

# Big Data & AI in Business

How to become a data driven business

Sessions 10&11: Data Strategy & AI Operations

David G Pisano

The **AI strategy** aims to obtain the maximum potential of this technology for our businesses by **solving all the barriers** to its implementation.

1	Business Value	<p><b>Usage:</b> What are the main drivers of AI value in each business? <b>AI ambition:</b> How much do we believe AI can transform our business? <b>AI maturity:</b> How innovative in the use of AI does the company need to be?</p>	
2	Data Access	<p><b>Types of sources:</b> What data will be needed at the company to develop AI? <b>Data Platform:</b> What technological tools will we use to facilitate access to data at the company? <b>Data Governance:</b> How are we going to ensure that they are of the right quality?</p>	
3	AI Operations	<p><b>Organization:</b> Where should we locate the AI specialists at the company? <b>Sourcing Model:</b> Who develops and maintains the AI models? <b>Process industrialization:</b> How efficient and scalable is AI development, deployment and maintenance?</p>	
4	AI Culture	<p><b>Skills training &amp; hiring:</b> How can we train the company staff? How do we attract &amp; retain talent? <b>Communication:</b> how do we make it easier for everyone to understand what AI brings to the company? <b>Government:</b> How do we guarantee the proper use of AI (security, privacy, ethics) at the company?</p>	
5	Are the legal/regulatory issues sufficiently resolved?	6 Is society ready to accept the transformation of the value proposition?	7 It may be legal, profitable and working - but is it ethical?

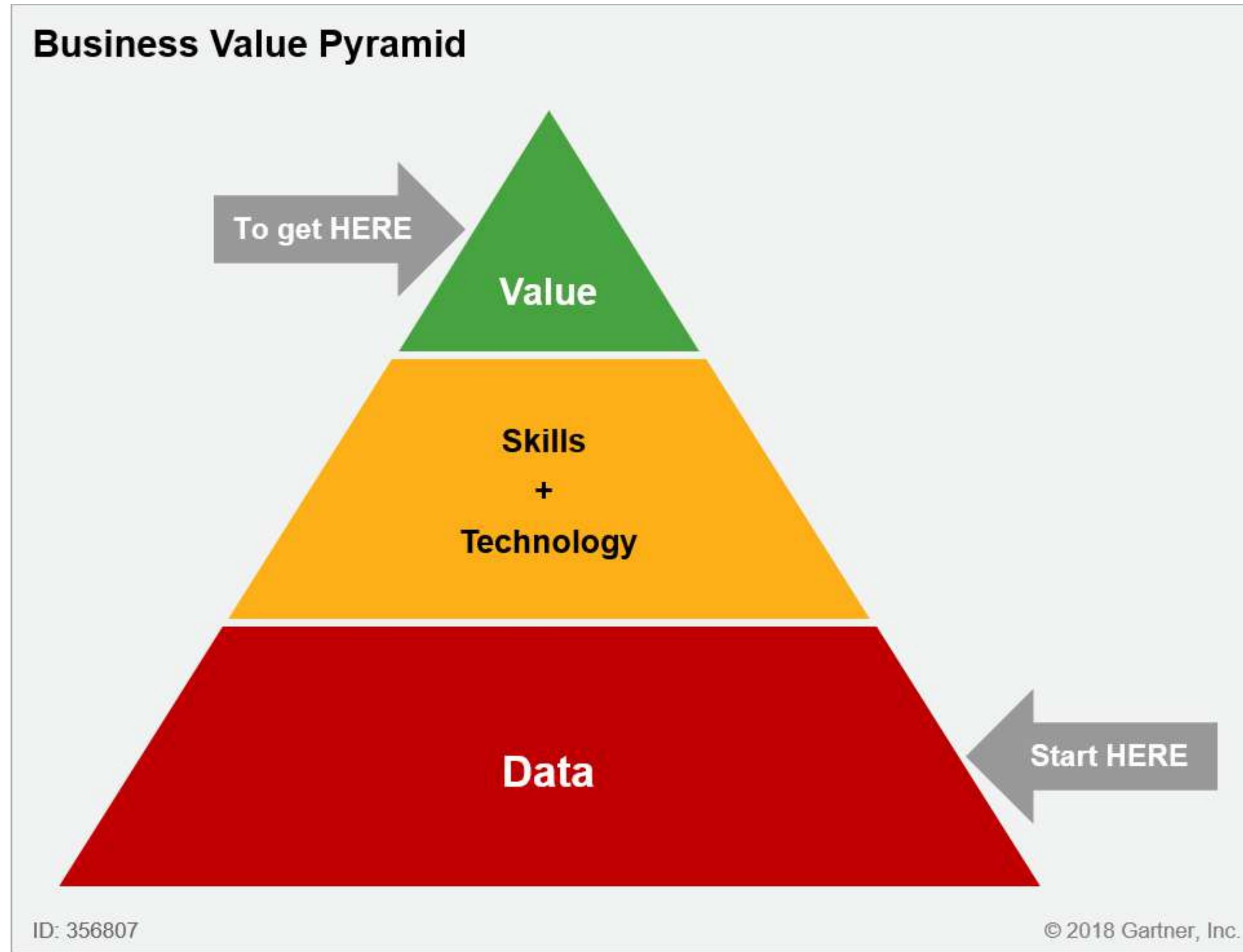
The **AI strategy** aims to obtain the maximum potential of this technology for our businesses by **solving all the barriers** to its implementation.

1	Business Value	<p><b>Usage:</b> What are the main drivers of AI value in each business? <b>AI ambition:</b> How much do we believe AI can transform our business? <b>AI maturity:</b> How innovative in the use of AI does the company need to be?</p>	
2	Data Access	<p><b>Types of sources:</b> What data will be needed at the company to develop AI? <b>Data Platform:</b> What technological tools will we use to facilitate access to data at the company? <b>Data Governance:</b> How are we going to ensure that they are of the right quality?</p>	
3	AI Operations	<p><b>Organization:</b> Where should we locate the AI specialists at the company? <b>Sourcing Model:</b> Who develops and maintains the AI models? <b>Process industrialization:</b> How efficient and scalable is AI development, deployment and maintenance?</p>	
4	AI Culture	<p><b>Skills training &amp; hiring:</b> How can we train the company staff? How do we attract &amp; retain talent? <b>Communication:</b> how do we make it easier for everyone to understand what AI brings to the company? <b>Government:</b> How do we guarantee the proper use of AI (security, privacy, ethics) at the company?</p>	
5	Are the legal/regulatory issues sufficiently resolved?	6 Is society ready to accept the transformation of the value proposition?	7 It may be legal, profitable and working - but is it ethical?

# 2nd Barrier: Data



# AI is the **evolution** of data through the various stages of analysis



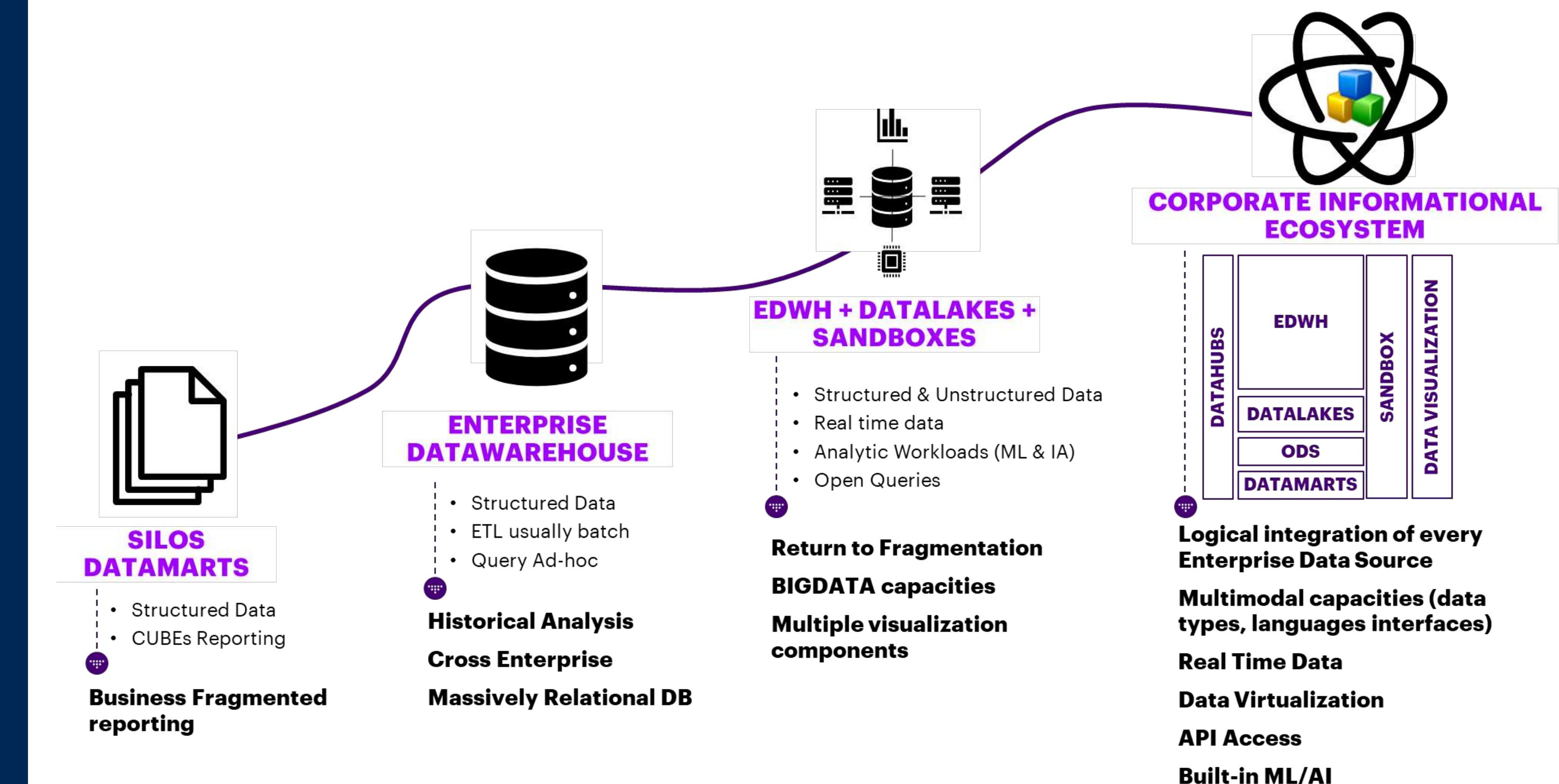
- Everyone talks about data, but no one emphasises the way **data needs to be managed systematically**, which allows **data to be used** by machine learning algorithms to **build intelligent systems** that **automate business processes**
- Think of **AI** as an **extension** of your **existing analytics** platform

# Data Access

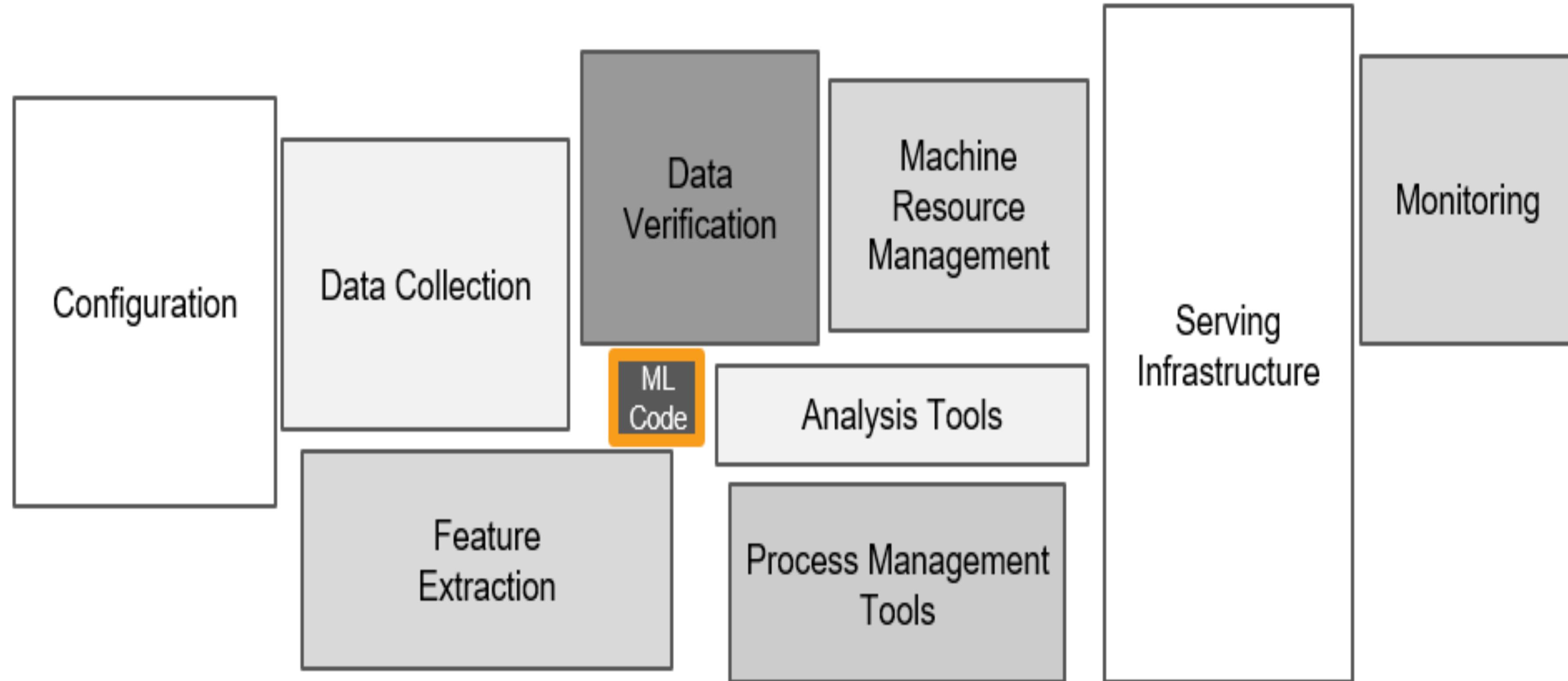


*For organizations that are more mature in the use of AI, data availability is the main problem to solve*

O'Reilly 2020



# Deployment of **AI/ML** is not (only) about the algorithms



A lot of the development work associated with AI/ML is not related to the AI/ML algorithms themselves. Rather, it is more concerned with the **preparation of data**, its **governance**, **performance** considerations and **delivery** of results according to service-level agreements (SLAs).

# Data Platform

¿What is the data management platform (**DMP**) used for collecting and managing data? ¿It allows to use big data and artificial intelligence algorithms to process and analyze large data sets? ¿Can it grow and acquire new functionalities as the data grows and the business needs evolve?

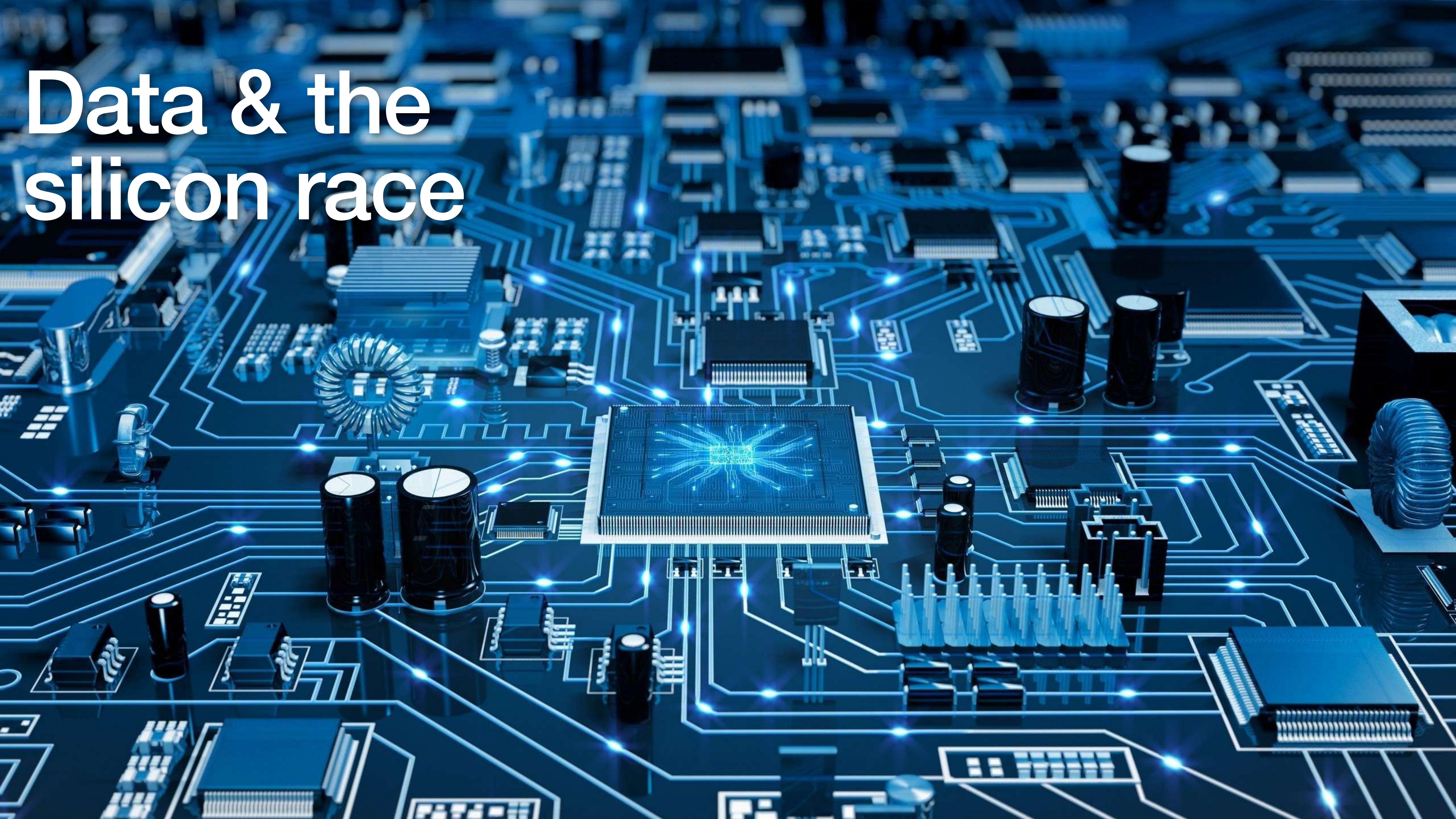


## Data Governance

¿What are the processes and policies in place that ensure that high **data quality** exists through its complete lifecycle? ¿What data **controls** there are that support business objectives? ¿It is the data **available**, **usable**, **consistent** and **secure**?



# Data & the silicon race



# Gordon E Moore

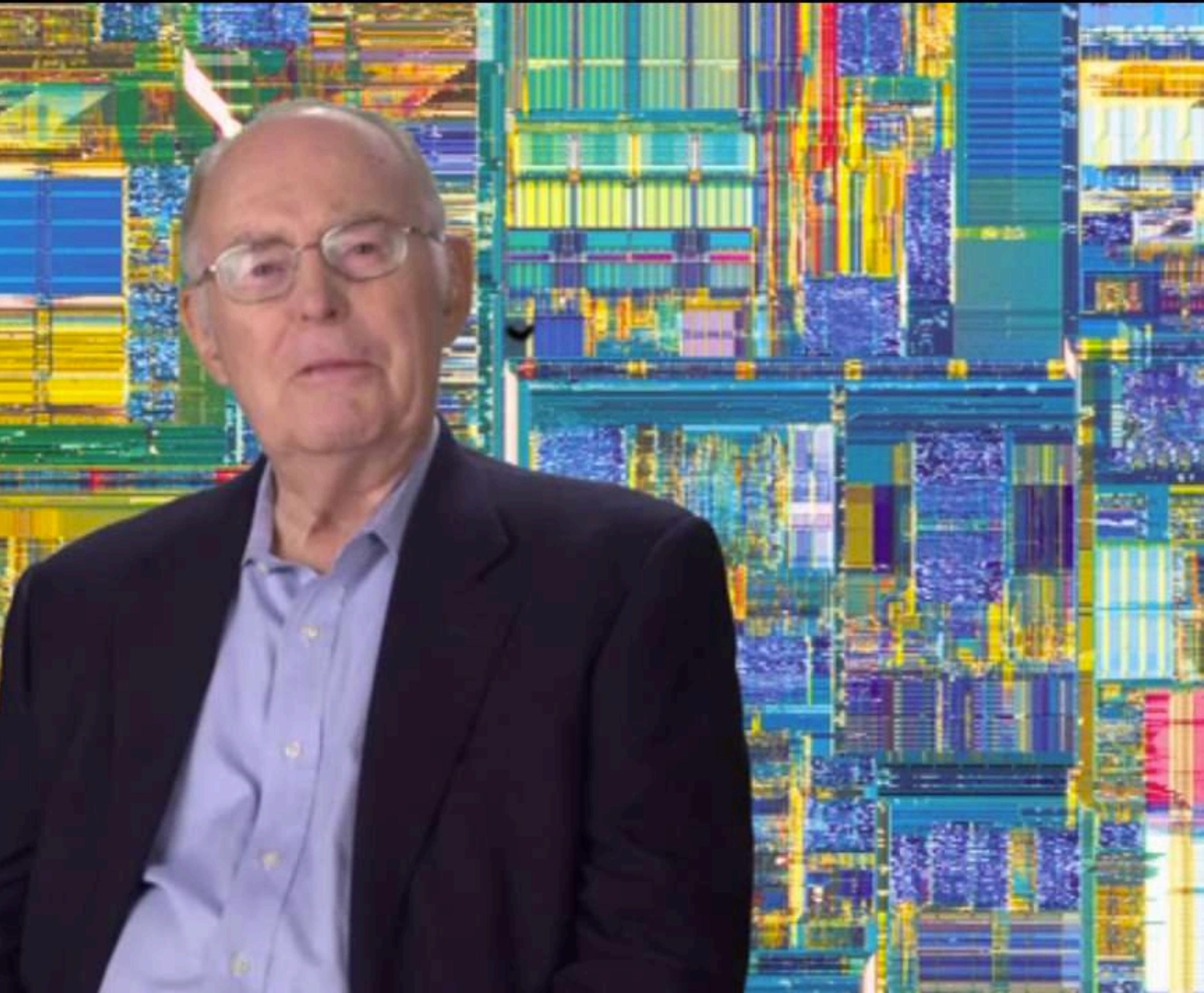
## Intel Co-founder

In an article published in 1965, Moore observed that the number of components (transistors, resistors, diodes, or capacitors) in a dense integrated circuit had **doubled approximately every year** and speculated that it would continue to do so for at least the next ten years.

In 1975, he revised the forecast rate to approximately **every two years**.

