

Unified Platform for Big Data & AI



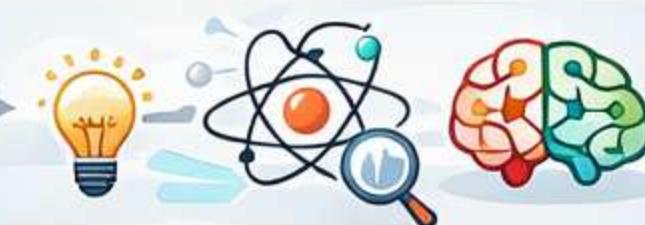
# Databricks is Needed?

## Handles Big Data Efficiently



Apache Spark processes large-scale data quickly

## All-in-One Platform



Combines Data Engineering, Analytics, Machine Learning

## Faster Analytics & Insights



In-memory processing makes everything faster

## Databricks

Built on Apache Spark®

## Scales Automatically



Cloud-based & scales up/down automatically

## Scales Automatically



Cloud-based & scales up/down automatically

## Easy Collaboration



Using Notebooks & Shared Data

## Supports Advanced Analytics & AI

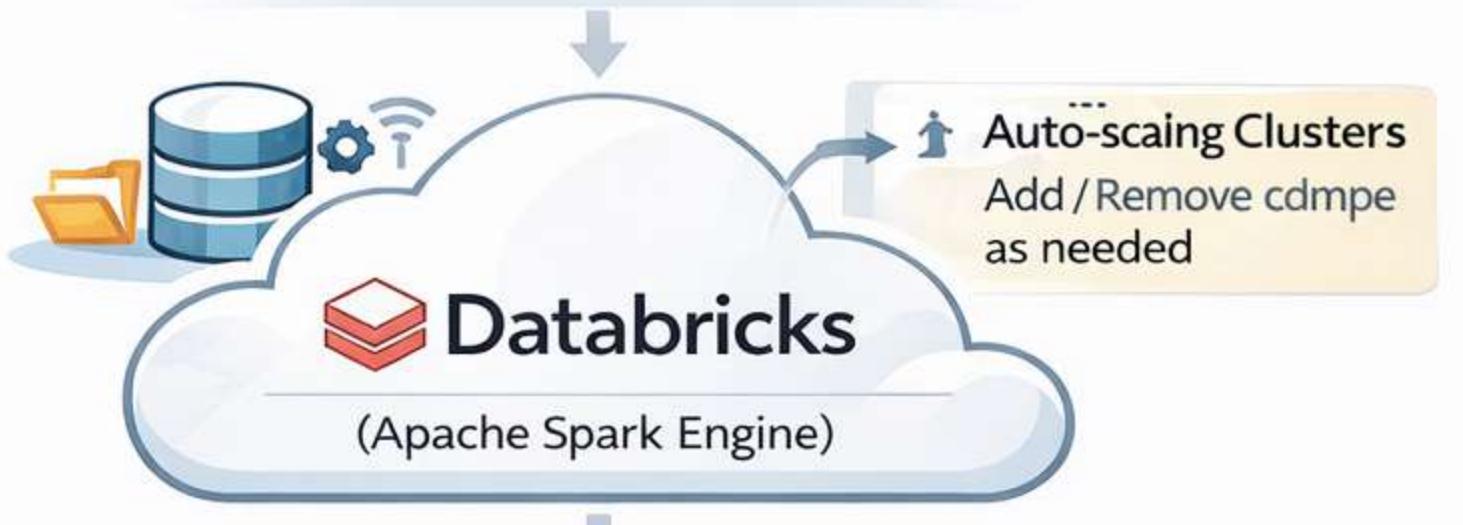


For Real-time Analytics & Machine Learning

Unified Data Platform for Big Data & AI

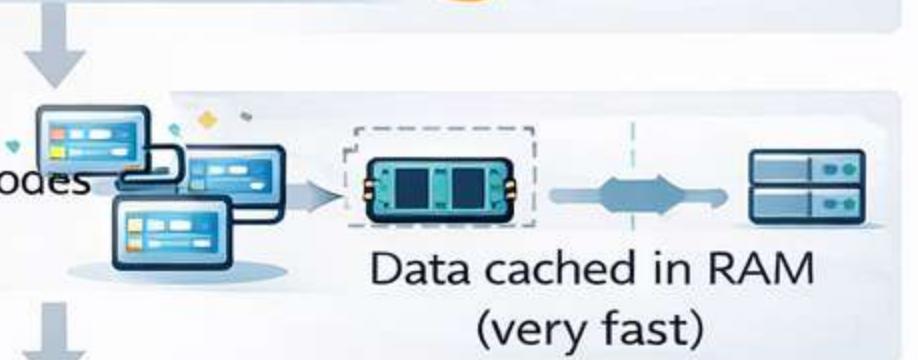
# How Databricks Handles Scale

## Data Sources (Files | Databases | Streams)



### Distributed Processing

- Split data into smaller chunks
- Process chunks in parallel across nodes
- More machines = faster processing



### Cloud-Based Compute

aws | Azure | GCP

### In-Memory Processing

Data cached in RAM (very fast)

