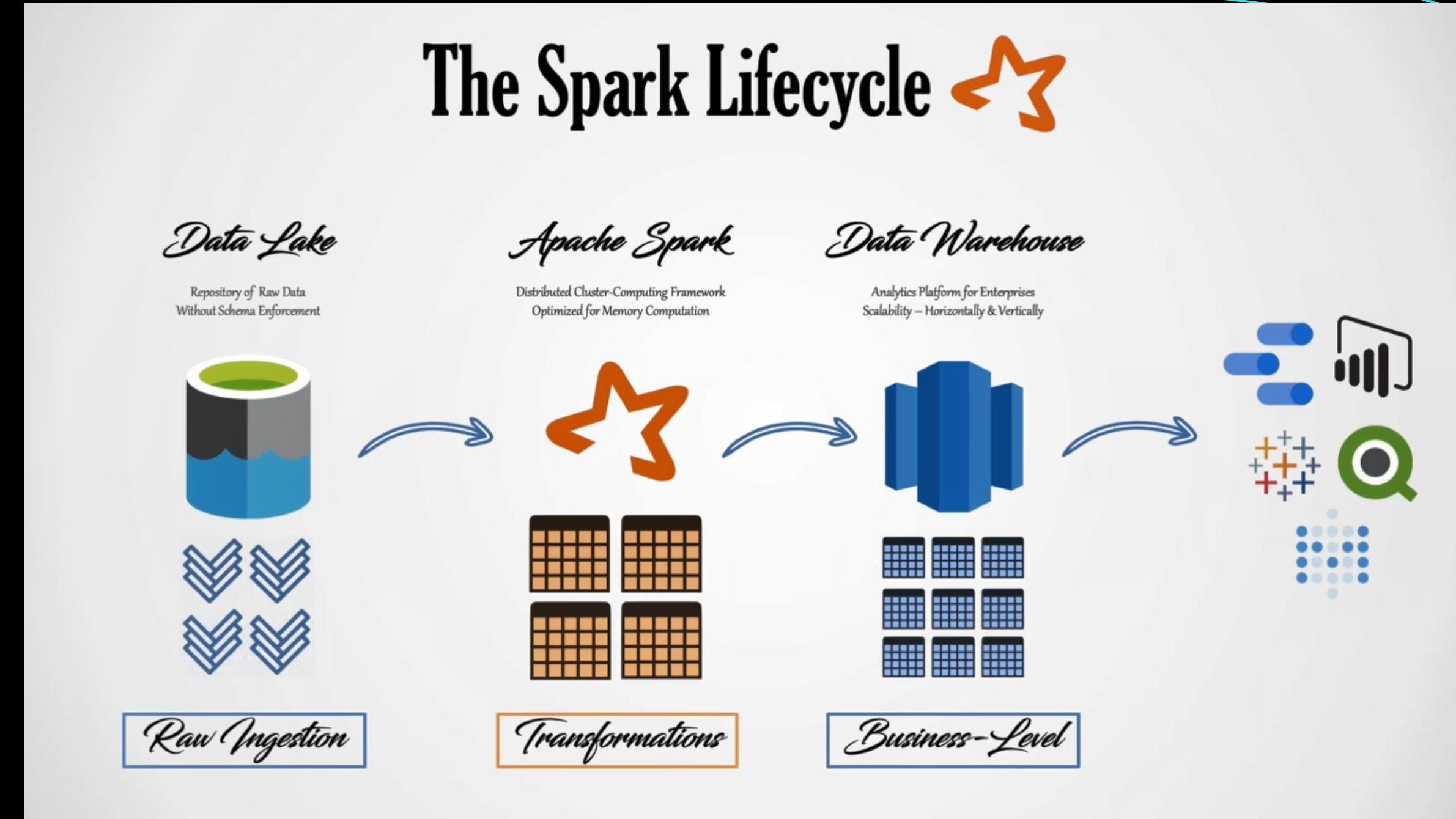


Updates & Deletes

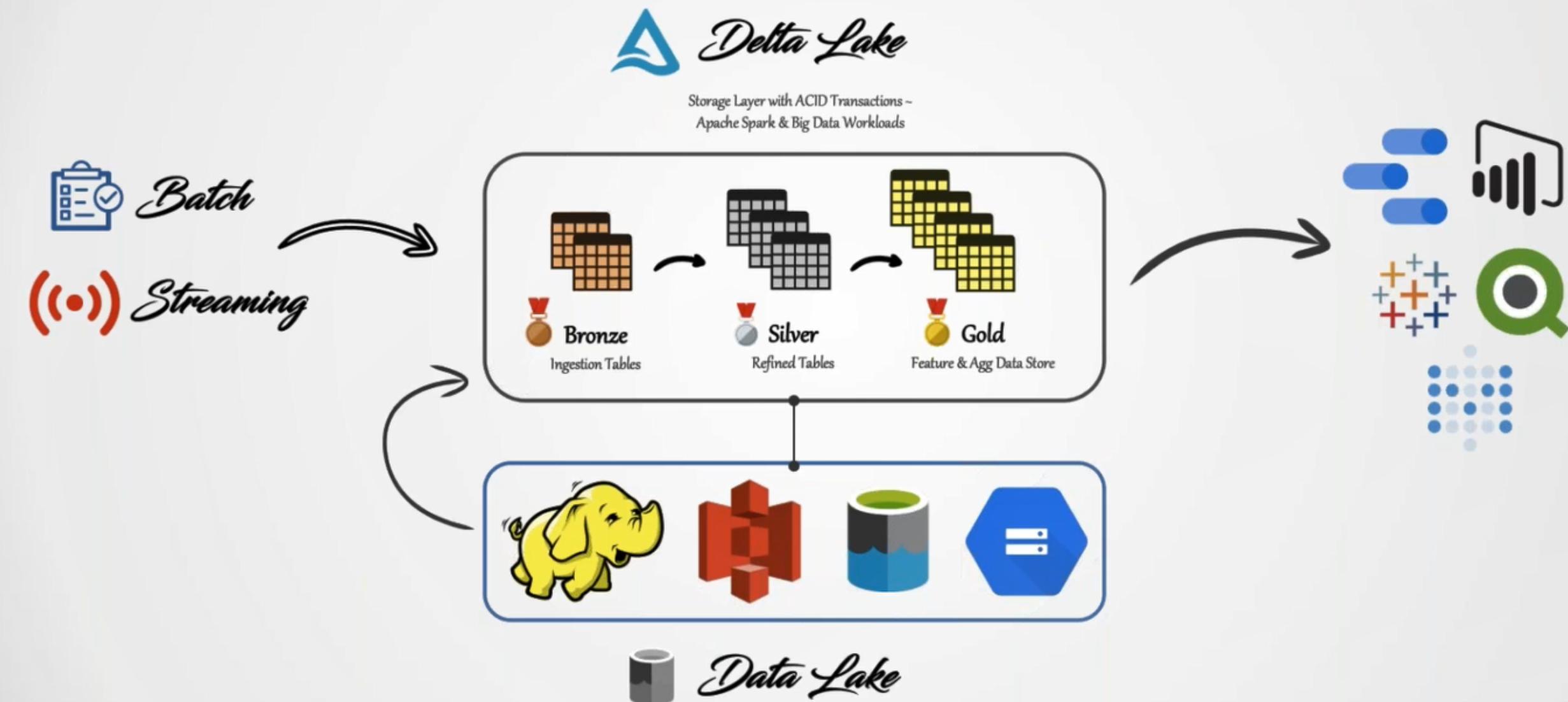
delta lake supports Scala, Java, Python, SQL APIs to merge, update and delete datasets. this allows you to easily comply with GDPR and CCPA and simplifies use cases like change data capture.