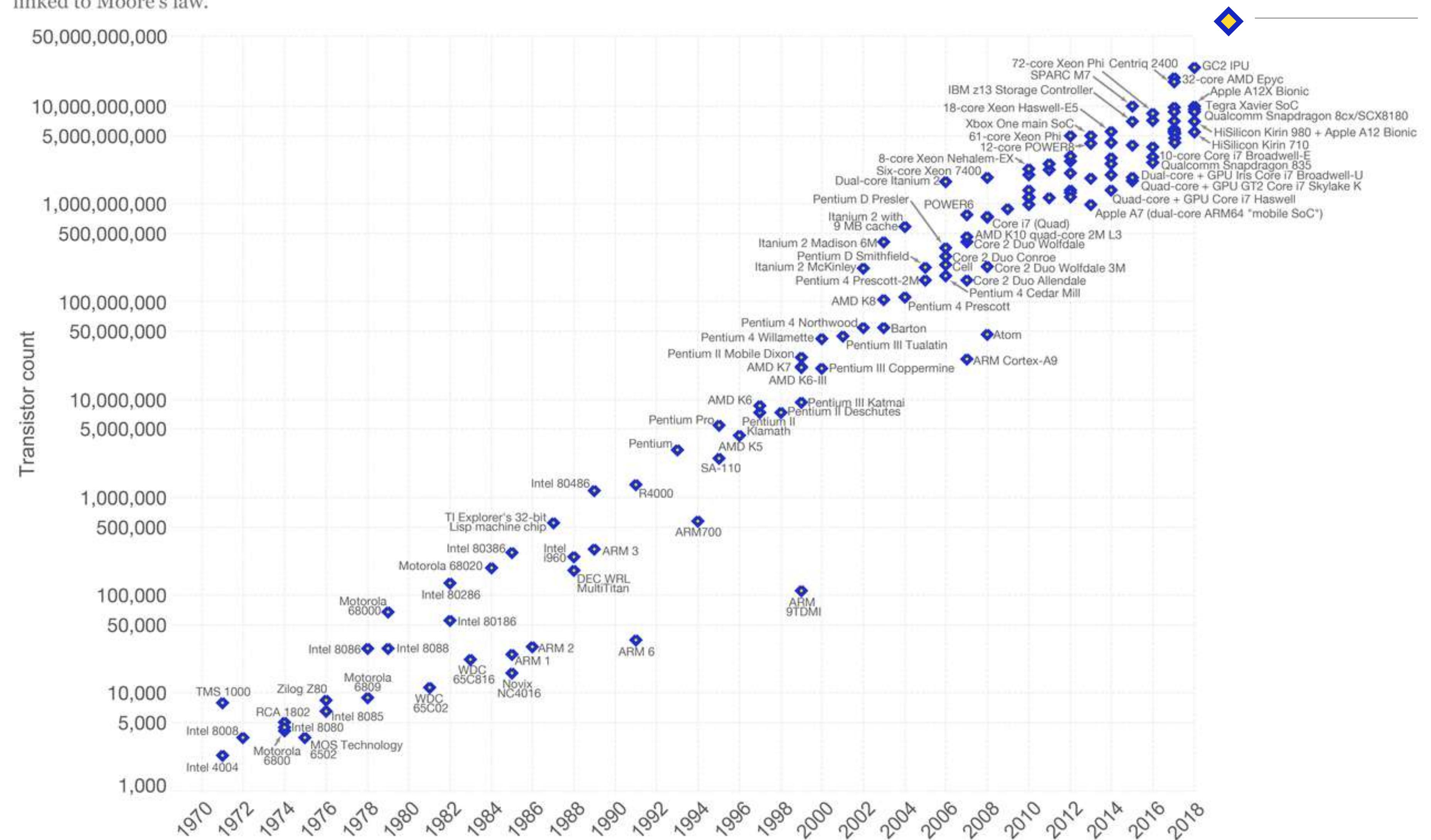
The phrase **Moore's law** was popularised.

The prediction has become a **target for miniaturization** in the semiconductor industry and has had widespread **impact in many areas of technological change**.



# Moore's Law – The number of transistors on integrated circuit chips (1971-2018)

Moore's law describes the empirical regularity that the number of transistors on integrated circuits doubles approximately every two years. This advancement is important as other aspects of technological progress – such as processing speed or the price of electronic products – are linked to Moore's law.

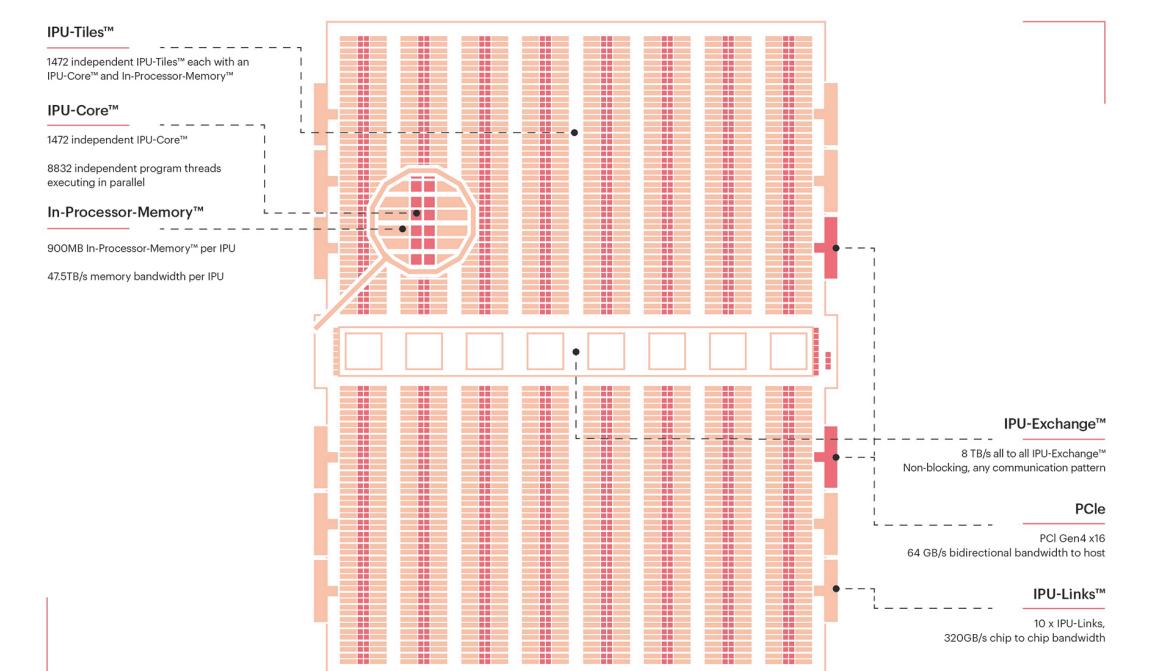


Data source: Wikipedia ([https://en.wikipedia.org/wiki/Transistor\\_count](https://en.wikipedia.org/wiki/Transistor_count))

The data visualization is available at OurWorldinData.org. There you find more visualizations and research on this topic.

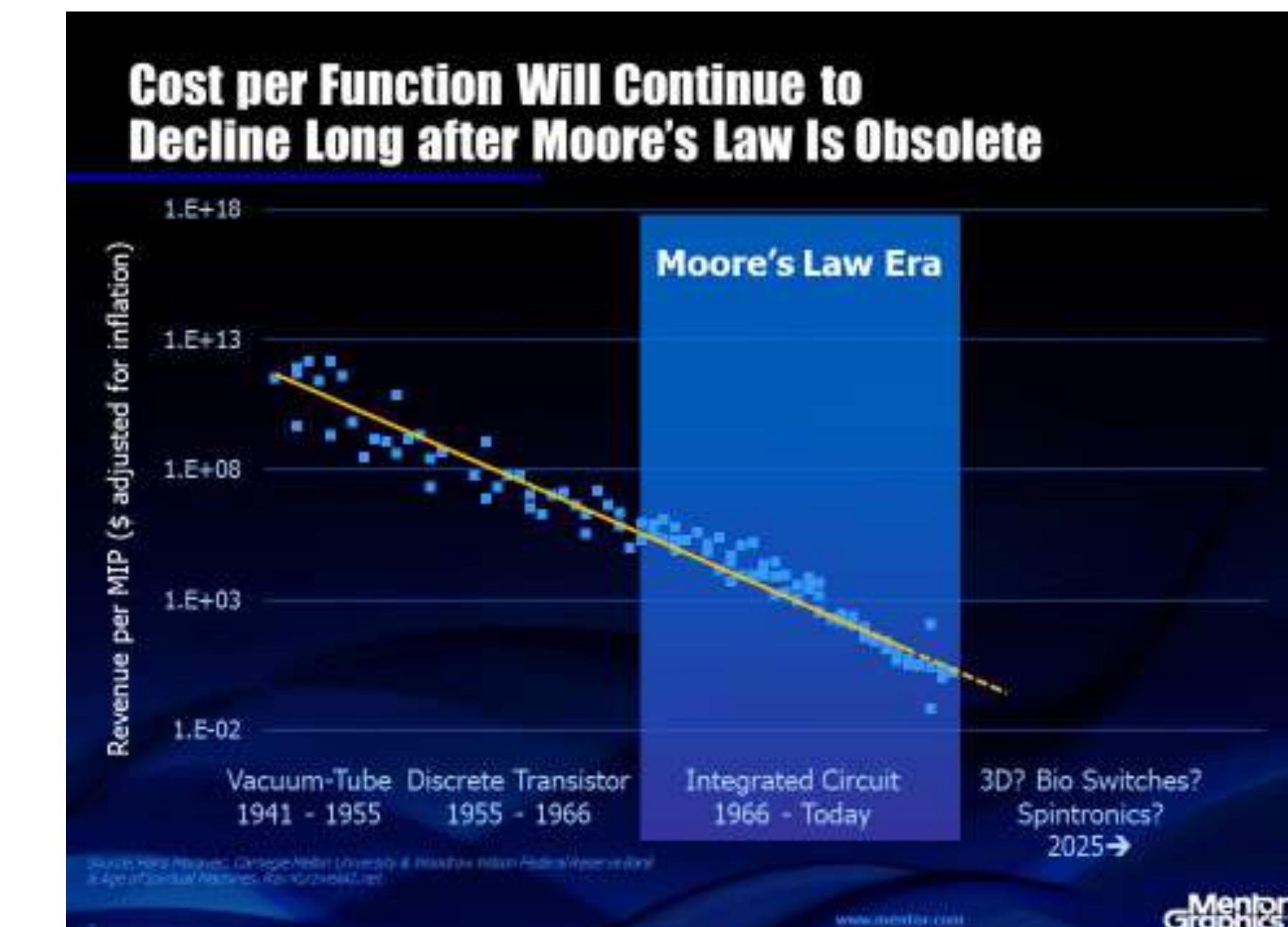
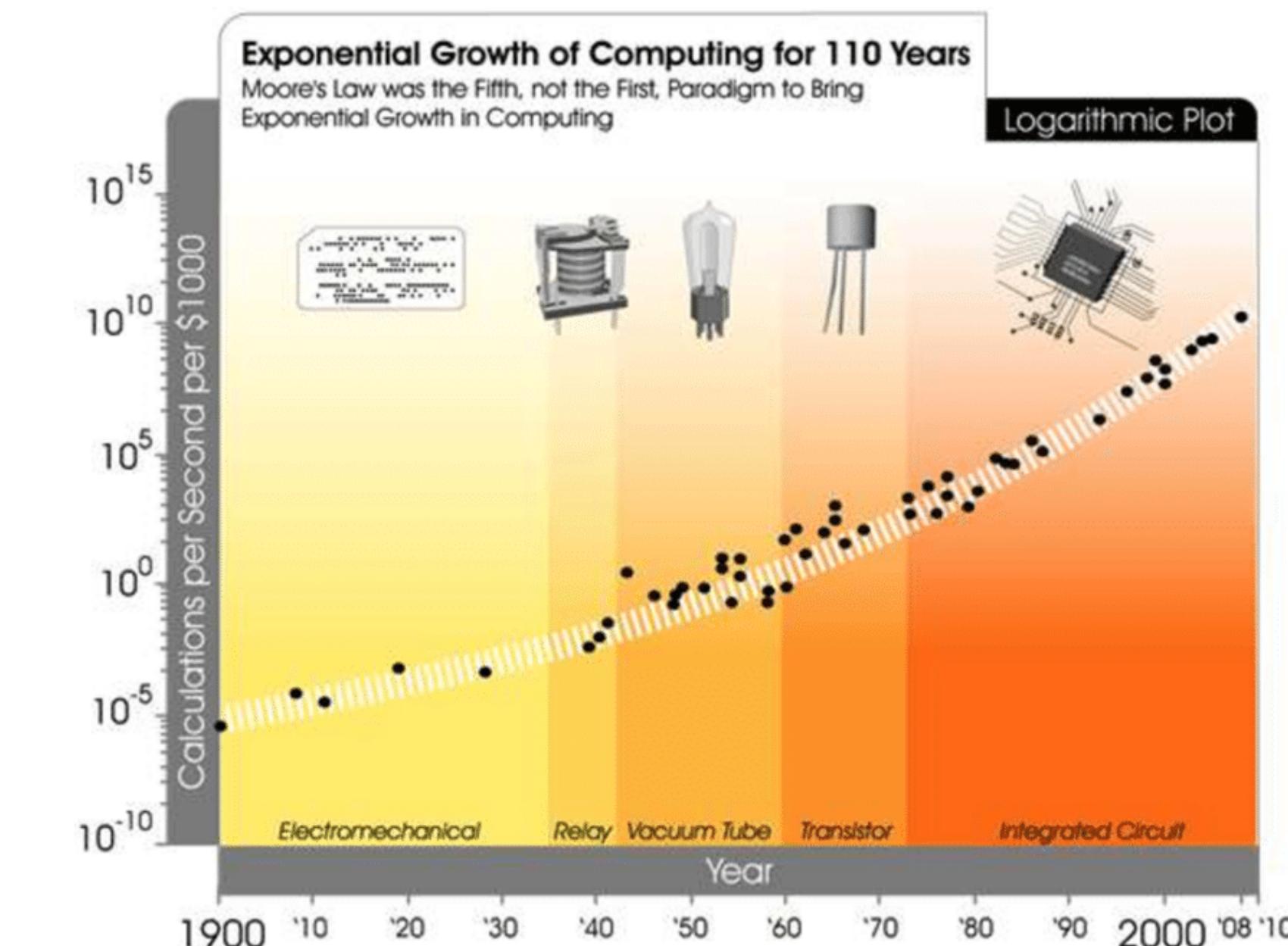
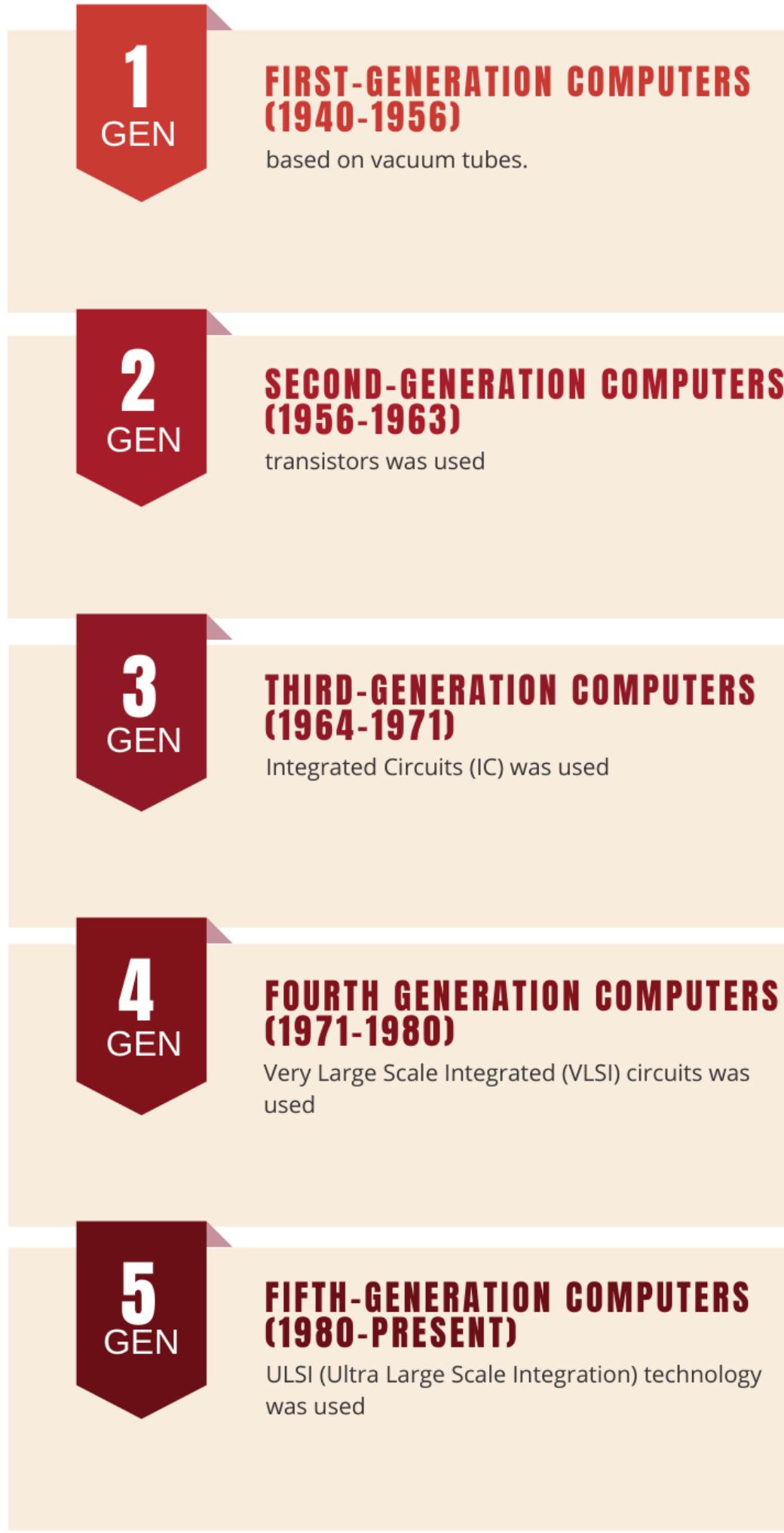
Licensed under CC-BY-SA by the author Max Roser.

Colossus™ MK2 GC200 IPU (2020)  
59.4Bn transistors  
1,472 cores  
250 TFlops per chip



# The **fifth** paradigm change: what will be the next one?

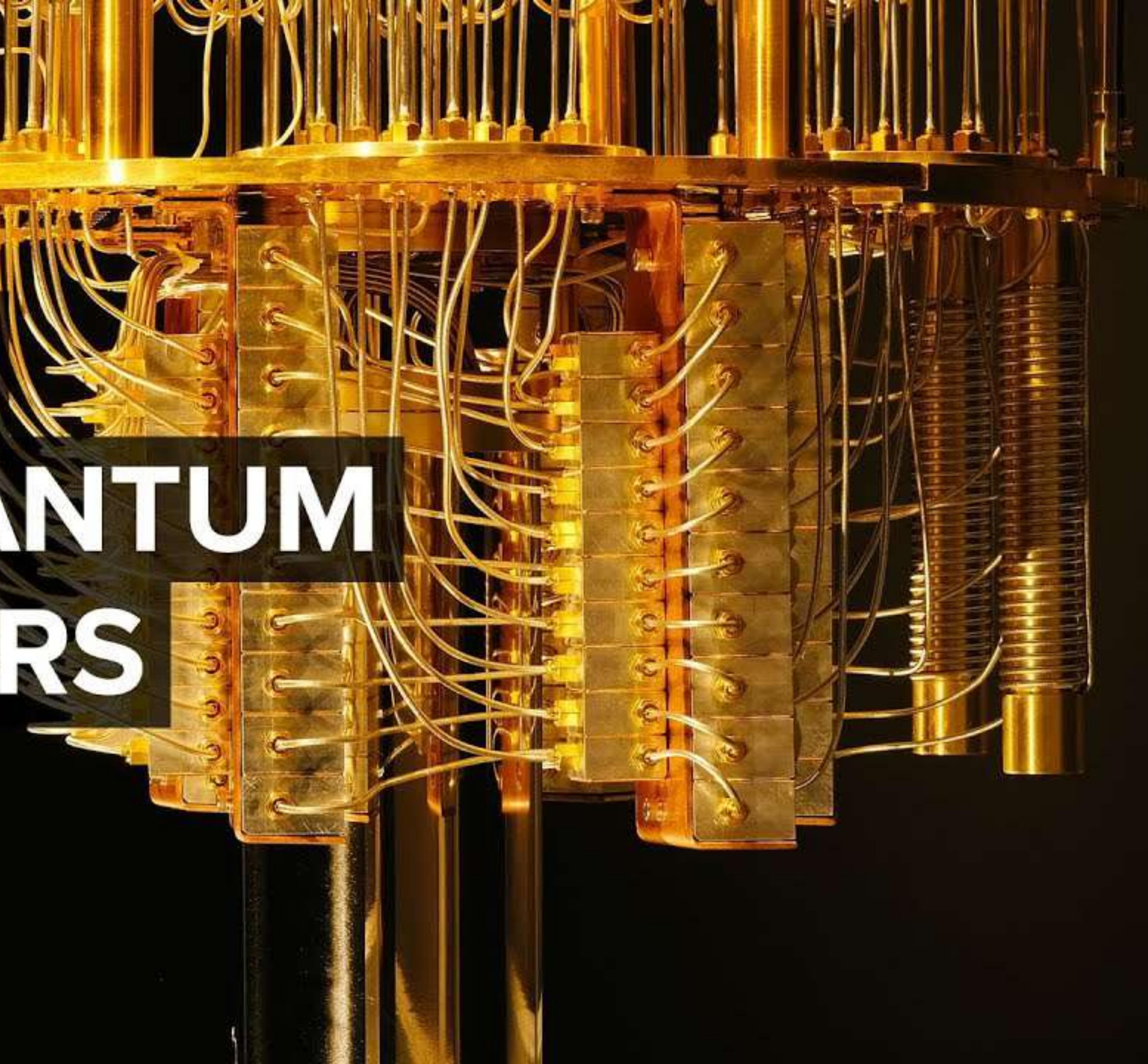
## GENERATION OF COMPUTERS



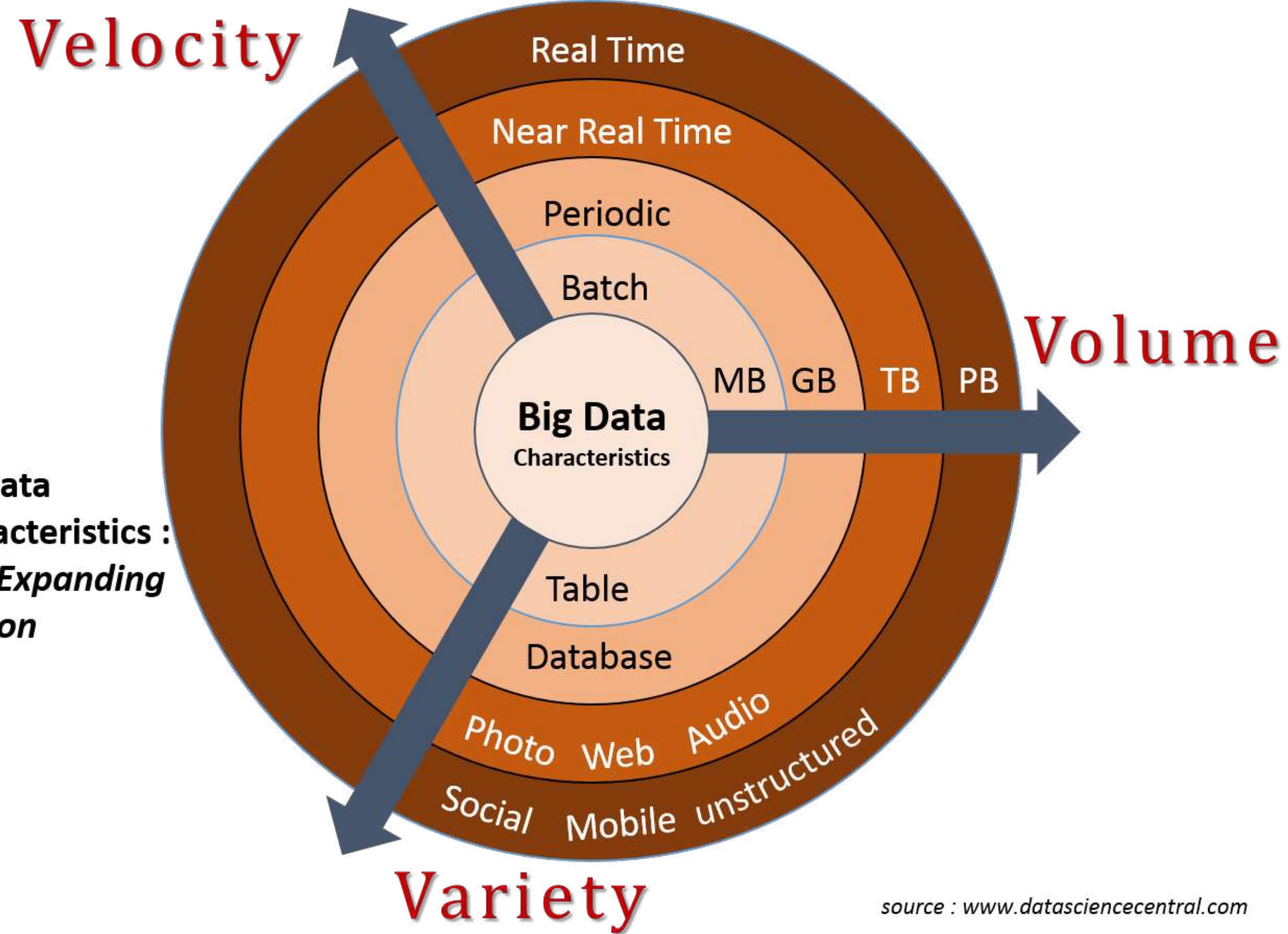


# THE HYPE OVER QUANTUM COMPUTERS

In October 2019, Google made a big announcement. It announced its 53-qubit quantum computer named Sycamore had achieved '**quantum supremacy**.' That's when quantum computers **can complete tasks exponentially more quickly** than their classical counterparts. In this case, Google said its quantum machine completed a task **in 200 seconds** that would have taken the world's most powerful computer **10,000 years\*** to complete.