### Distributed Processing

- Split data into smaller chunks
- Process chunks in parallel across nodes
- More machines = faster processing



### Storage ≠ Compute

Separation of compute clusters and storage (Data Lake)

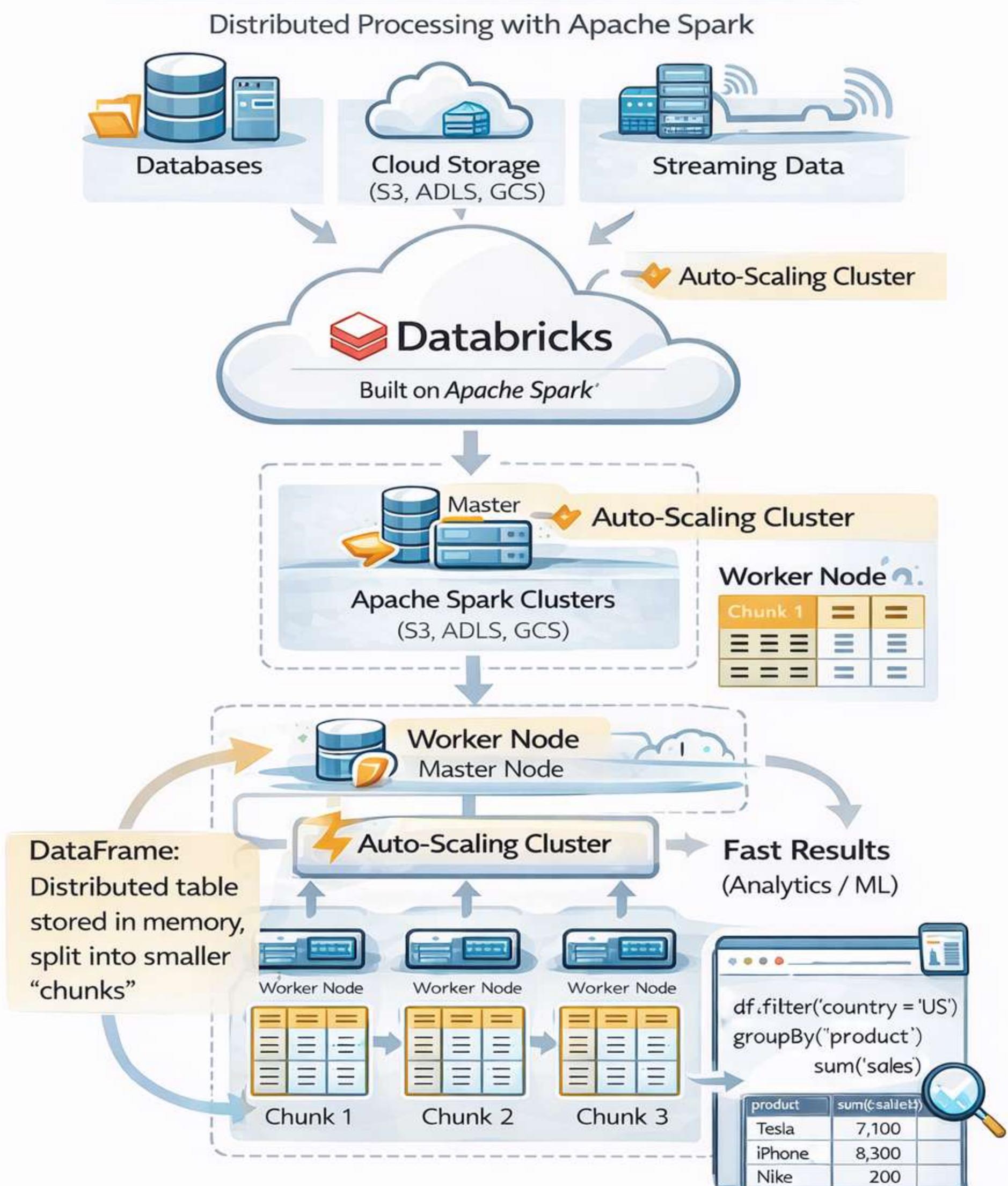
Auto Scaling Clusters

Scale compute up/down

Cloud-Based Compute

Runs on AWS, Azure, GCP

# Clusters & DataFrames in Databricks

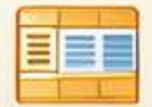


## Distributed Processing



Breaks big data into “chunks” and processes them in parallel across a cluster

## DataFrames in RAM



Keeps chunks in-memory for high-speed queries

## Auto-Scaling Clusters



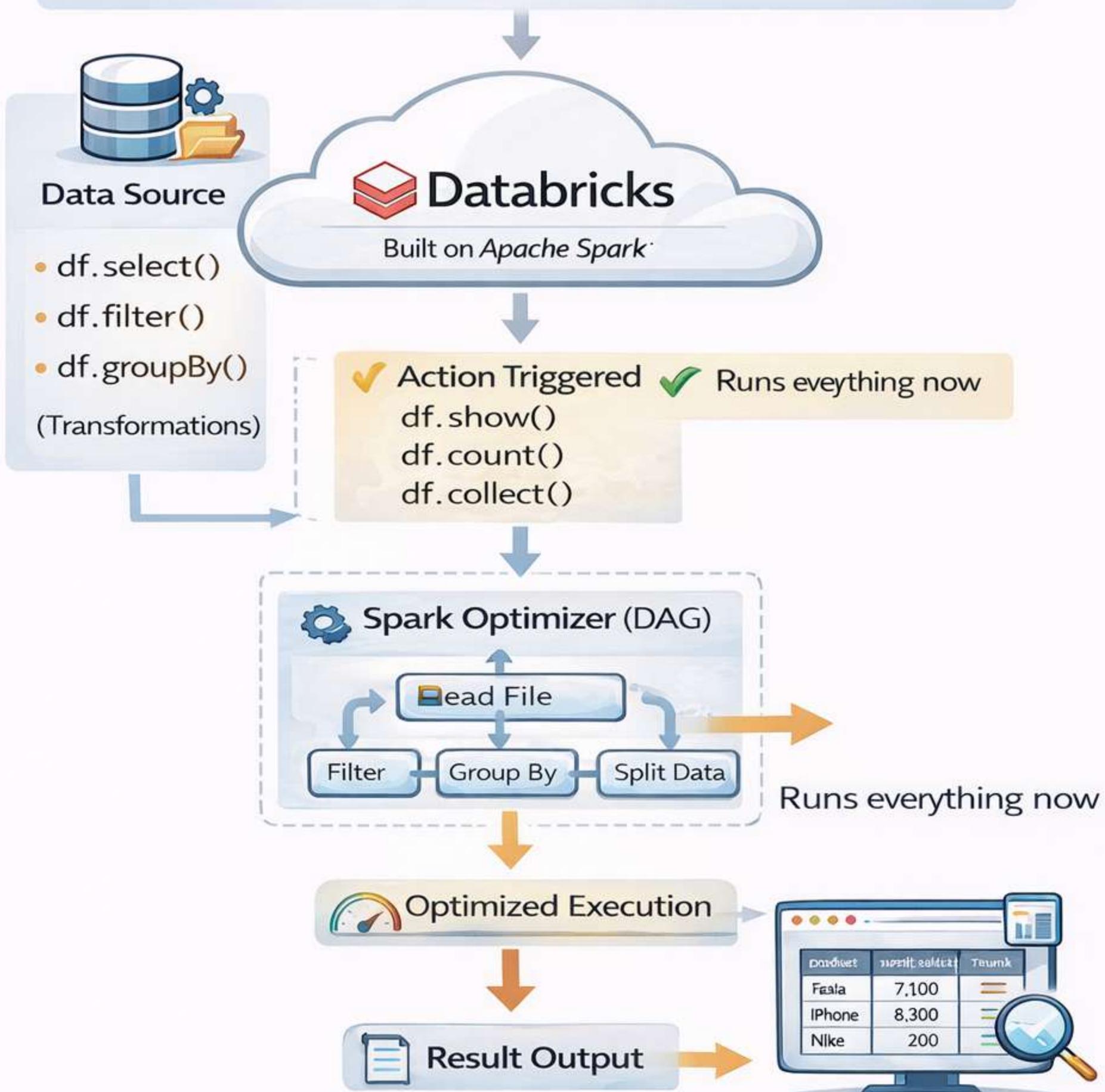
Adds or removes worker nodes based on workload

# Lazy Execution in Databricks

(Apache Spark)



Spark builds an **execution plan first and runs it only when a result is required.**



## Transformations (Lazy)

- `select()`
- `filter()`
- `groupBy()`
- `withColumn()`

Spark records steps

## Actions (Trigger Execution)

- `show()`
- `count()`
- `write()`

Spark runs all steps now



# Databricks SQL & Tables

## Use SQL Queries

```
SELECT product, SUM(sales)
FROM sales
JOIN products ON
sales.product_id = products.id
WHERE country = 'US'
GROUP BY product
ORDER BY SUM(sales) DESC;
```

Interactive Queries



Interactive Queries

## Analyze & Transform Data in Tables

product	sales	Growth
Tesla	5,200	20%
iPhone	6,800	15%
Nike	320	5%

UPDATE    INSERT

DELETE



Databases

## Build Fast Analytics on a Data Lake



Databricks

Data Lakehouse

Databases

sales

products

users

web\_logs



Cloud Storage

aws | Azure | GCP

Databricks combines fast analytics with data lake storage, enabling powerful data analysis and transformation with SQL.



# AI & Why This Matters?



## 1 Better Decisions

AI analyzes data to find patterns and insights

## 2 Faster Results

Queries and predictions much quicker

**Databricks**  
Unified Data & AI Platform

## 4 Cost Savings

Automating tasks reduces costs & errors

## 3 Innovation & Growth

AI drives new products and better services



## 4 Cost Savings

Automating tasks reduces costs & errors



**Databricks**



Innovation & Growth  
AI drives new products and better services

Databricks combines AI accessible by simplifying large-scale analytics, machine learning, and data engineering with SQL.

## Why Databricks?



Easily train and deploy AI models on your data