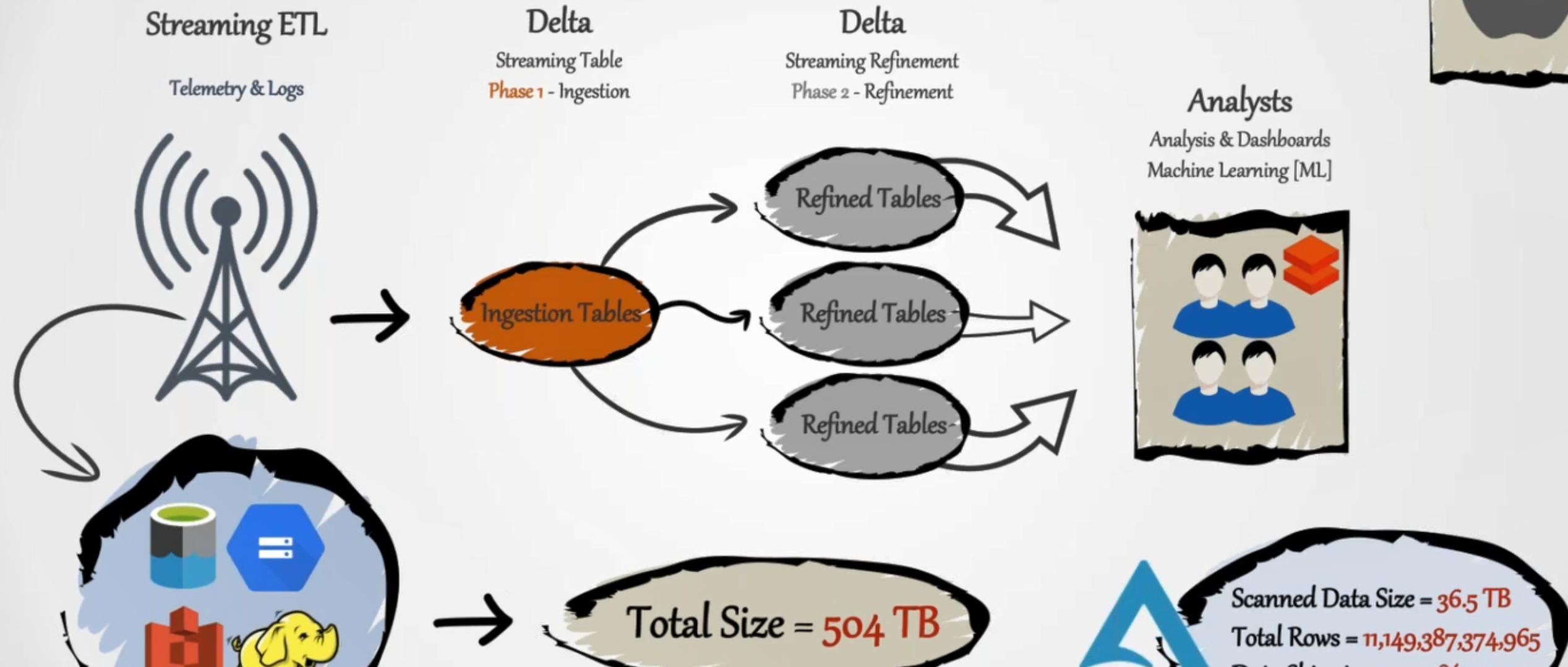
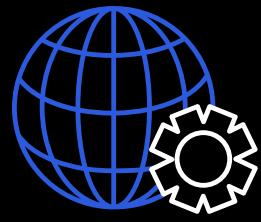
Data Lake vs. Delta Lake

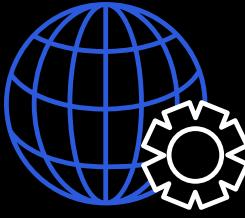


The Delta Architecture



Delta Lake [Apple's Use-Case]





Data Warehouse [Dw]

Decision Support for BI Applications
Since 1980s evolving & MPP Architecture
Structured Data & No Cost Efficient

Data Lake [DL]

Large Amount of Data [DV] from Different Sources
Since 2000s being a Data Raw Repository
Without Transactions, Data Enforcement, Lack of Consistency
Batch & Streaming Jobs not Performing as Expected
AI Advances based on Unstructured Data (Text, Images, Video, Audio)

Lakehouse [LH]

Similar Structure of Dw with a Low-Cost Storage with Features:

- Transaction Support
- Schema Enforcement and Governance
- BI Support
- Decoupled Storage & Computation
- Openness with Structured & Unstructured Data
- End-to-End Streaming

Data Lakehouse

Best of a Data Warehouse and Data Lake

OLTP System
RDBMS

ETL
Extract | Transform | Load

Data Mart [DM]
Departmental Business Decision

Data Warehouse [Dw]
Ralph Kimball [Bottom-Up]



Data Sources
Applications, RDBMS & NoSQL

ELT
Extract | Load | Transform

Data Lake
Structured, Semi-Structured & Unstructured



Data Sources
Structured, Semi-Structured & Unstructured Data

ELT
Extract | Load | Transform

Lakehouse
BI, ML, Streaming Analytics & Data Science



Data Ingestion into Data Lakehouse Best Practices



Network of Partners [Connectors]

essential ecosystem of connectors to bring data into Delta Lake
Azure Data Factory, Fivetran, Qlik, Infoworks, StreamSets, Syncsoft



Cloud Storage [Auto Loader]

loading data continuously from cloud stores with Exactly-Once guarantees at low cost, low latency and minimal DevOps work. Auto Loader is one of the most expensive operations. Auto Loader is an optimized file source that uses Structured Streaming



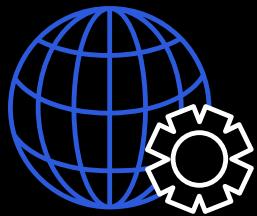
Streaming Load [Structured Streaming]