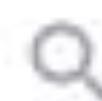


**Big Data  
Characteristics :**  
*Ever Expanding  
horizon*



source : [www.datasciencecentral.com](http://www.datasciencecentral.com)

# Google



big data



Google Search

I'm Feeling Lucky

# MapReduce: Big Data distribution

## MapReduce: Simplified Data Processing on Large Clusters

Jeffrey Dean, Sanjay Ghemawat

OSDI'04: Sixth Symposium on Operating System Design and Implementation, San Francisco, CA (2004), pp. 137-150

### Motivation: Large Scale Data Processing

Many tasks: Process lots of data to produce other data

Want to use hundreds or thousands of CPUs

- ... but this needs to be easy

MapReduce provides:

- Automatic parallelization and distribution
- Fault-tolerance
- I/O scheduling
- Status and monitoring

### Programming model

Input & Output: each a set of key/value pairs

Programmer specifies two functions:

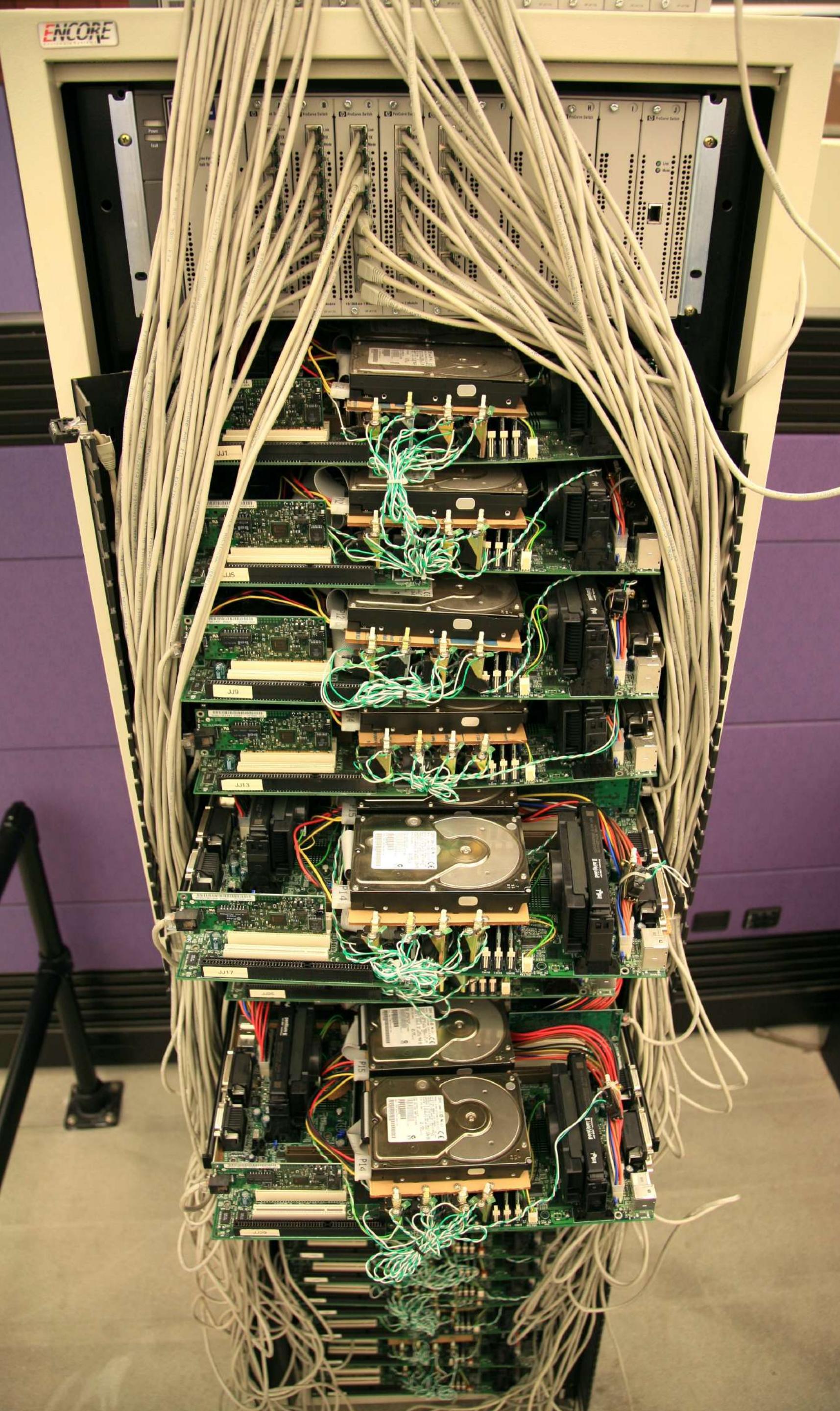
```
map (in_key, in_value) -> list(out_key, intermediate_value)
```

- Processes input key/value pair
- Produces set of intermediate pairs

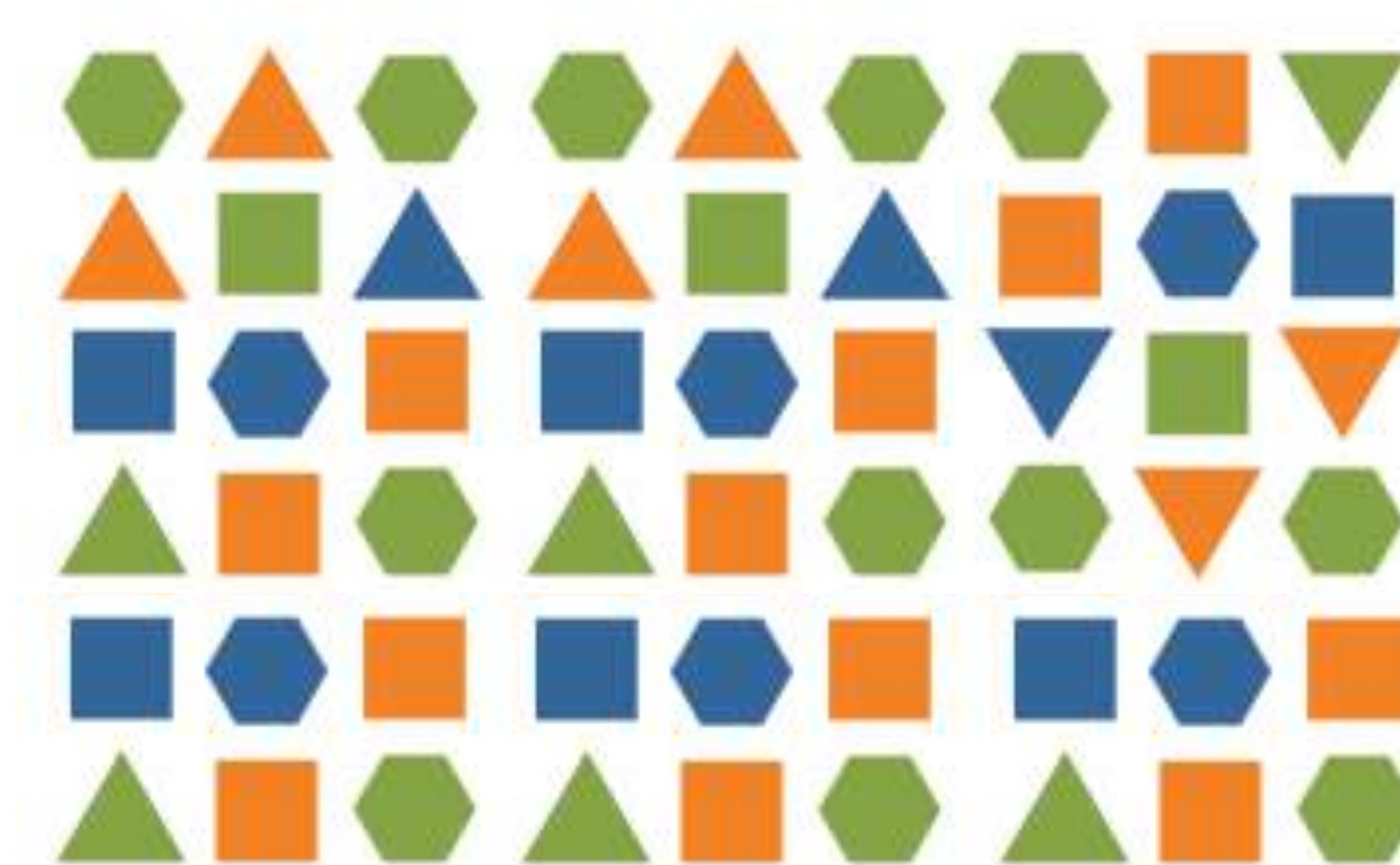
```
reduce (out_key, list(intermediate_value)) -> list(out_value)
```

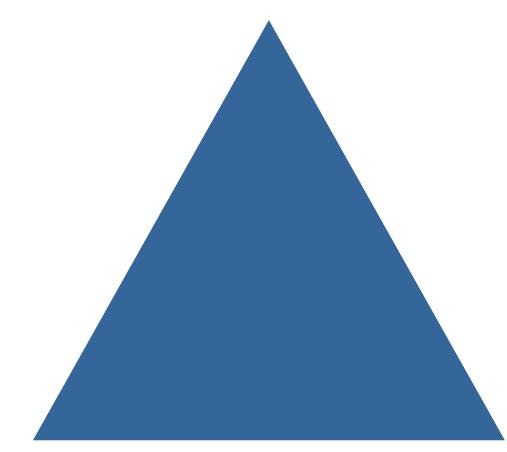
- Combines all intermediate values for a particular key
- Produces a set of merged output values (usually just one)

Inspired by similar primitives in LISP and other languages



# Map Reduce example

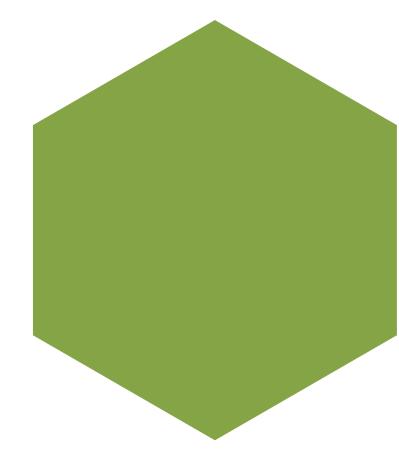




**15**



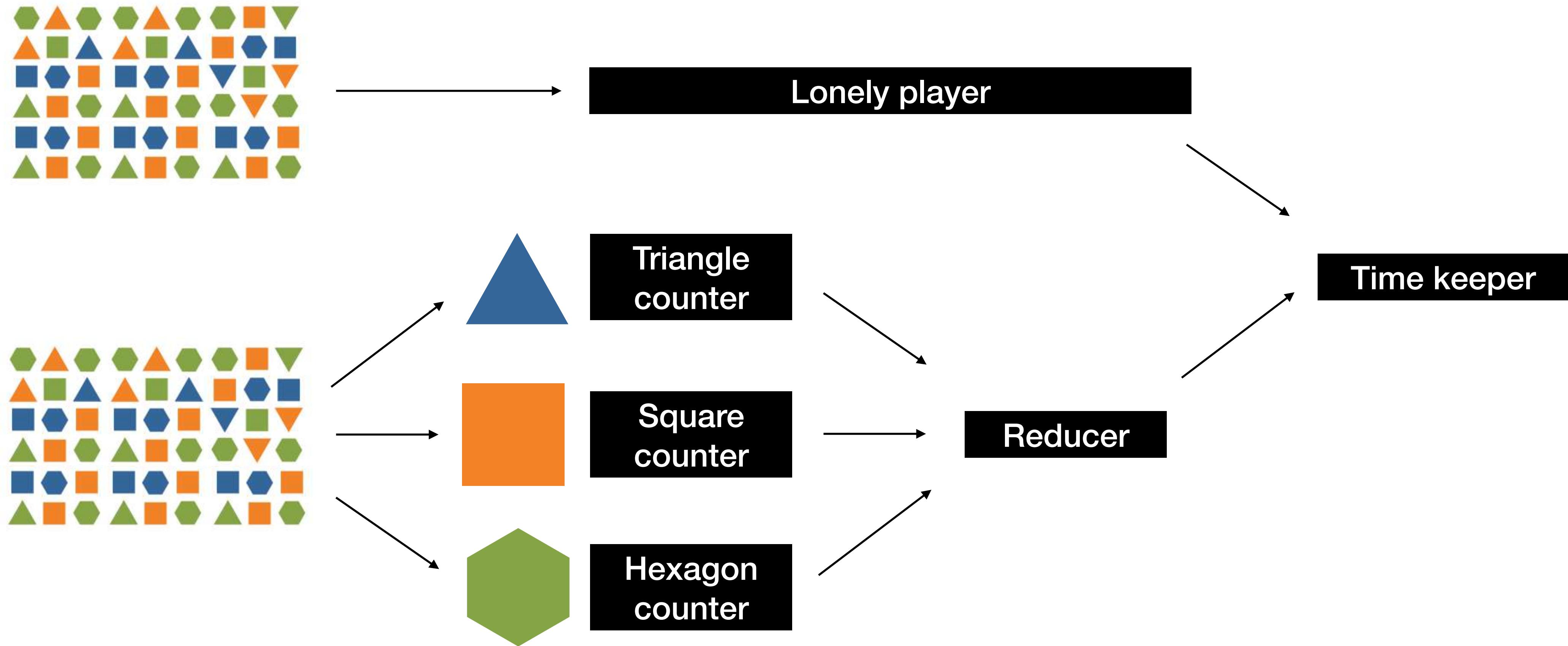
**21**

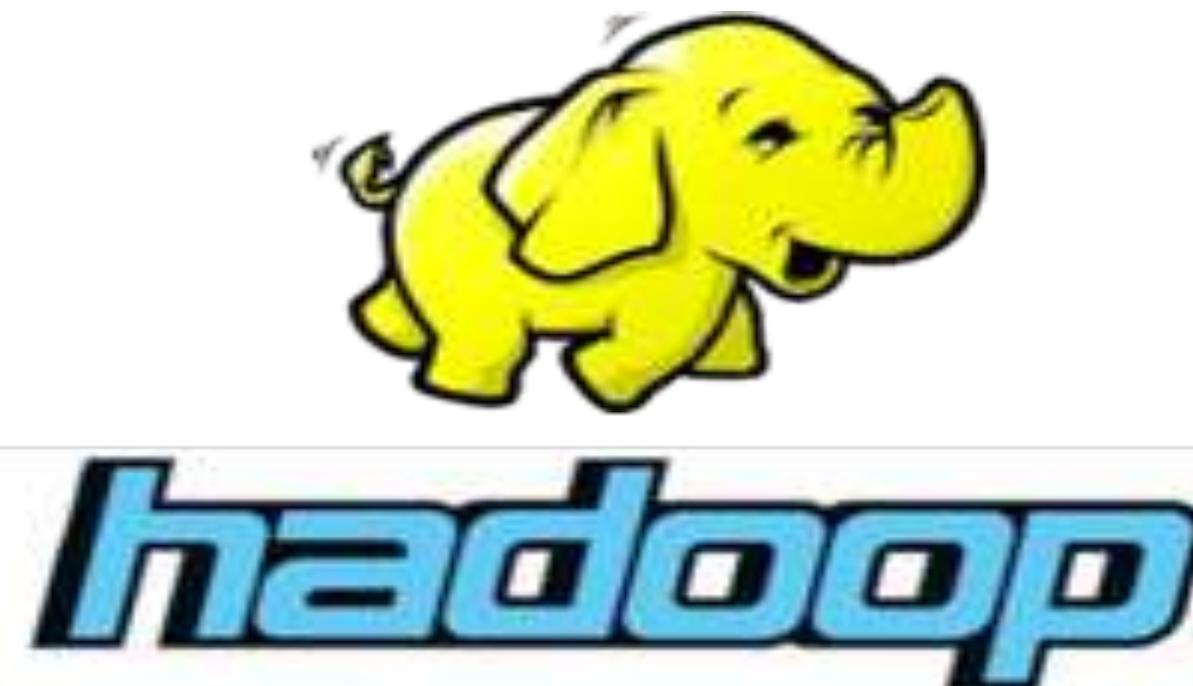


**18**

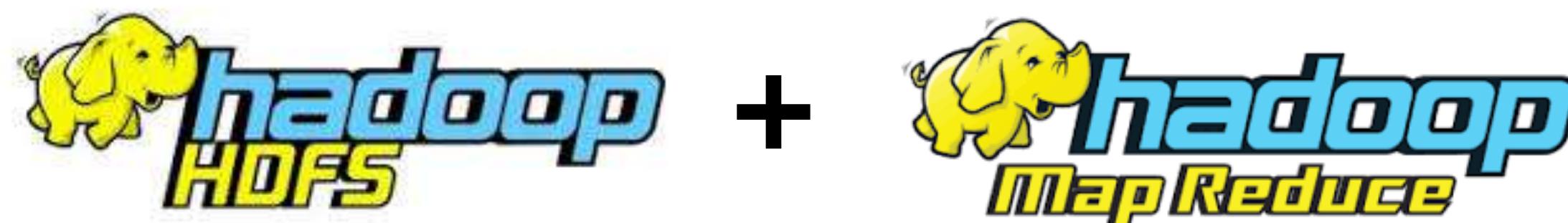
# How to do it

Organize groups of 6 people  
In each group assign the following roles:





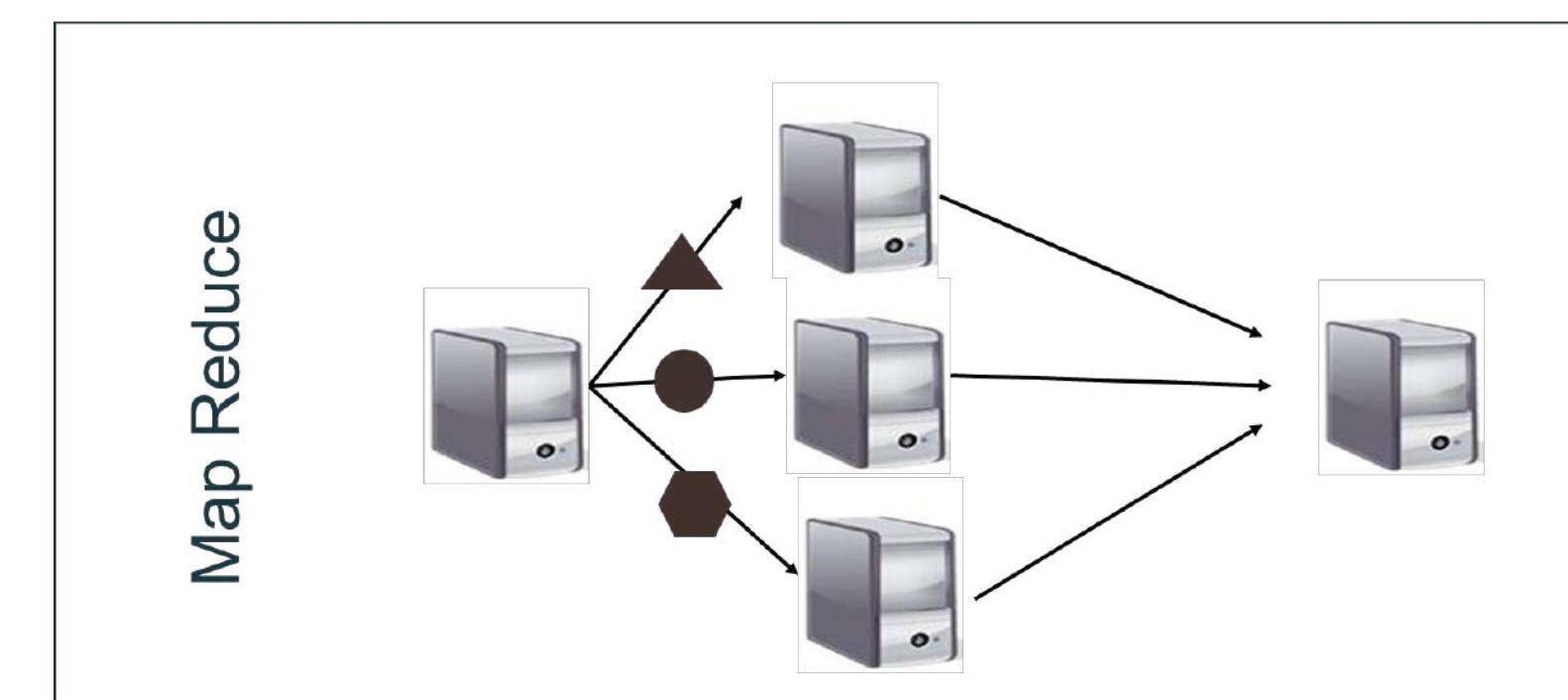
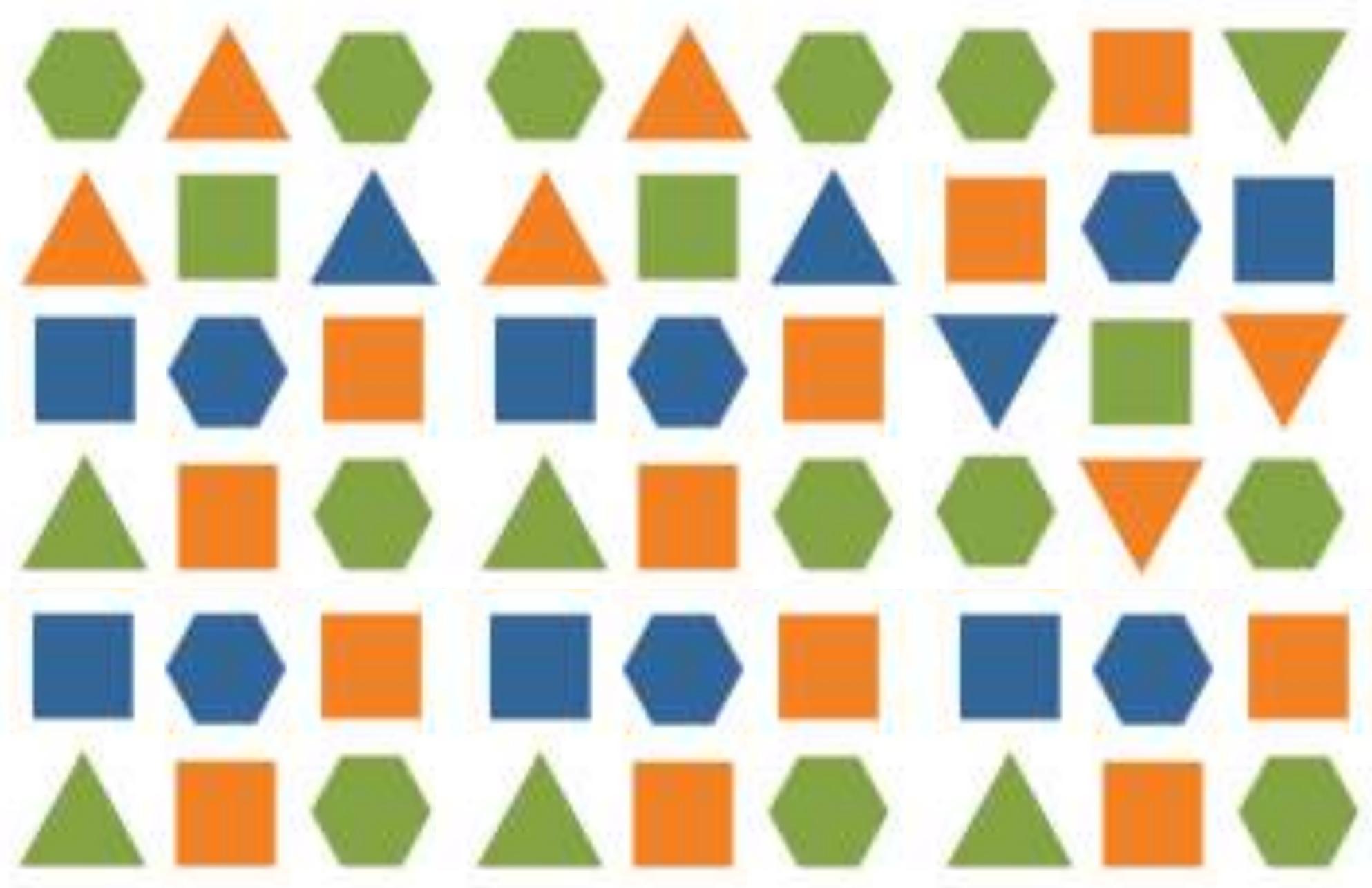
Doug Cutting, 2006



Splits the data between the nodes

Code is transferred to the nodes  
Data is processed in parallel

DATA LOCALITY: The nodes manipulate the data they have access to

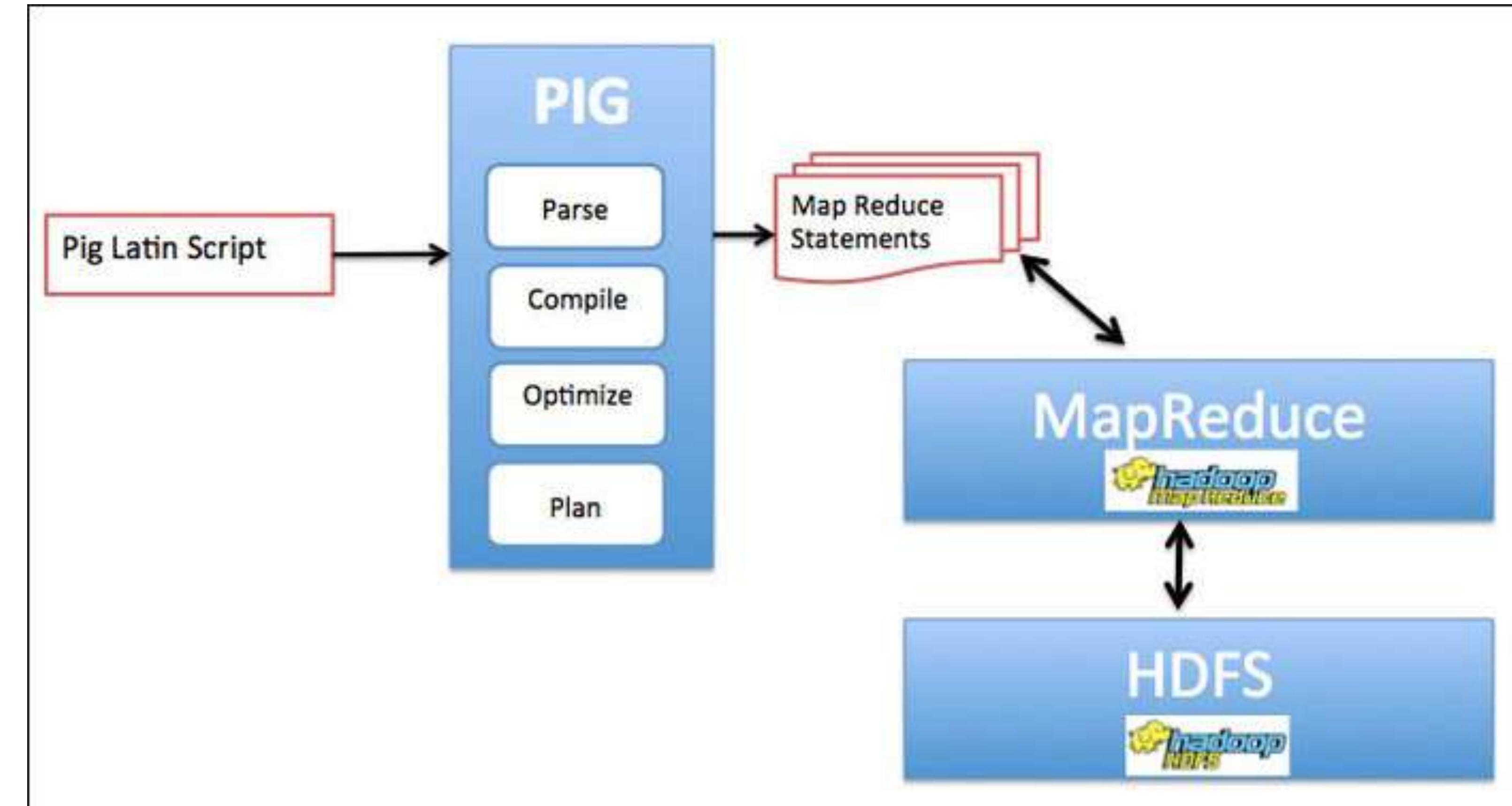




Apache Pig

# Hadoop gets a programming language

- Makes it simpler to write MapReduce programs
- Pig Latin abstracts you from specific details and allows you to focus on the data processing





# Distributed Database over Hadoop





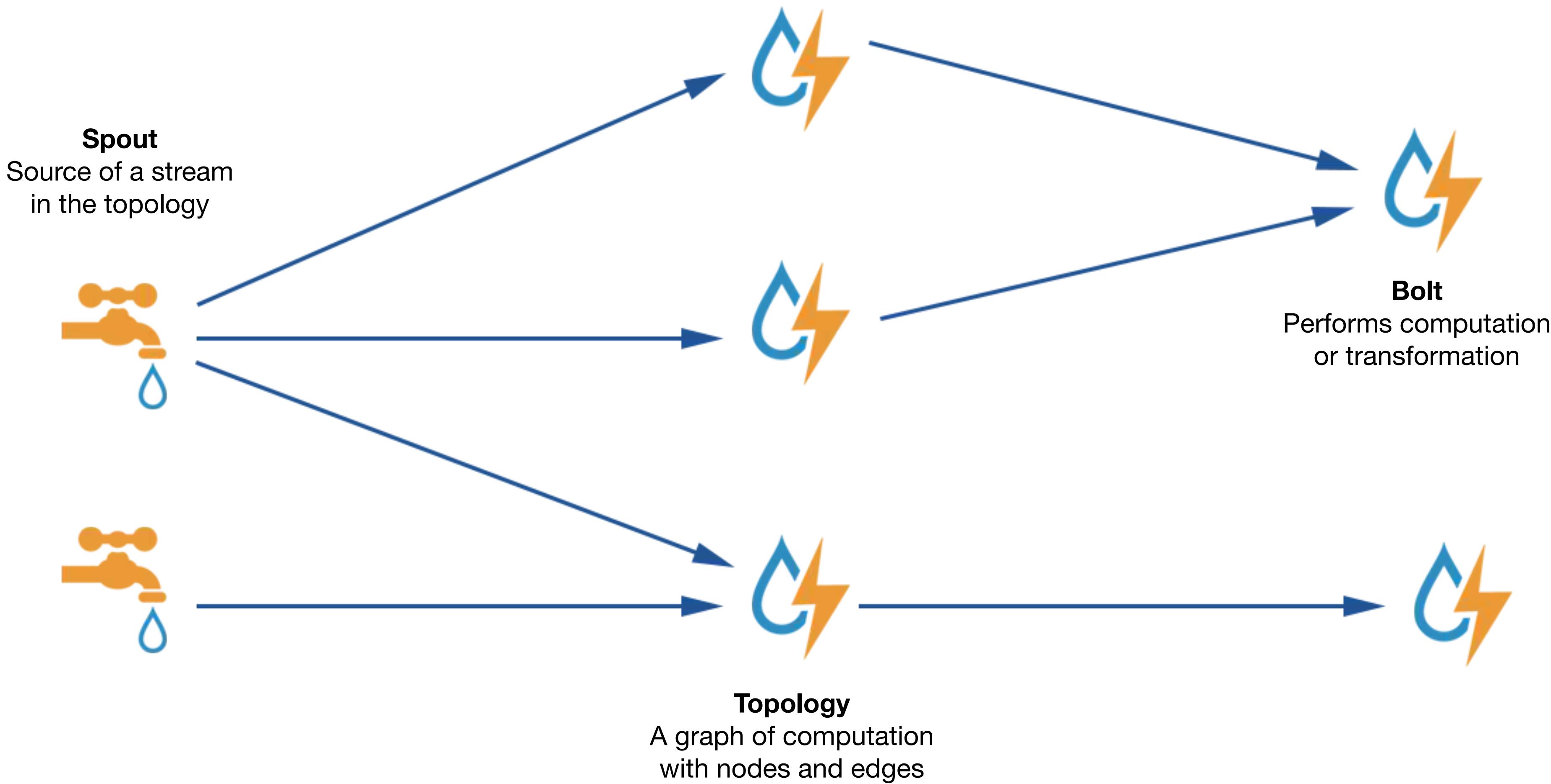
APACHE  
**STORM**<sup>TM</sup>  
Distributed · Resilient · Real-time



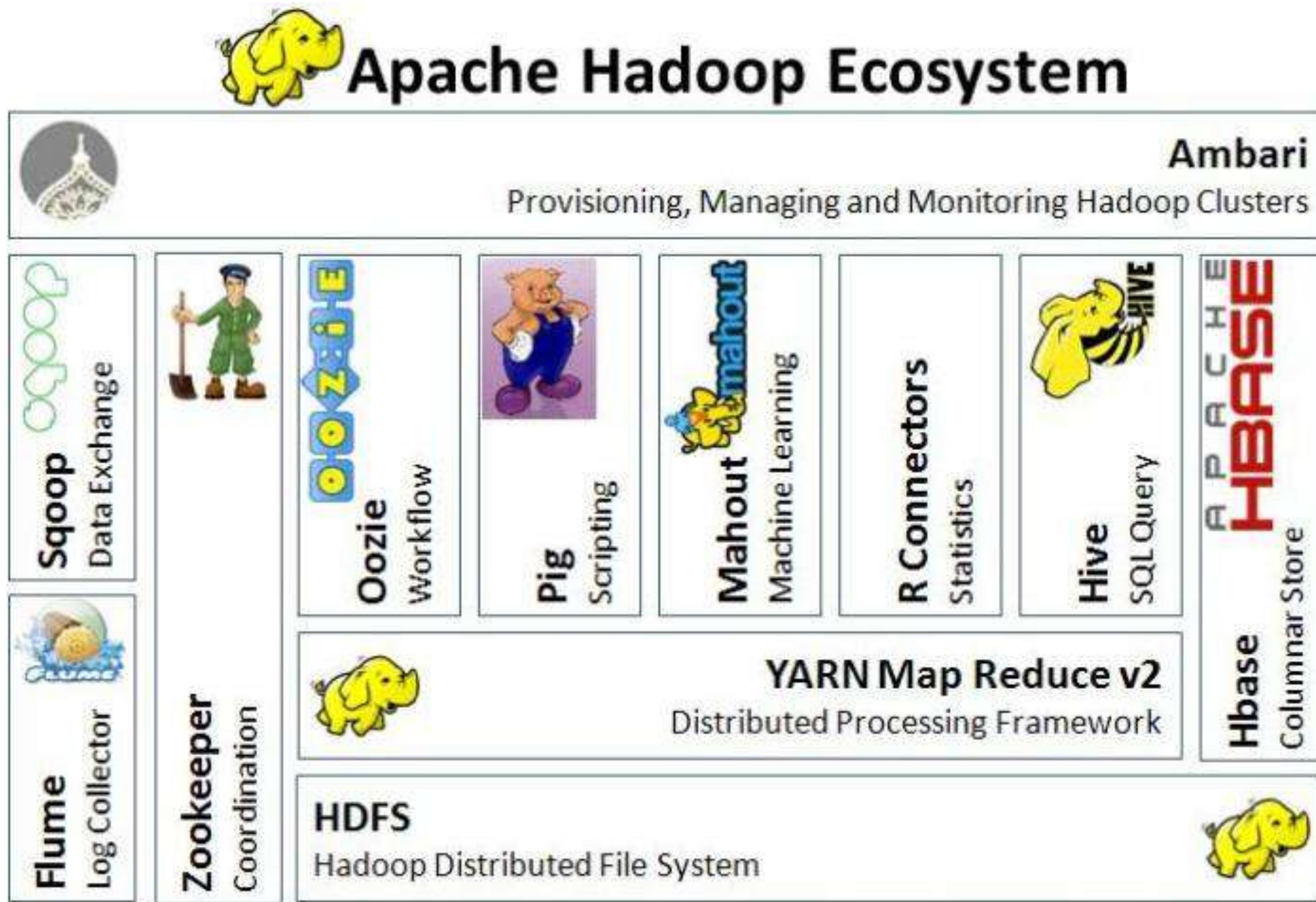
Open sourced after  
acquired by Twitter

# Real-time engine for distributed stream processing

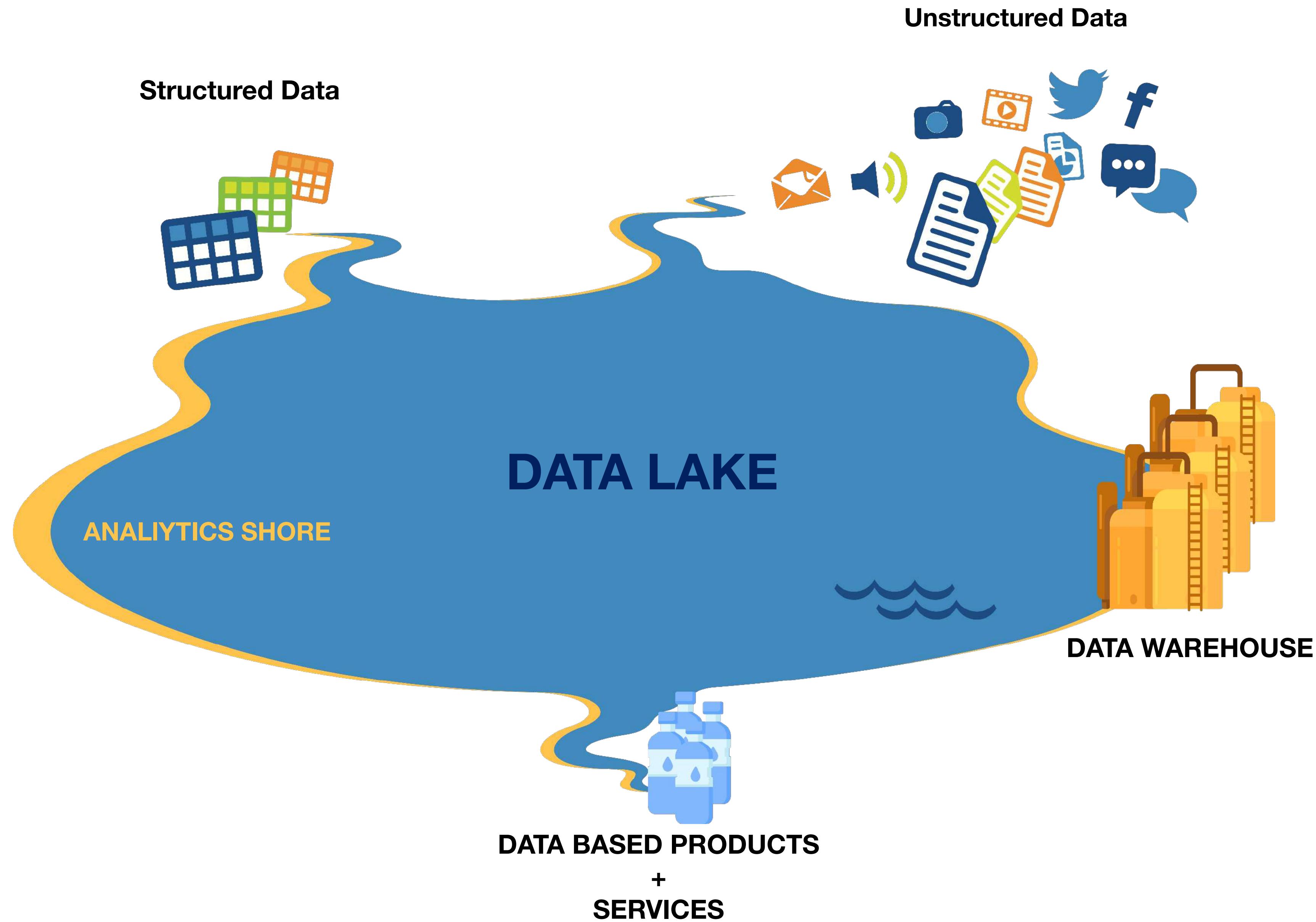
It did for realtime processing what Hadoop did for batch processing



# Hadoop as a Big Data Ecosystem









Machine  
Learning



Analytics



On-premises  
Data Movement



Real-time Data  
Movement



# Enterprise Ready

Cloud / On premises enterprise big data management platforms

**CLOUDERA**

(merged 2019)



(shut down 2019)

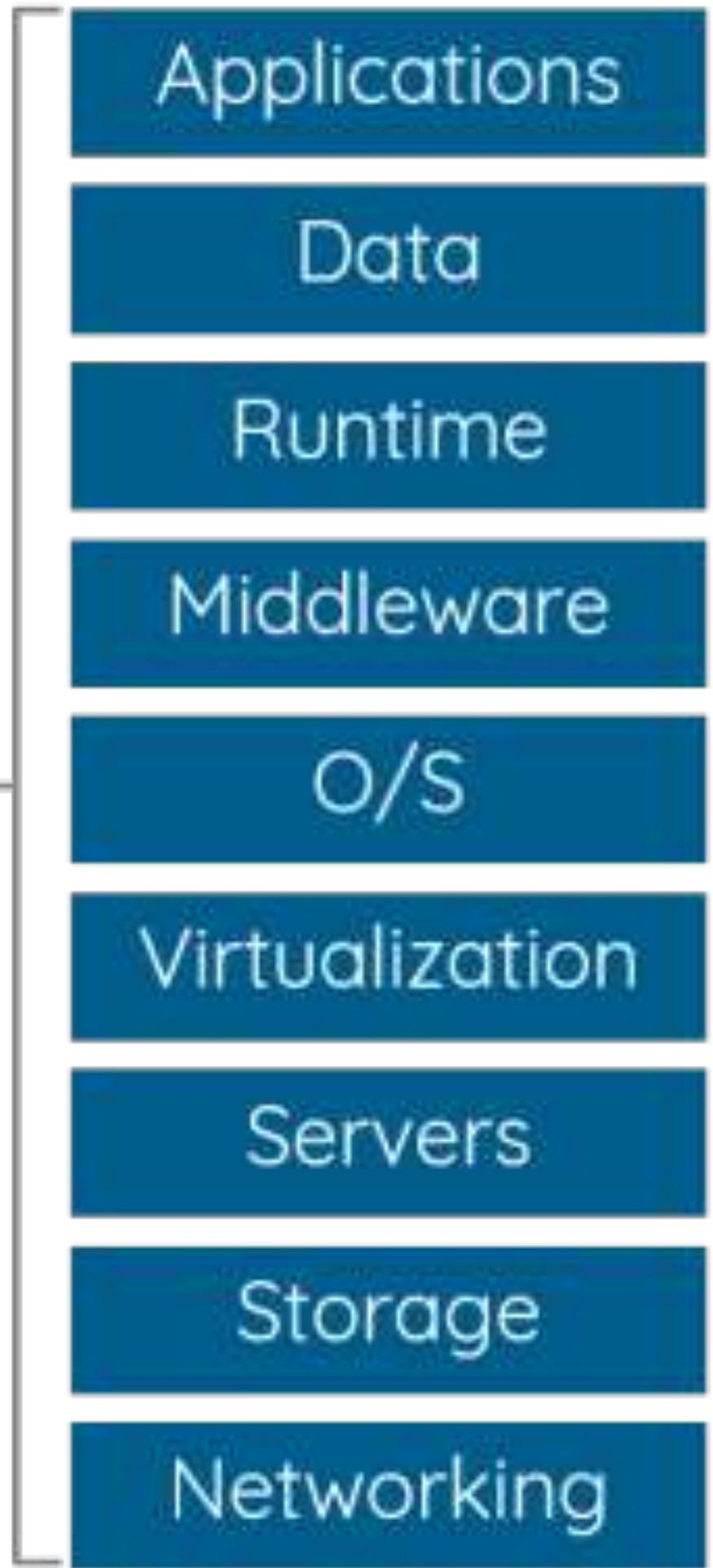




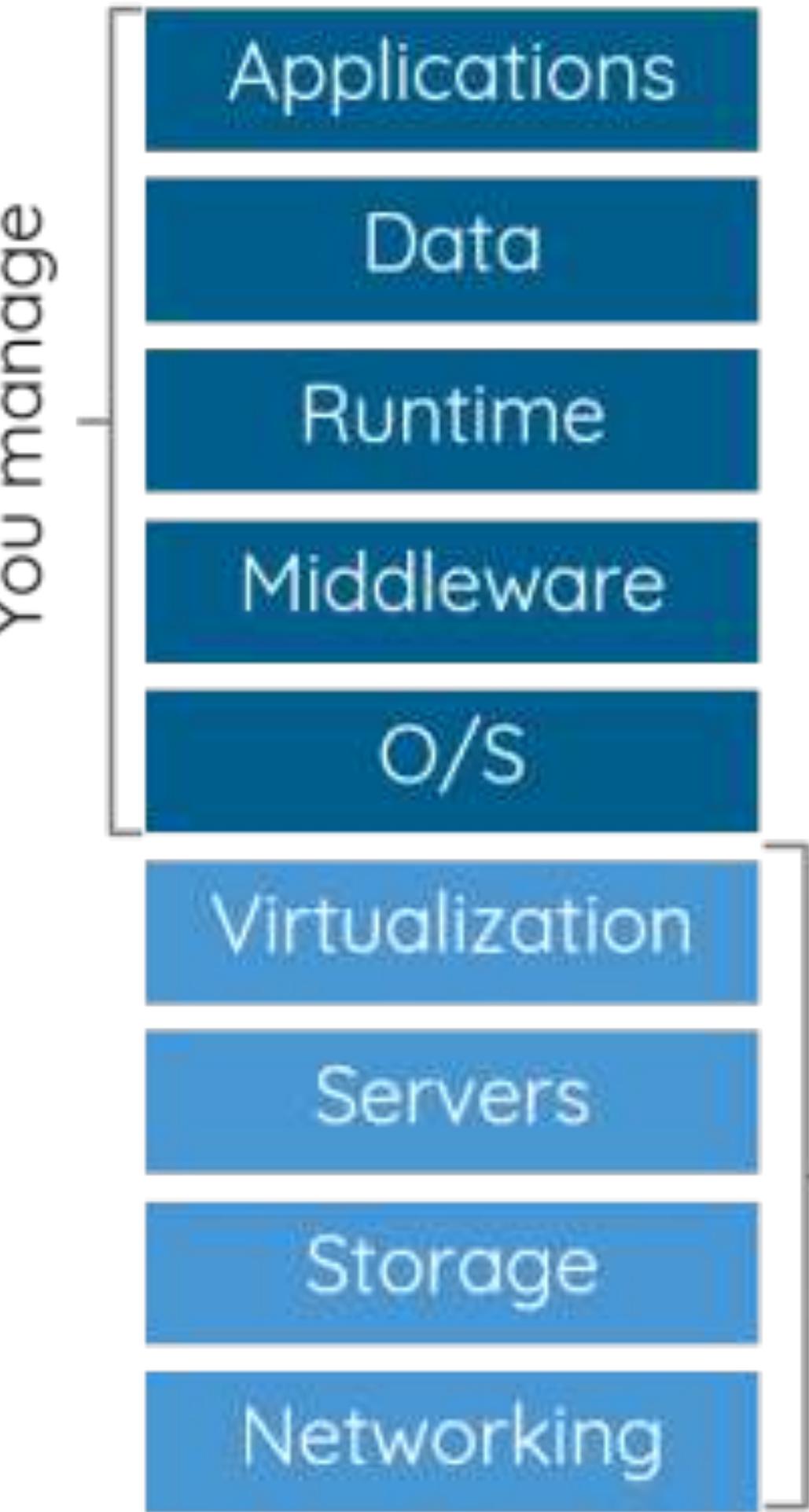
amazon

## Hybrid Cloud

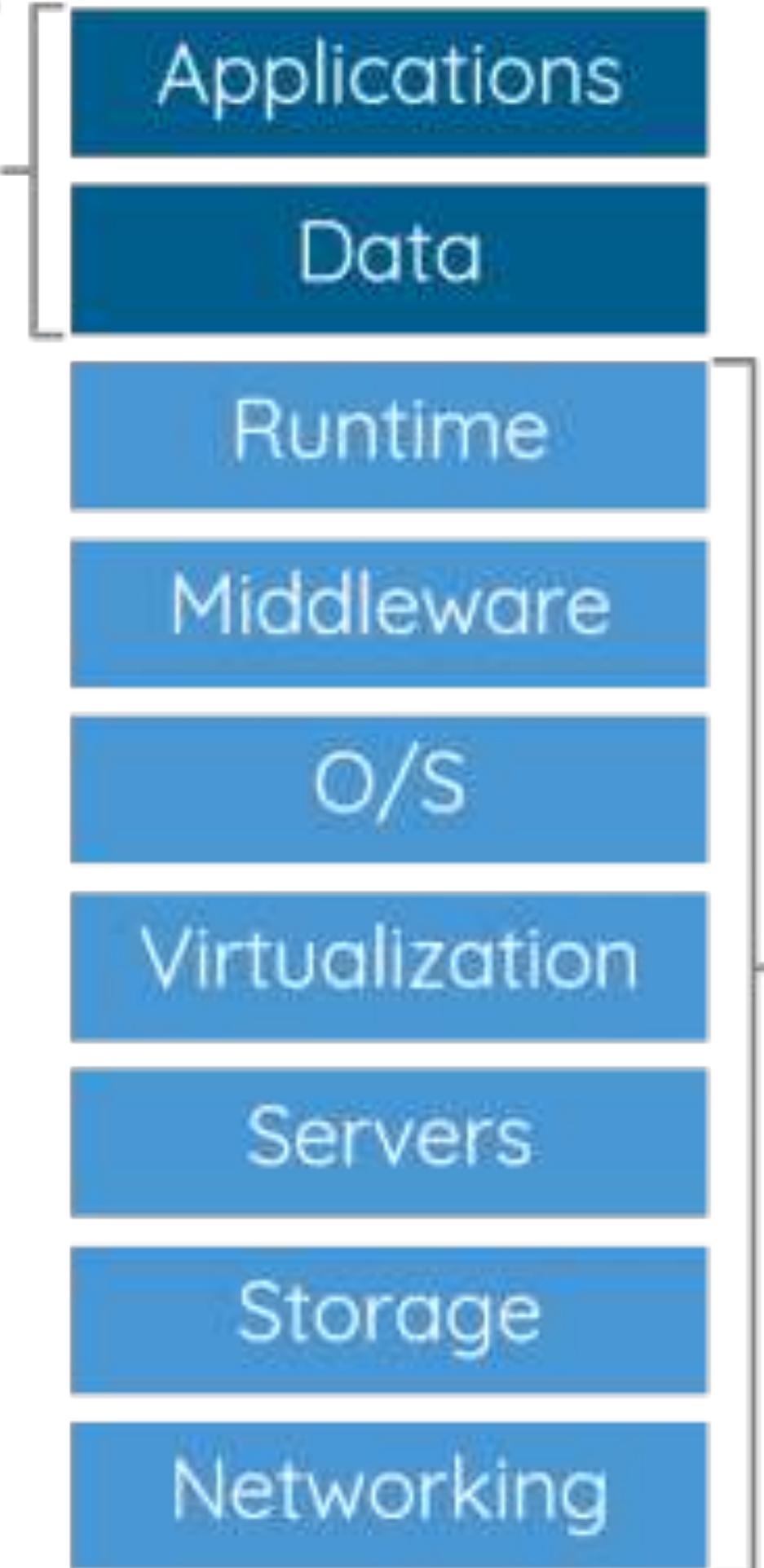
On Premises  
(own server)