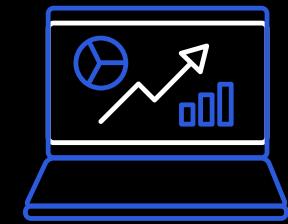
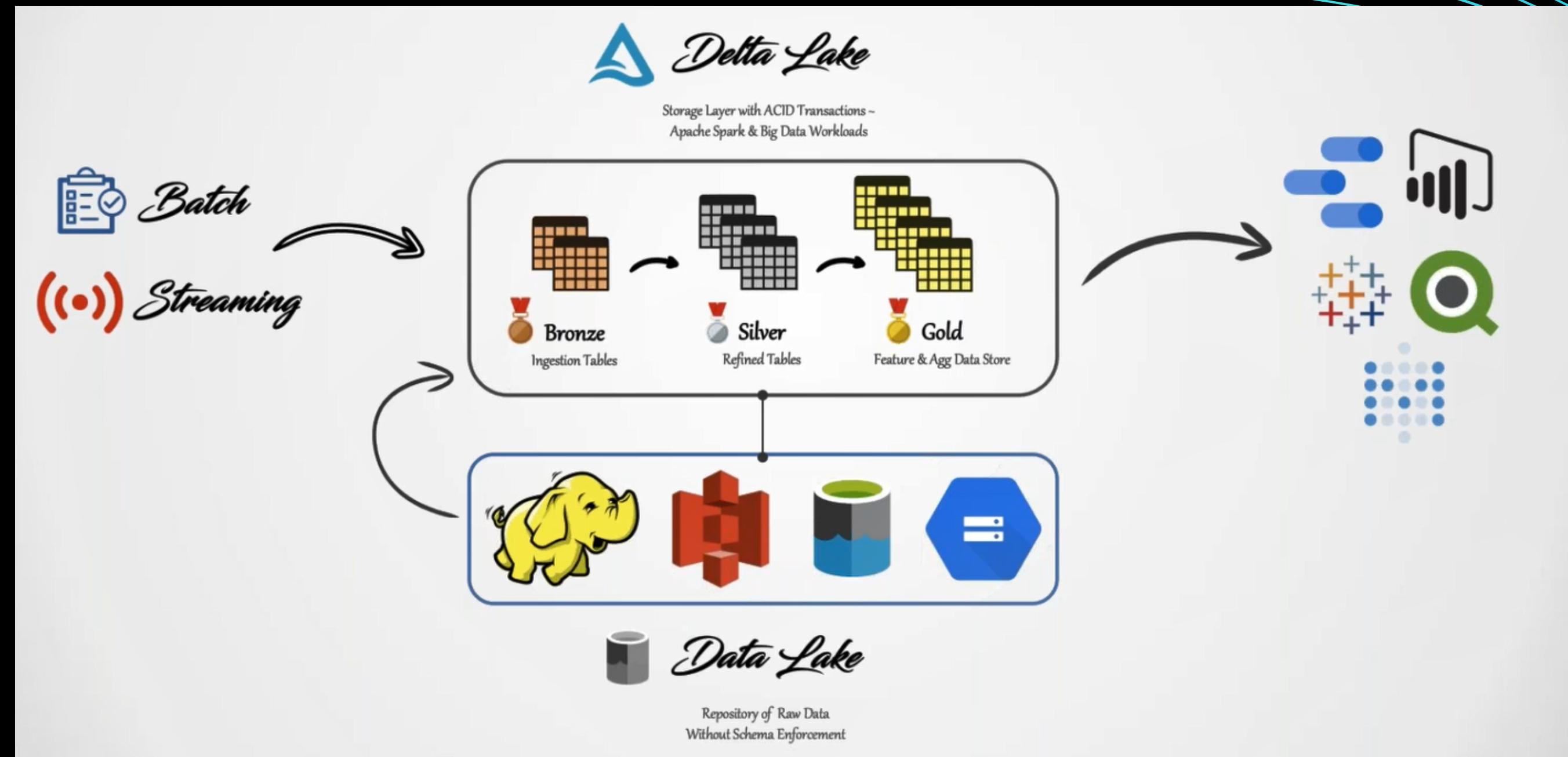
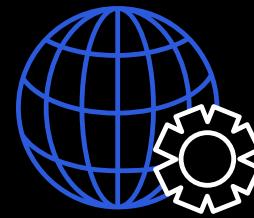
loading data continuously from Apache Kafka with Exactly-Once using the Structured Streaming API