IAAS  
(Infrastructure as a Service)



PAAS  
(Platform as a Service)



SAAS  
(Software as a Service)

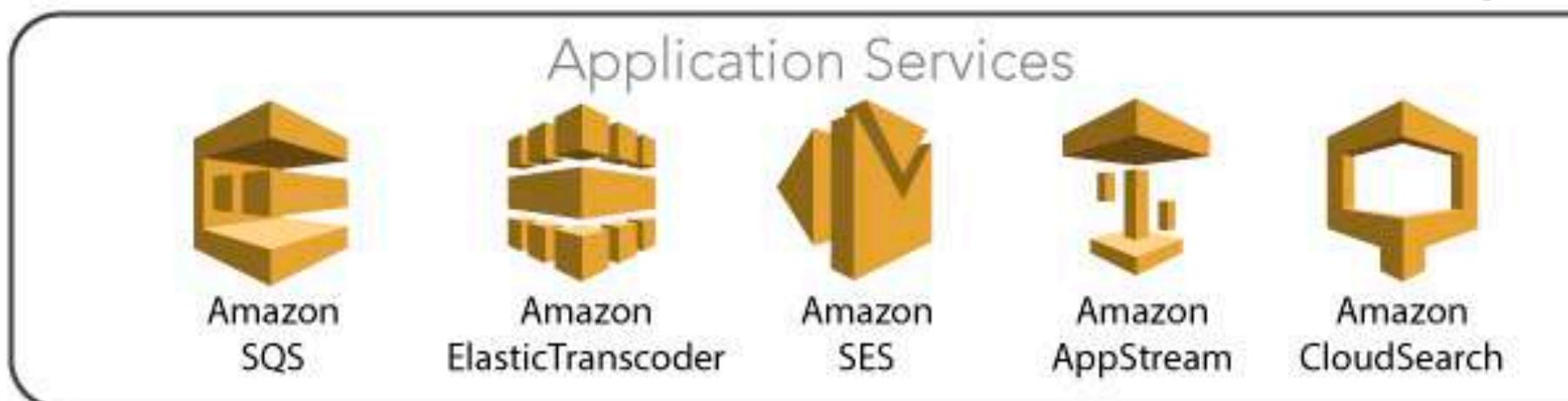


Managed by cloud service provider

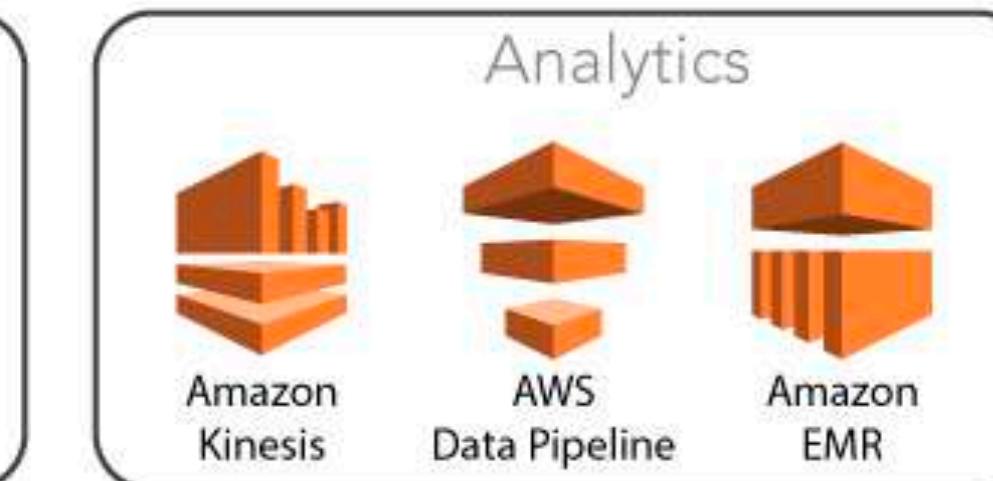


# AWS Services

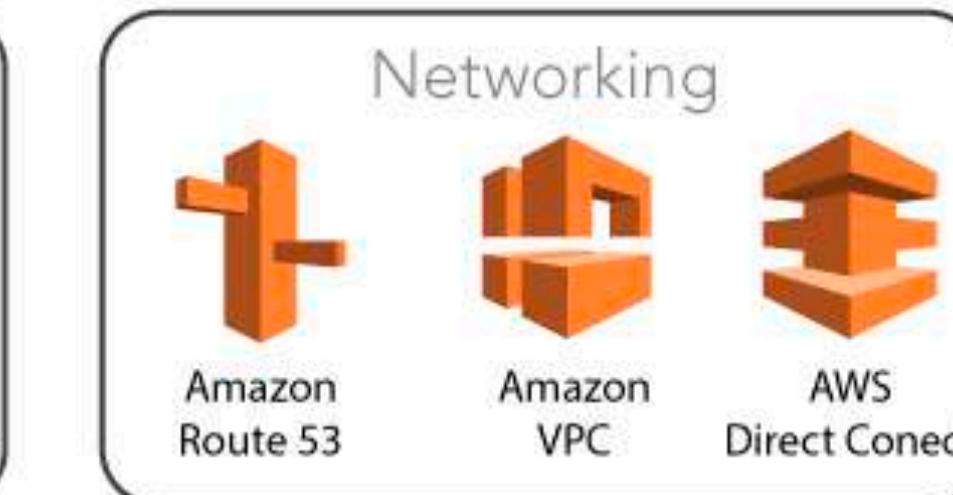
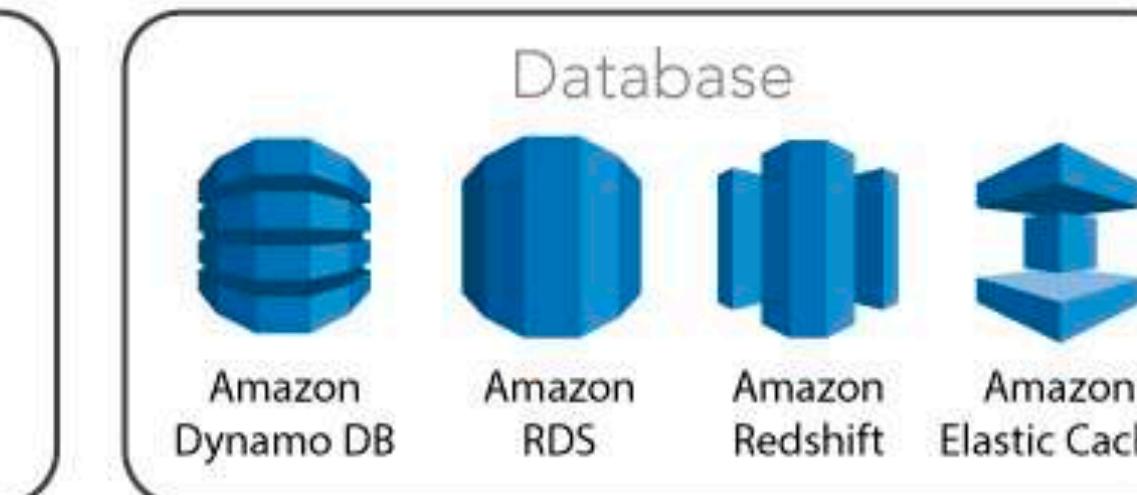
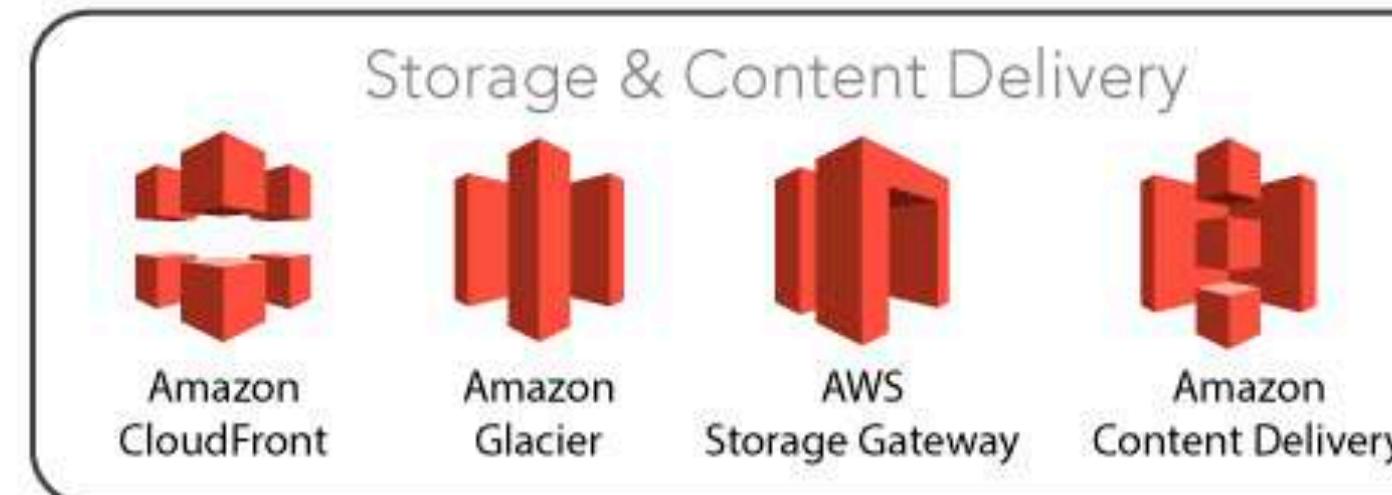
## Deployment & Management

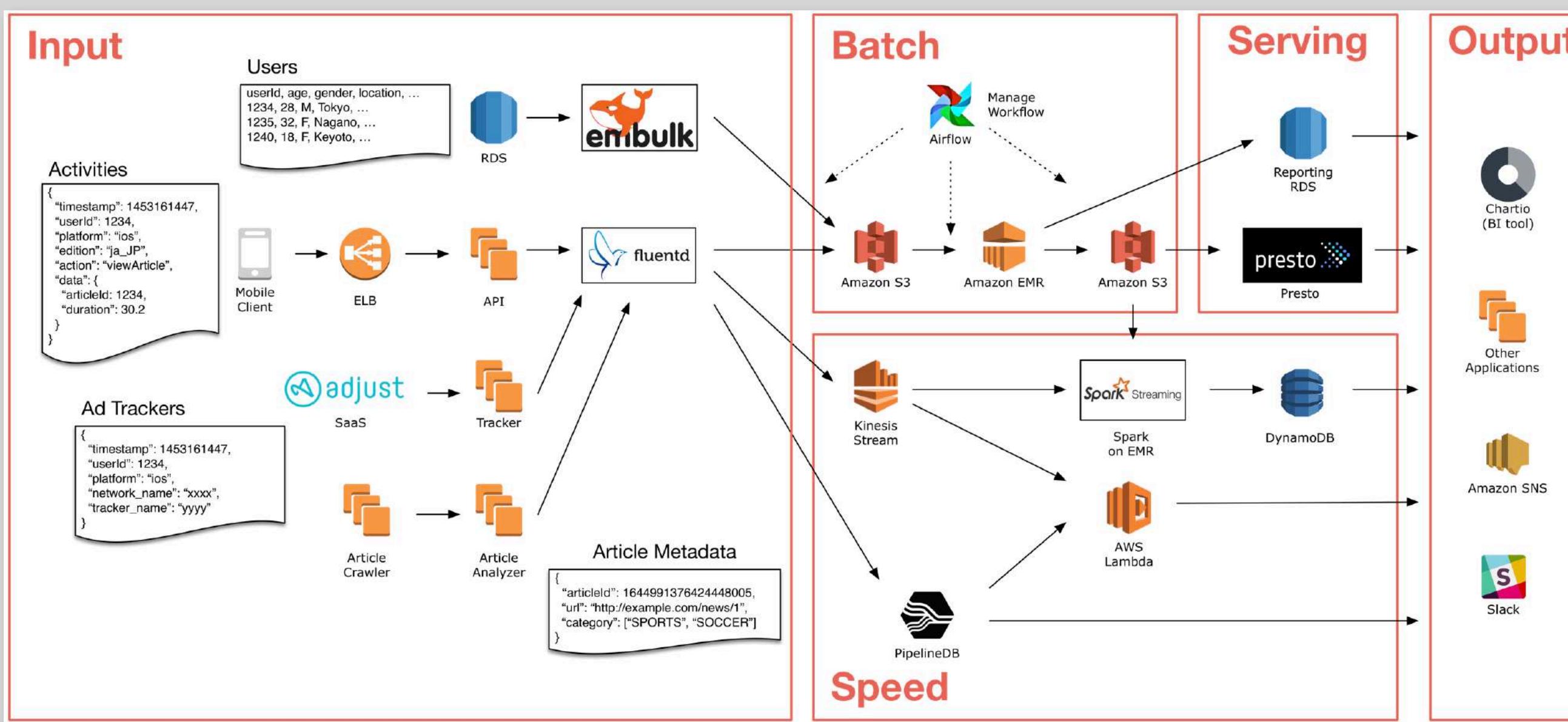
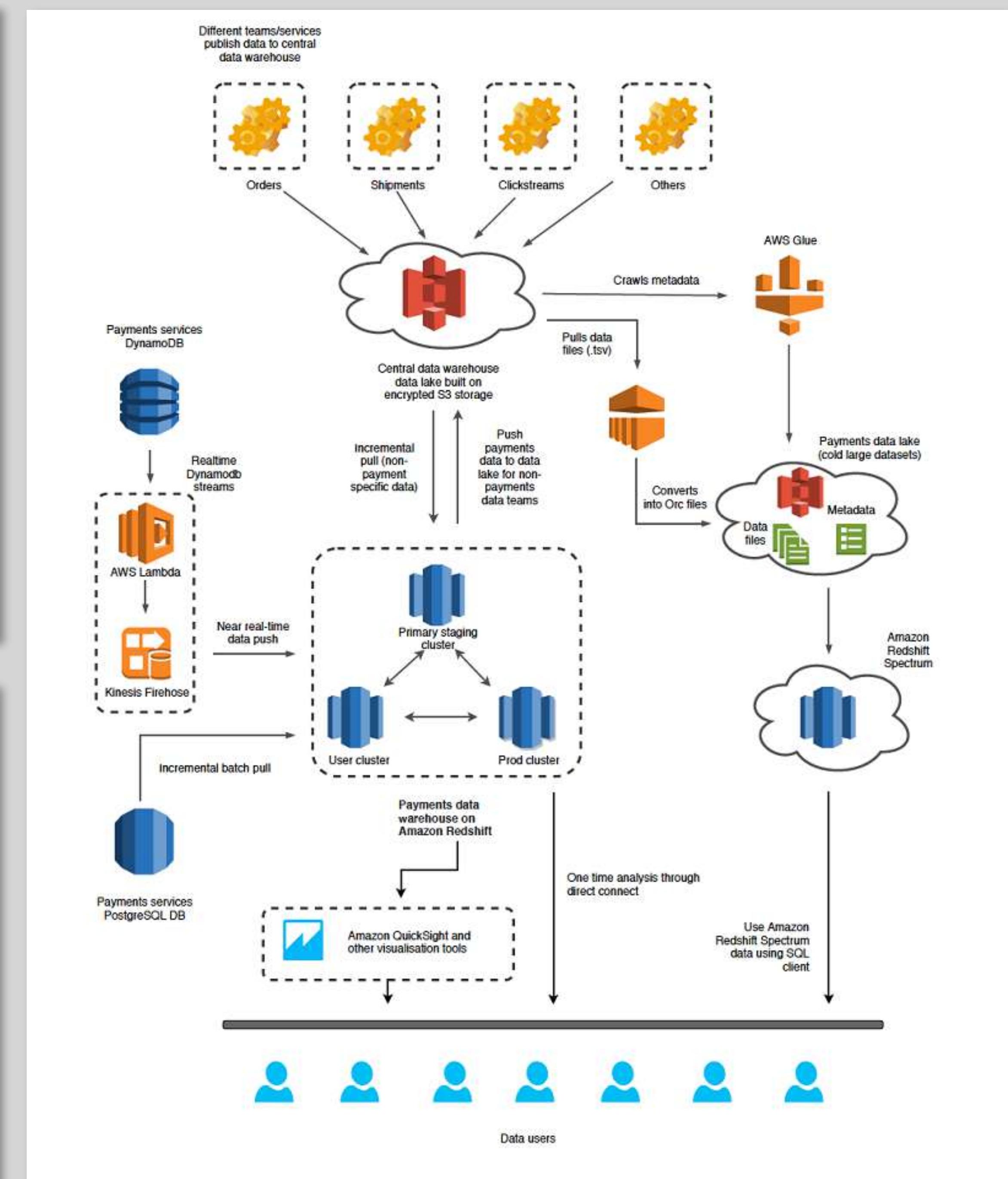
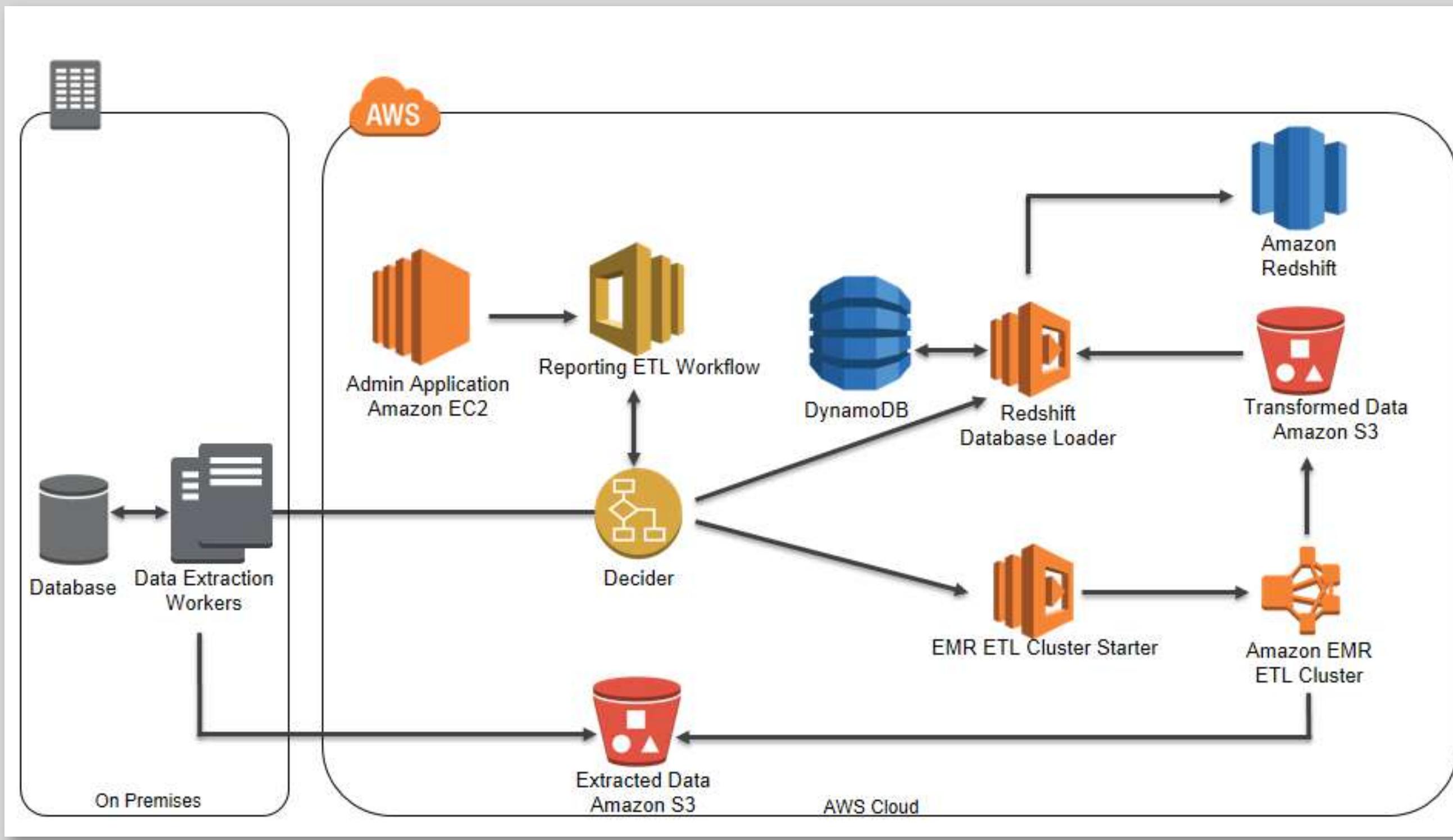


## Application Services



## Foundation Services

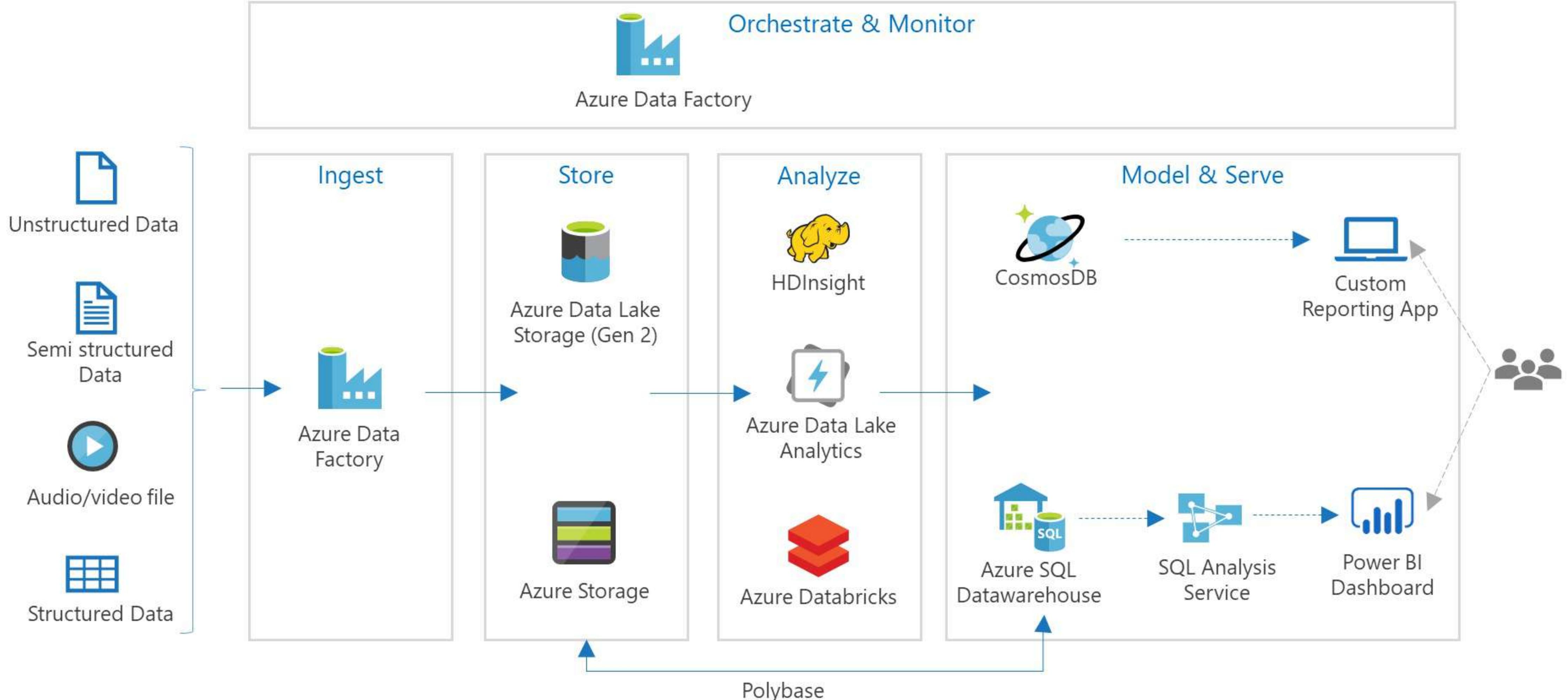




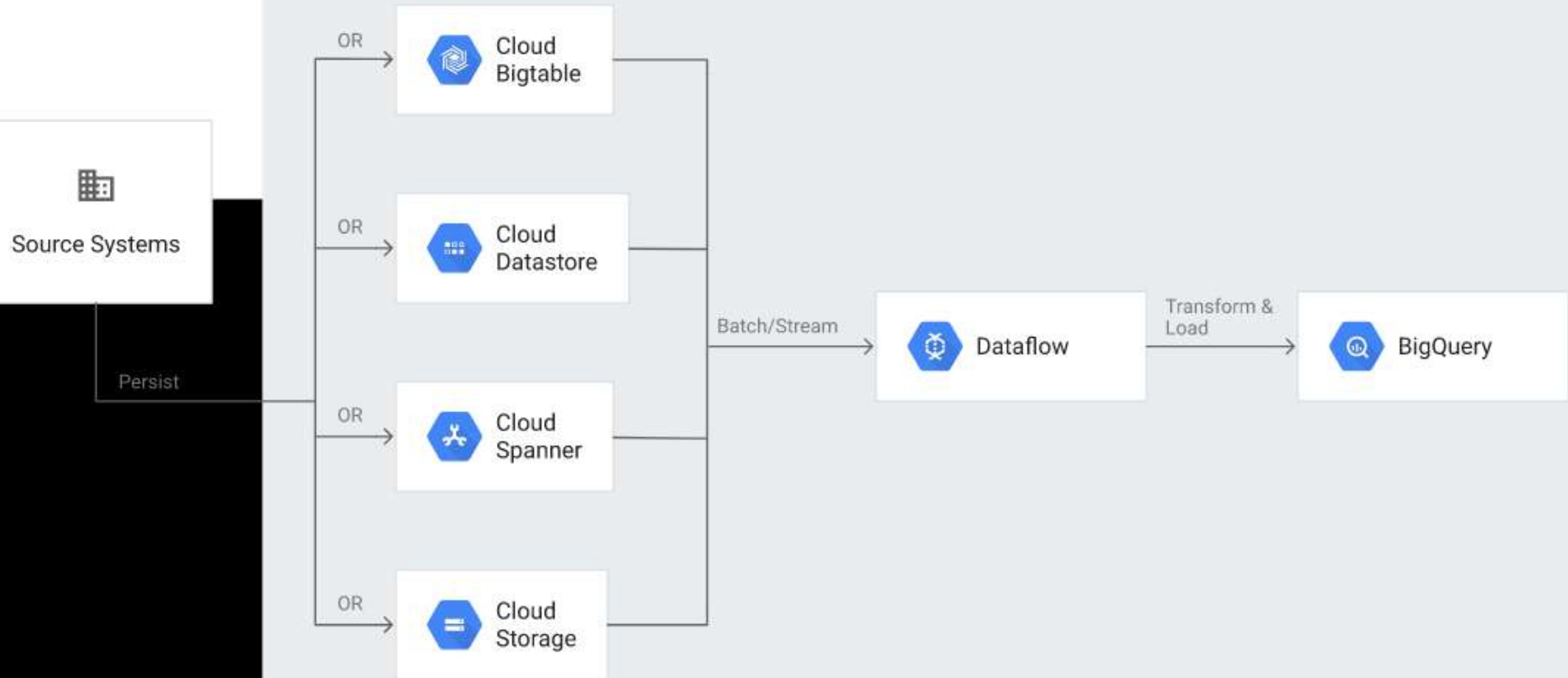
# The Blockbuster Growth of Amazon's Cloud Business

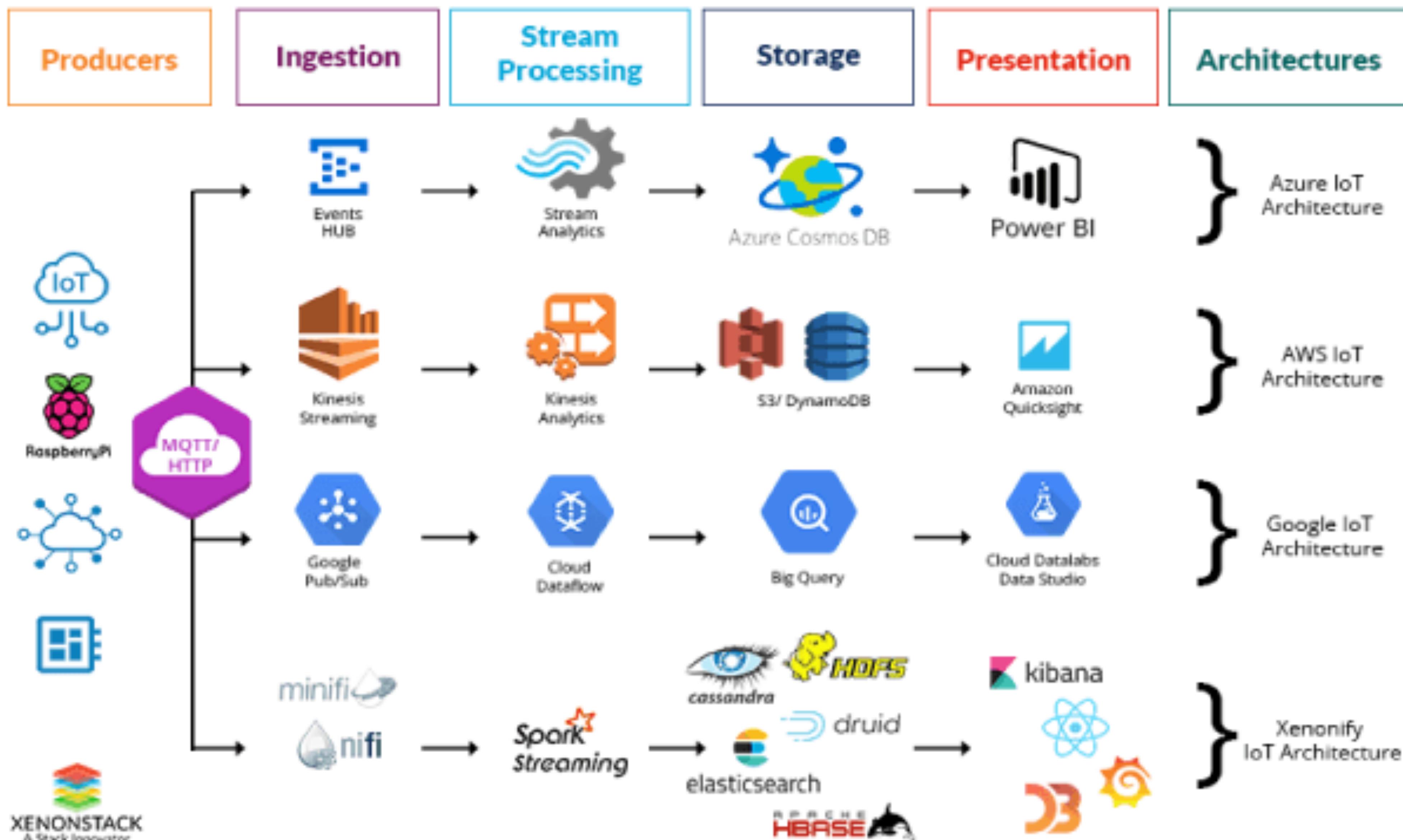
Quarterly revenue of Amazon Web Services\*





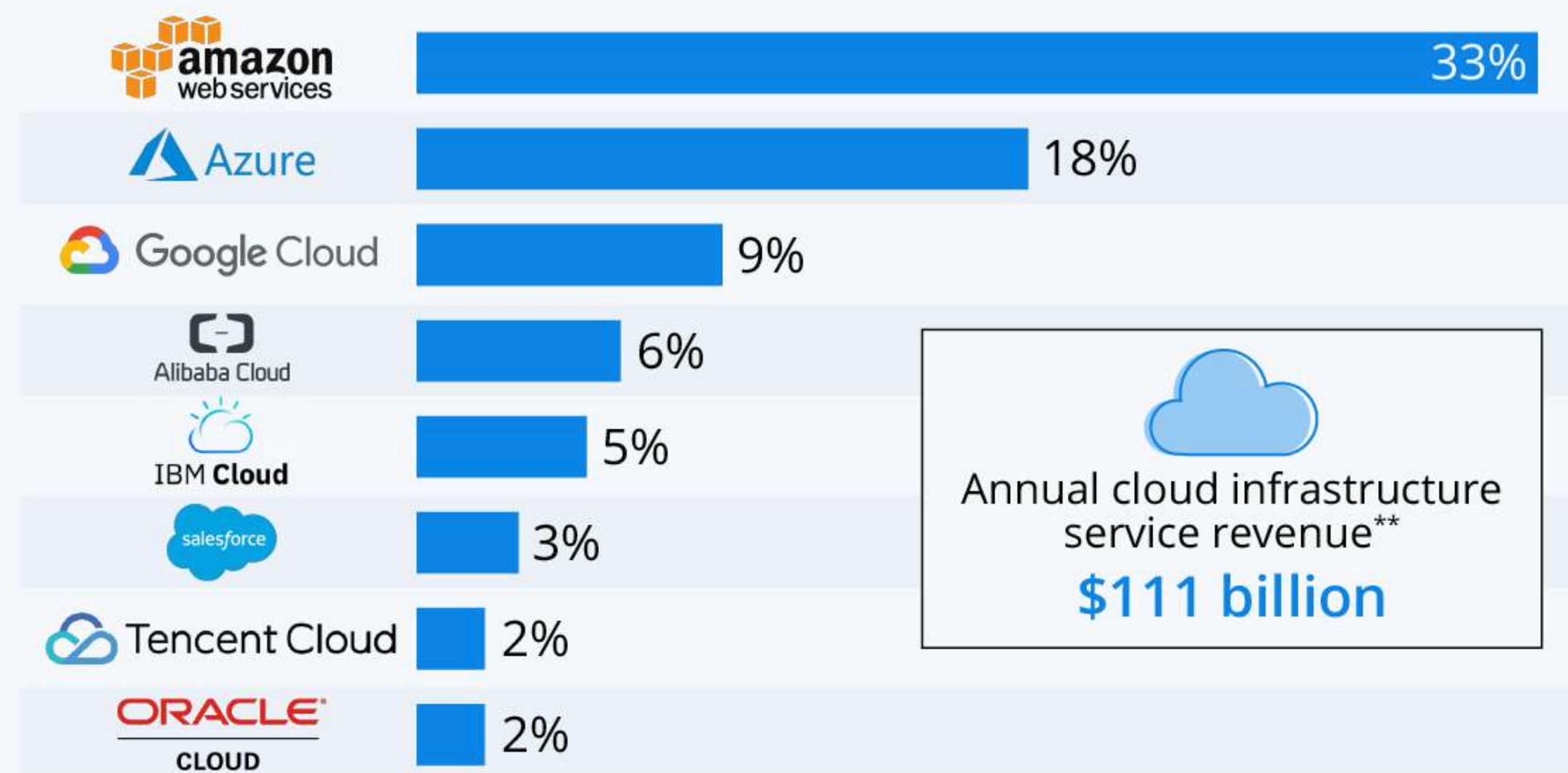
## Google Cloud Platform





# Amazon Leads \$100 Billion Cloud Market

Worldwide market share of leading cloud infrastructure service providers in Q2 2020\*



\* includes platform as a service (PaaS) and infrastructure as a service (IaaS)  
as well as hosted private cloud services

\*\* 12 months ended June 30, 2020

Source: Synergy Research Group



CNIO  
(On premise)

**2010**

2 years

SANITAS  
(IaaS)

**2016**

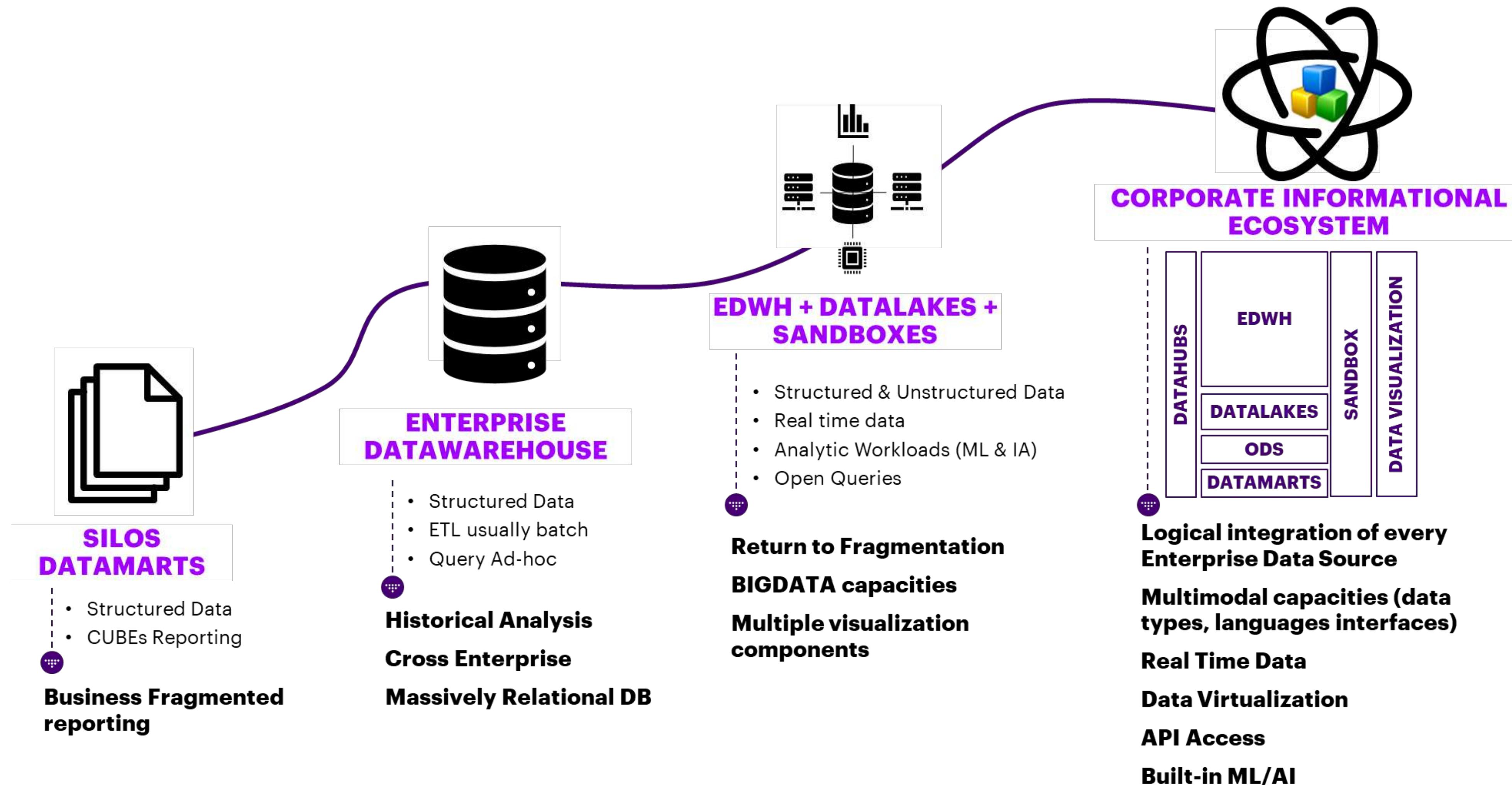
3 months

Genomcore  
(SaaS)

**2020**

<1 hour

# Data Platform evolution



# ARIA



# Our Technological Approach

Cloud infrastructure that yields flexibility, scalability, and rapid deployment

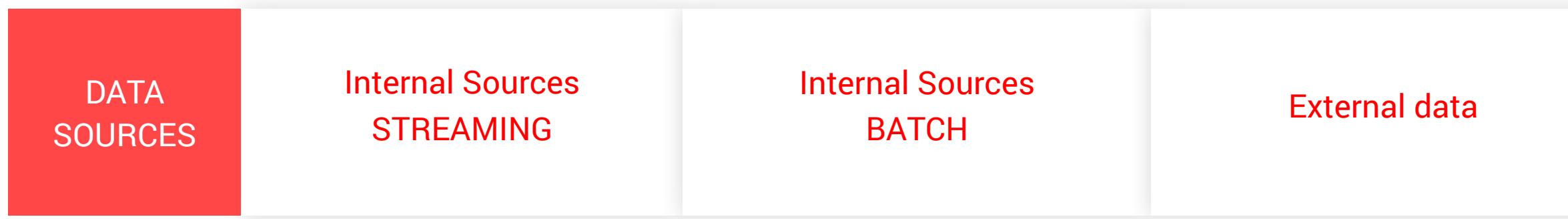


Reducing maintenance and operation costs

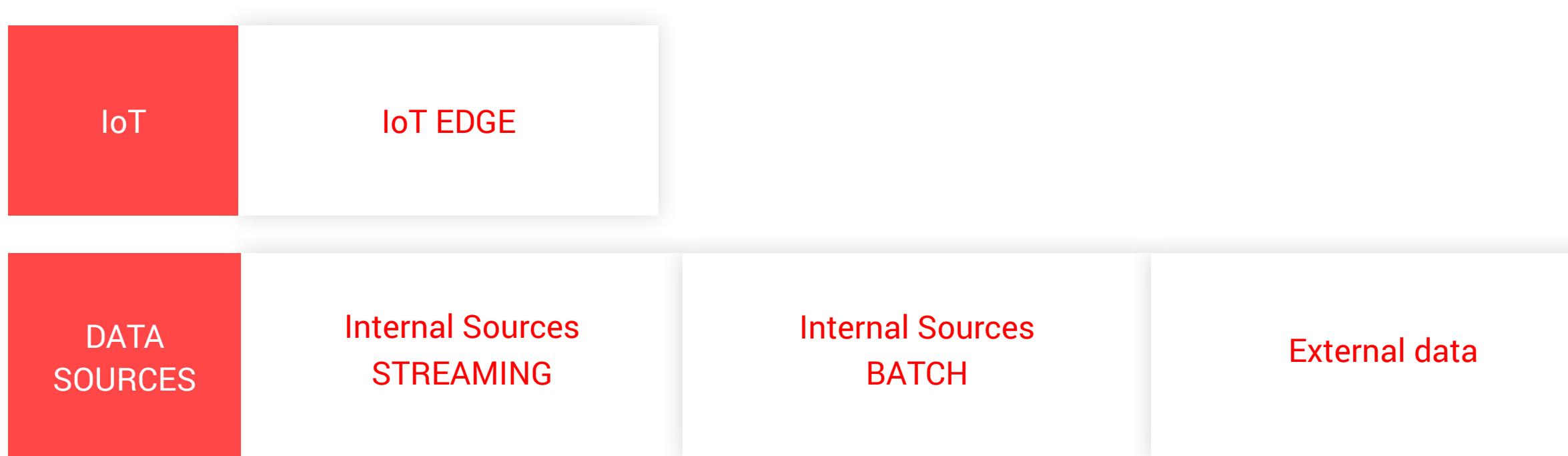


Mostly PaaS that makes possible to increase speed and decrease maintenance.