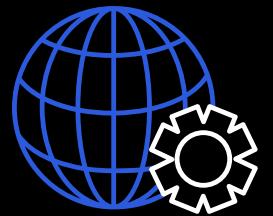


Delta Lake

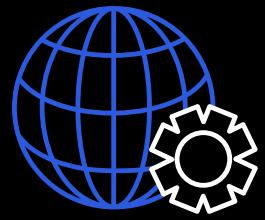
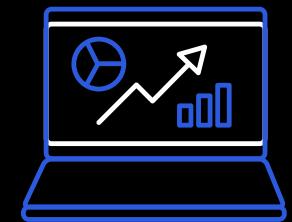
storage system with ACID capabilities that is designated to build a Data Lakehouse







Best Practices Organizing a Data Lake Repository

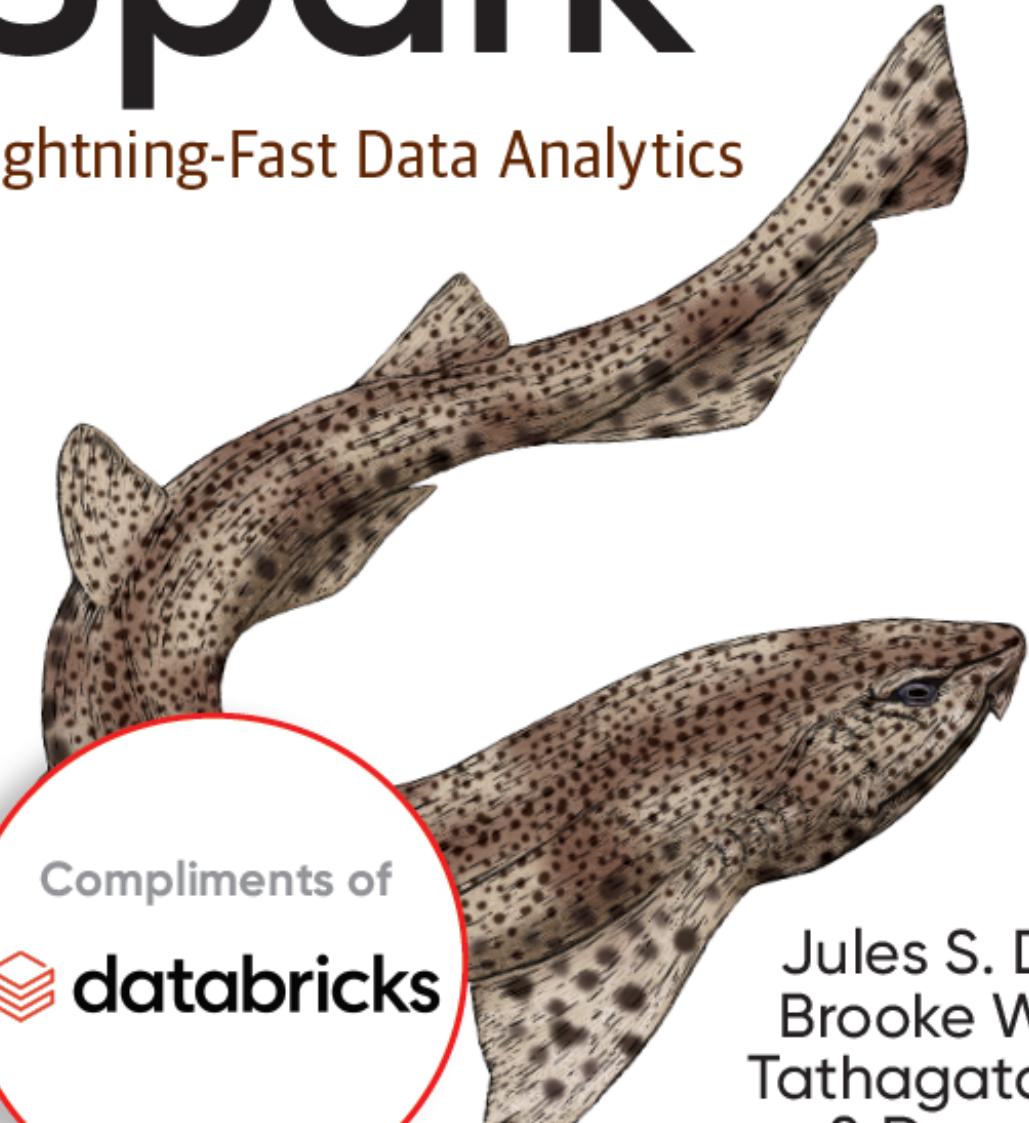


O'REILLY®

Learning Spark

Lightning-Fast Data Analytics

2nd Edition
Covers
Apache Spark 3.0



Compliments of



Jules S. Damji,
Brooke Wenig,
Tathagata Das
& Denny Lee

Foreword by Matei Zaharia