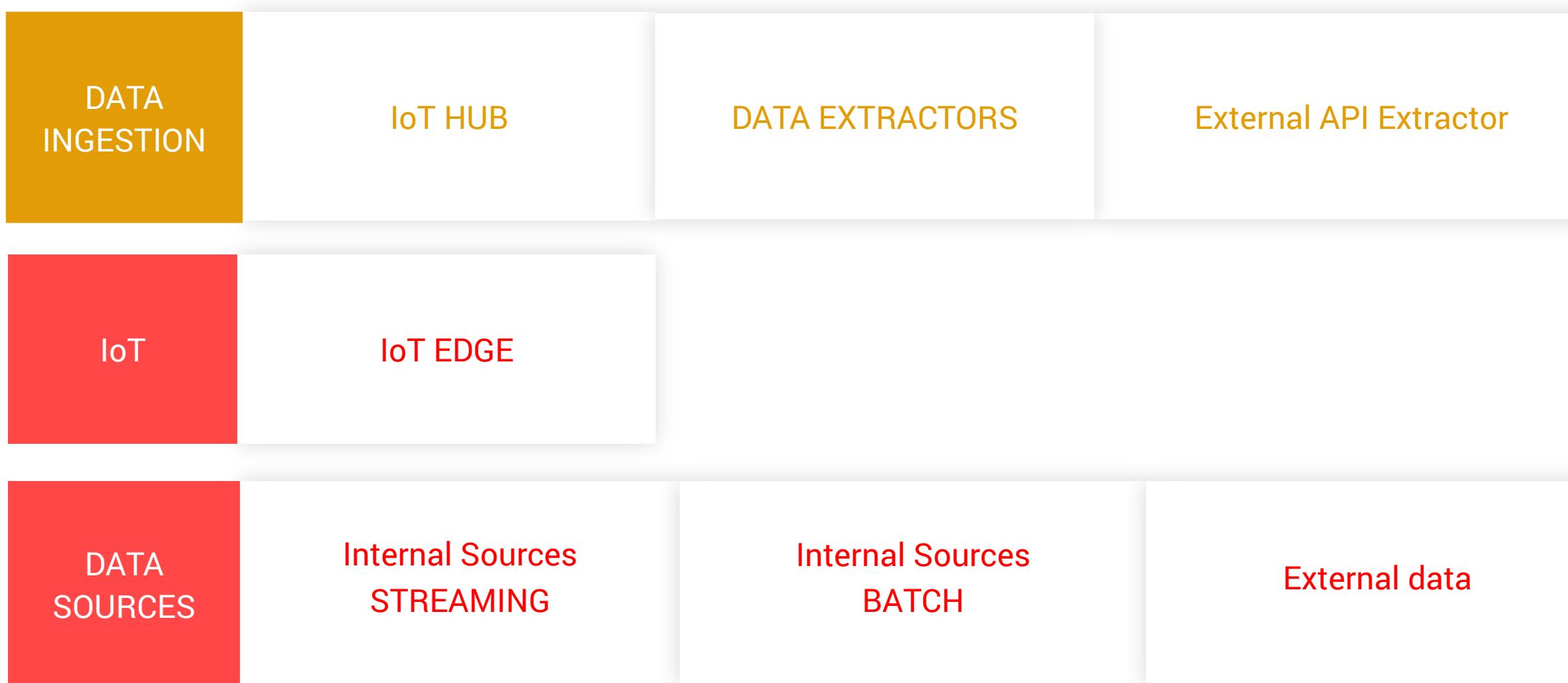
Is it built case by case selecting those component that will be reusable for other use cases



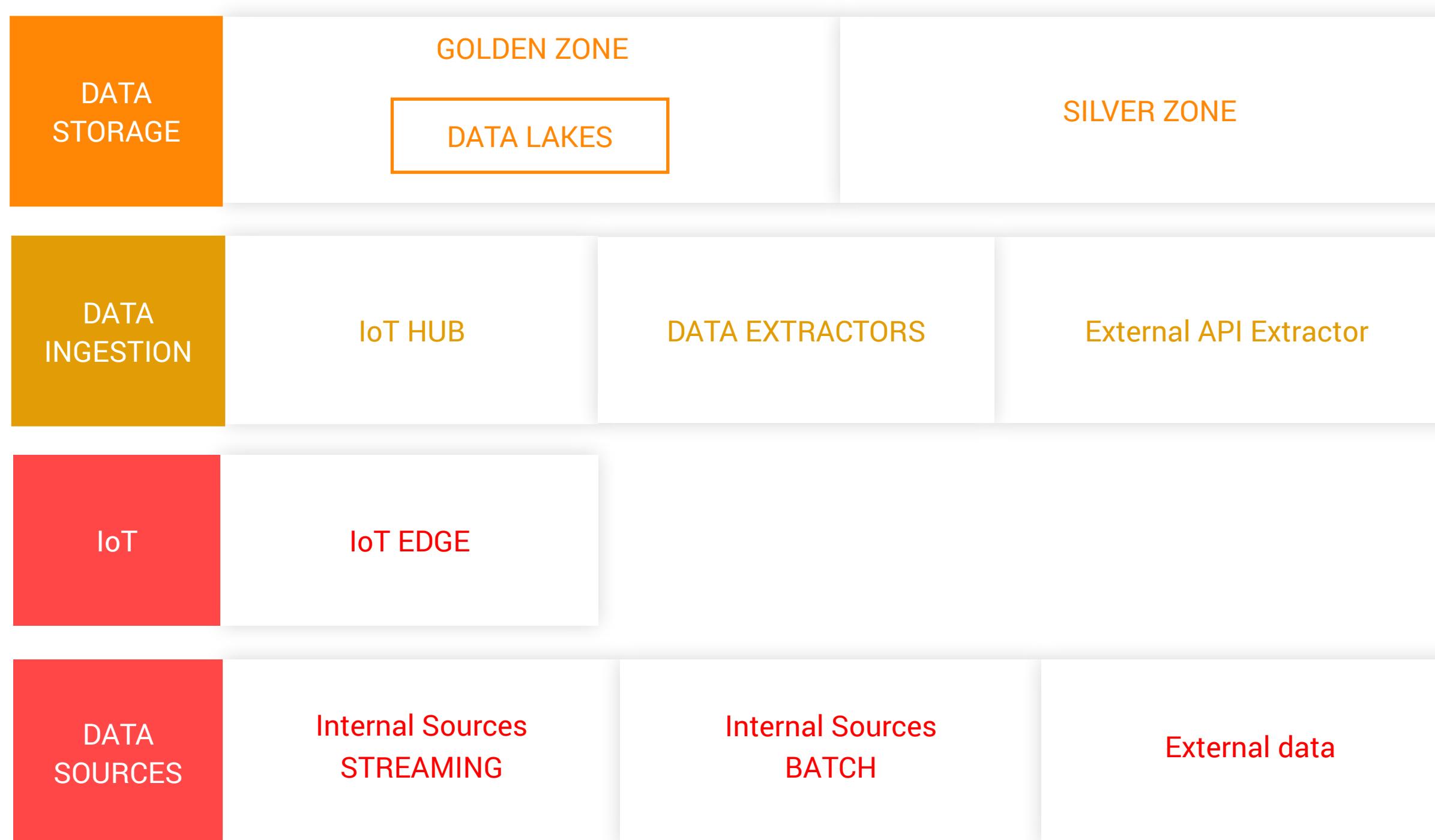
ARIA



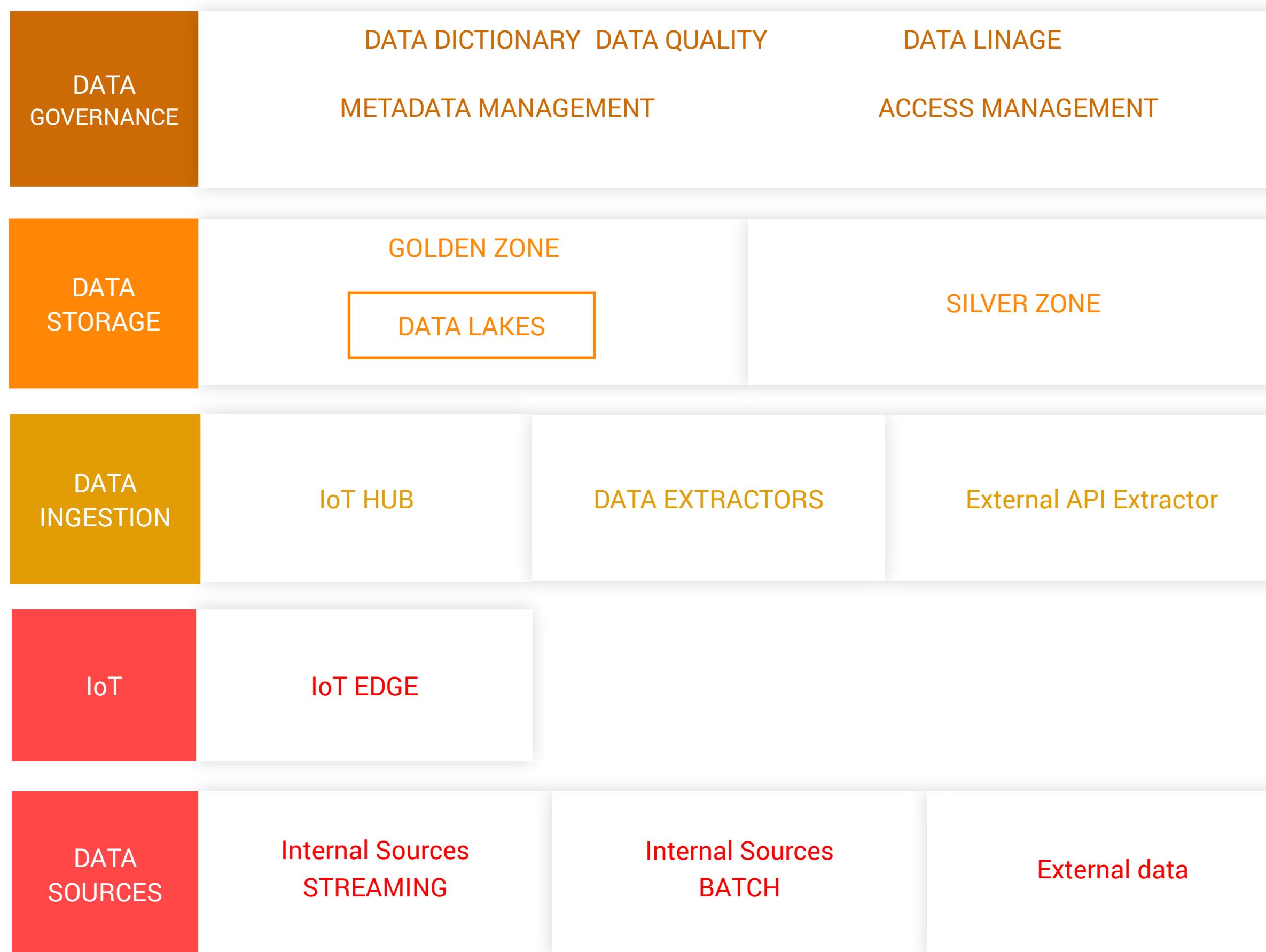
ARIA



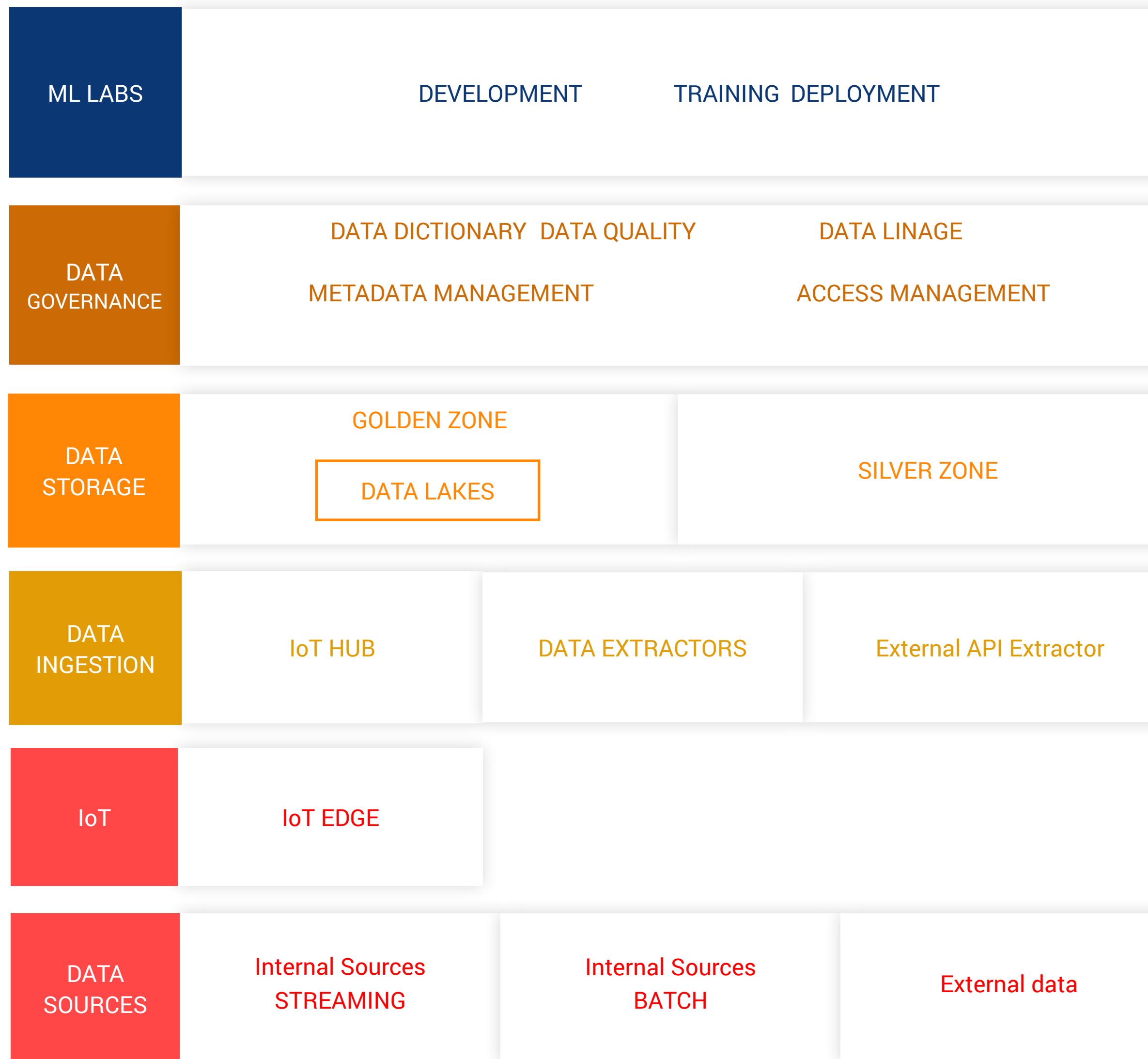
The ARIA logo, consisting of the letters "ARIA" in a stylized, blocky font where each letter has a different color: orange for "A", dark blue for "R", light blue for "I", and dark blue for "A".



ARIA



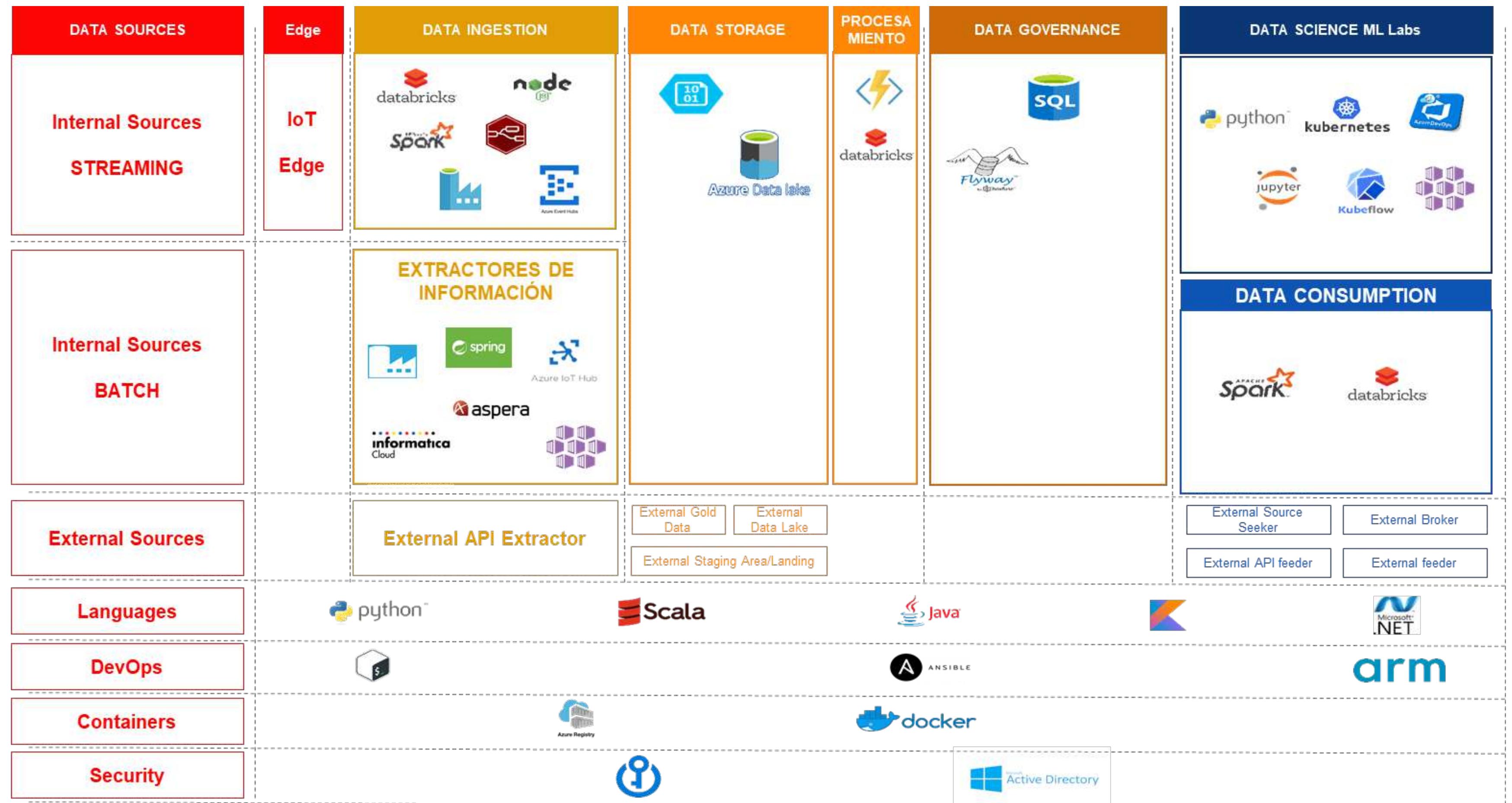
The ARIA logo, consisting of the letters "ARIA" in a stylized, blocky font, with "ARIA" in orange and "RIA" in dark blue.



ARIA



**ARIA**



+40

# Data Platform

¿What is the data management platform (**DMP**) used for collecting and managing data? ¿It allows to use big data and artificial intelligence algorithms to process and analyze large data sets? ¿Can it grow and acquire new functionalities as the data grows and the business needs evolve?



## Data Governance

¿What are the processes and policies in place that ensure that high **data quality** exists through its complete lifecycle? ¿What data **controls** there are that support business objectives? ¿It is the data **available**, **usable**, **consistent** and **secure**?



# Data Governance





## DATA GOVERNANCE FOR ANALYTICS

DATA GOVERNANCE is a set of guidelines and mechanisms based on a series of principles and best practices that provide the company with the tools to guarantee the knowledge, security, quality and availability of information in accordance with organization policies, and with third parties.

### GOALS:



**INVENTORY**  
*Have a complete information inventory avoiding information silos and duplicate data sets*



**LANGUAGE**  
*Establish a common language in the company that allows to univocally identify a data avoiding duplication of information*



**QUALITY**  
*Control and improve the quality of the information, taking care of its completeness, inconsistency and integrity.*



**COMPLIANCE**  
*Comply with regulations, avoiding sanctions based on the traceability of information and the establishment of timely access controls.*



**SECURITY**  
*Improve information security by establishing new roles and responsibilities in the company for the custody of information.*



**TRANSFERS**  
*Establish information transfer policies including roles and responsibilities of the parties involved.*



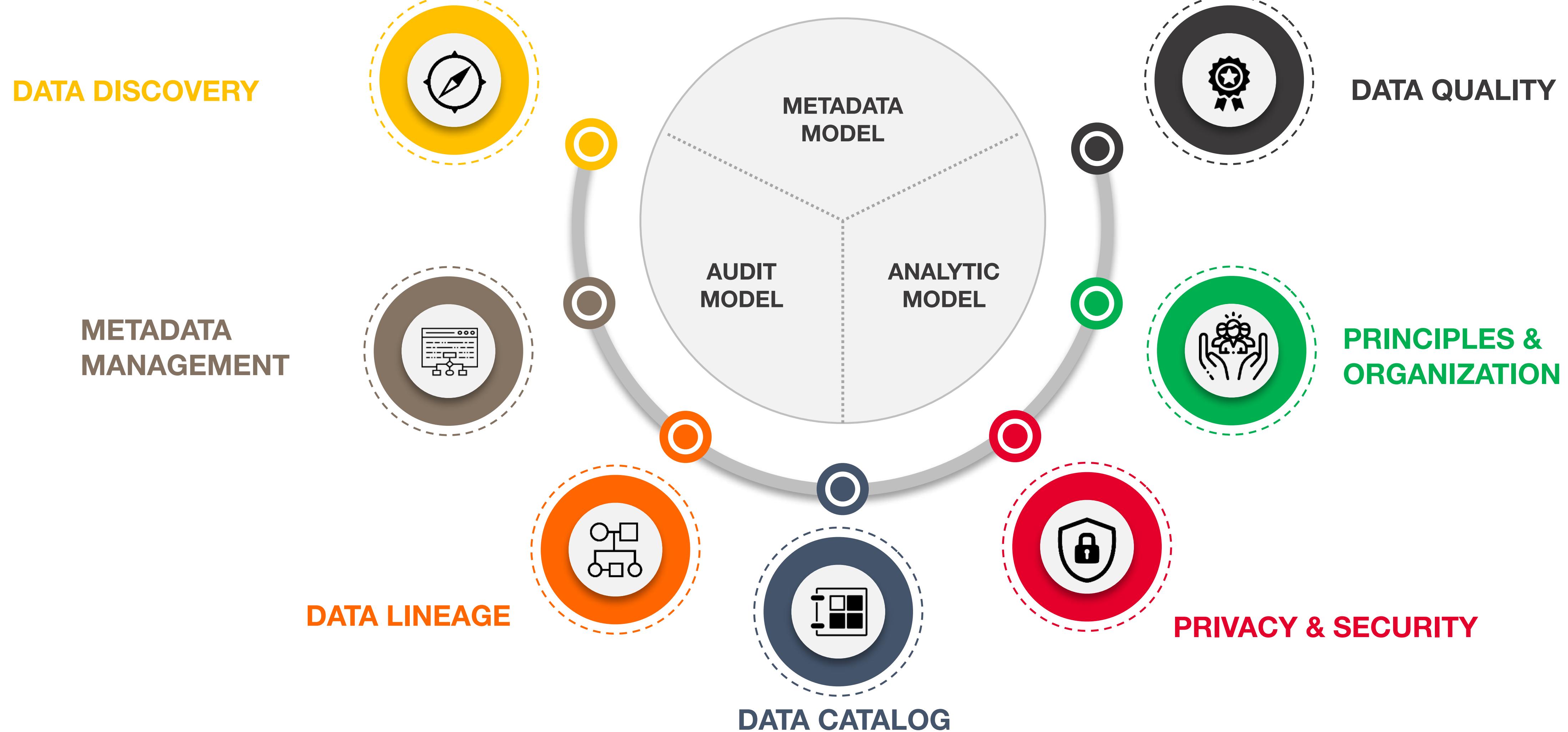
**TIME TO MARKET**  
*Optimize the execution of analytical projects avoiding the reprocessing of searches, transfers and cleaning of information.*



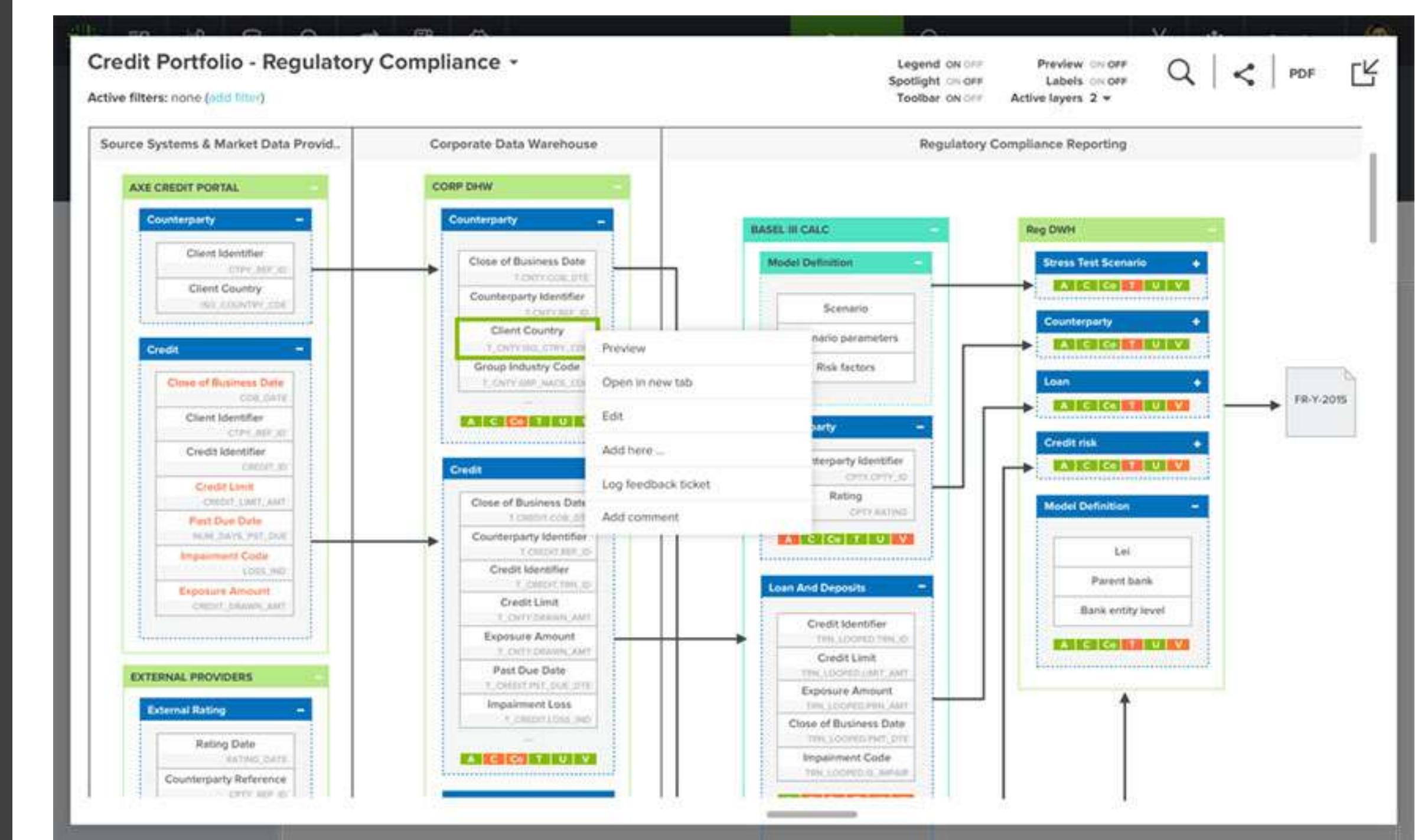
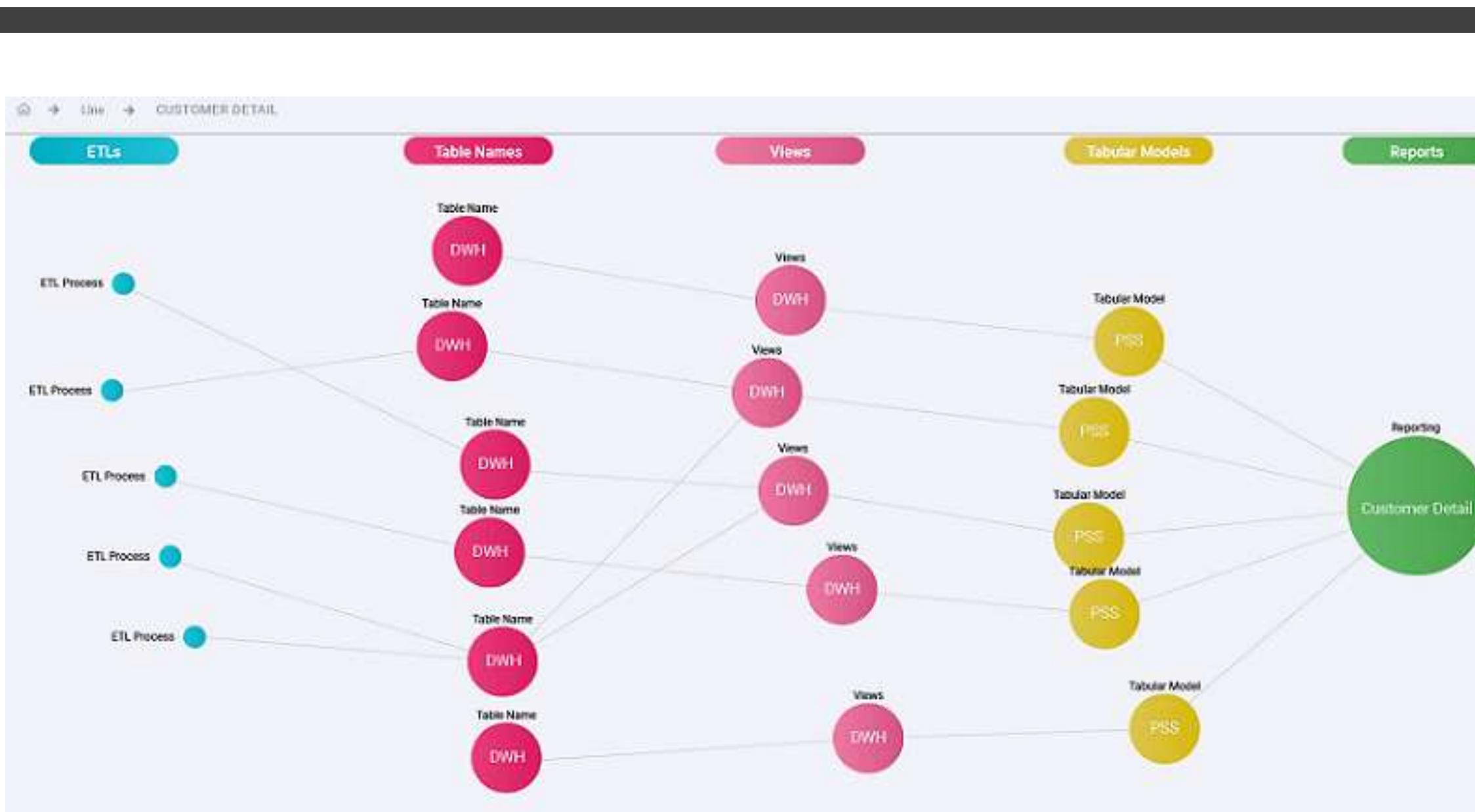
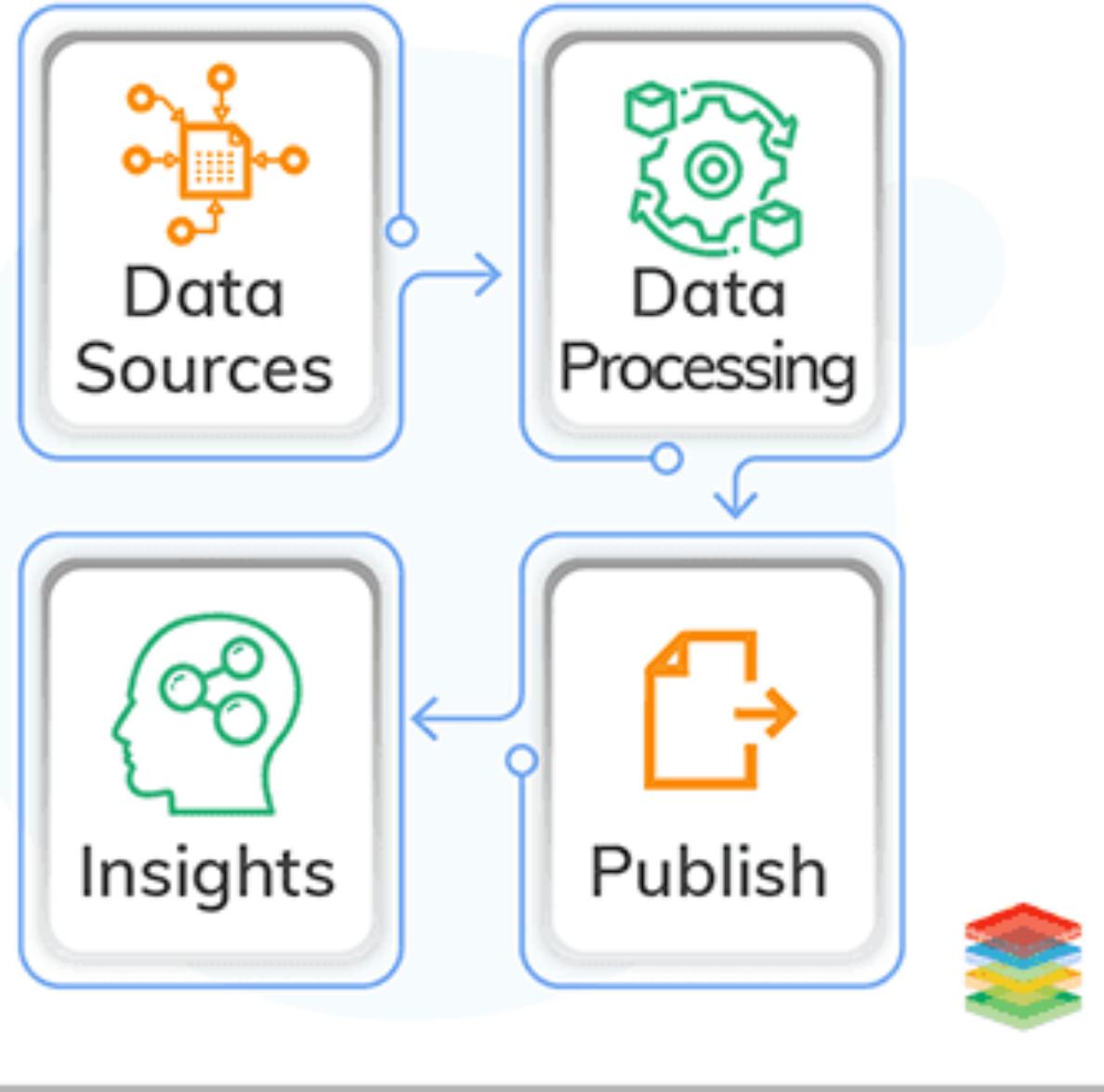
**DEMOCRATIZATION**  
*Democratize access to information so that business decisions are made based on global information.*

# Data Governance for Analytics

## Components Map

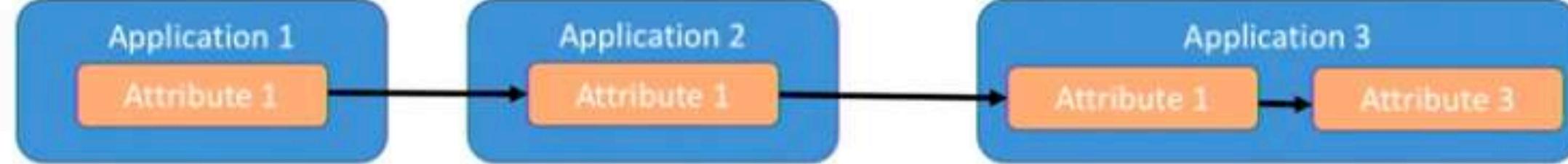


# Data Lineage



**What is Data Lineage?**

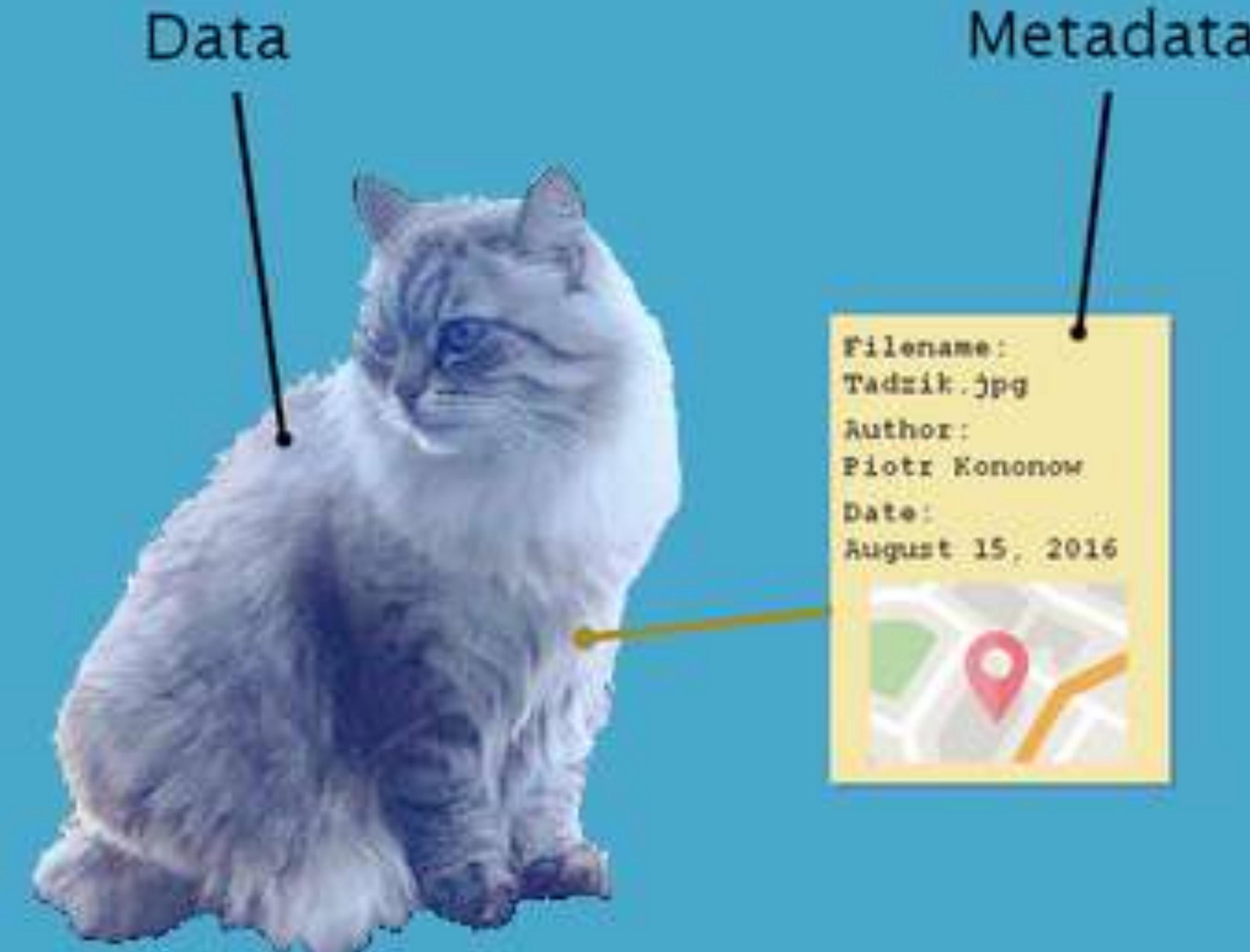
"Data lineage is generally defined as a kind of data life cycle that includes the data's origins and where it moves over time. This term can also describe what happens to data as it goes through diverse processes. Data lineage can help with efforts to analyze how information is used and to track key bits of information that serve a particular purpose." - Techopedia



# What is Metadata

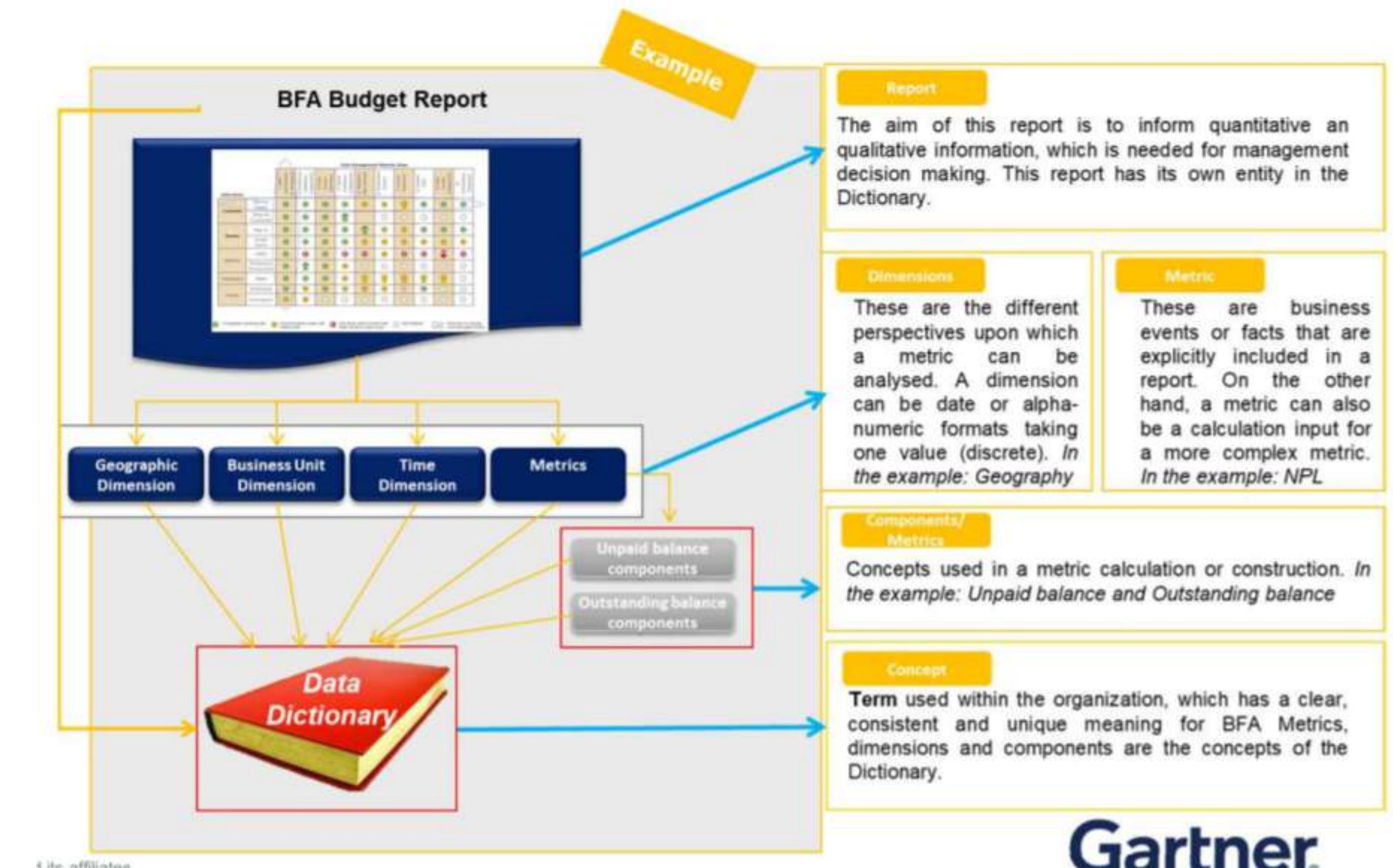
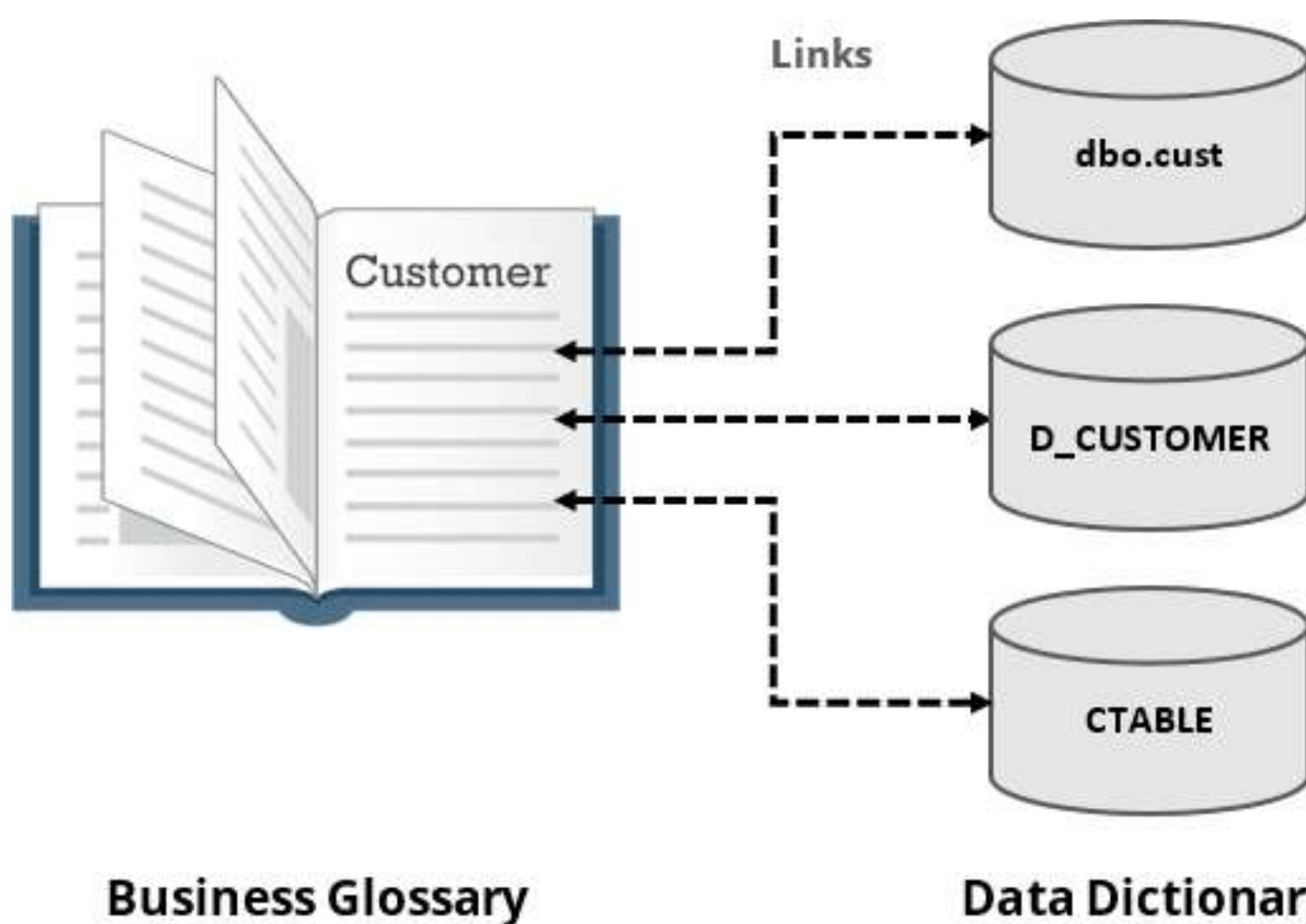
**Metadata is  
data about data**

- Information about what sort of things a piece of data contains
- Describes other data
- Used for indexing, discovery & search
  - For example, it may say an entry is a diagnosis, but does not say what the diagnosis is
  - Useful for governance
- It is not a summary or abstract
- Does not contain transaction data



# Data Catalog, Data Dictionary and Business Glossary

	Summary	Scope	Owner
Business Glossary	Business Term Definitions	Enterprise or Division-Wide	Business
Data Dictionary	Technical Descriptions	Database-Specific	Technology
Data Catalog	Location Directory	Enterprise or Division-Wide	Technology



**Gartner**



# Airbnb Dataportal

Search all Airbnb **data**

## Search

The screenshot shows a search results page for 'trips launch'. At the top, there are 'Top-level search filters' for All, Data sources, Charts, Knowledge Posts, Groups, and People. Below the search bar, there's an 'Advanced search' section with an 'Order by' dropdown set to 'Most relevant'. The results list includes:

- Trips Launch Metrics**: SUPERSET DASHBOARD. Shows metrics like 9.2% conversion rate, 275 users, and Tamar Eterman as author.
- Trips Web Launch Dashboard**: SUPERSET DASHBOARD. Shows metrics like 2.1% conversion rate, 2 users, and Daren Zhou as author.
- Trips Launch Detailed Search Metrics**: SUPERSET DASHBOARD. Shows metrics like 3.1% conversion rate, 3 users, and Gregory Dyer as author.
- Trips Web Launch Search Metrics**: SUPERSET DASHBOARD. Shows metrics like 0.1% conversion rate, 0 users, and Paul Liu as author.

Annotations point to the search filters, advanced search, and the individual resource details.

Top-level search filters

Resource details & metadata

Context, context, & context

## Context and metadata

The screenshot shows a detailed view of a data table for 'core\_data.dim\_listings'. The table has columns for ID, listing\_id, name, address, city, state, country, and other properties. Annotations point to various parts of the interface:

- Description, external link, social**: Points to the table header and a 'View on Airtable' button.
- Metadata & consumption**: Points to the 'Related Content' sidebar, which lists other dashboards and reports.
- Surface relationships, everything's a link to promote exploration**: Points to the overall layout where every element is a link.

core\_data.dim\_listings

A summary table for integrating Airbnb. All attributes and aggregates are reported in UTC time.

Related Posts

Created Nov 25, 2016

Last updated 1 hr ago

1 post

Recently commented 1 hr ago

Related Content

CEO Dashboard 2016-17

CEO Dashboard 2017-18

8 Dashboard

CEO Data Dashboard Test

## Employee-centric Data

The screenshot shows a user profile for John Bodley, a Software Engineer. It includes:

- User details & metadata**: Points to the profile picture and basic info.
- What they make, what they consume**: Points to a table showing various projects and their status.

What they make,  
what they consume

## Team-centric Data

The screenshot shows a group overview for the 'econ' team. It includes:

- Curated + Popular content**: Points to a large thumbnail for 'Over 100+ Dashboards'.
- Thumbnails for maximum context**: Points to the 'Find us at' section.
- Basic organization functionality**: Points to the 'Members' and 'Actions' sections.

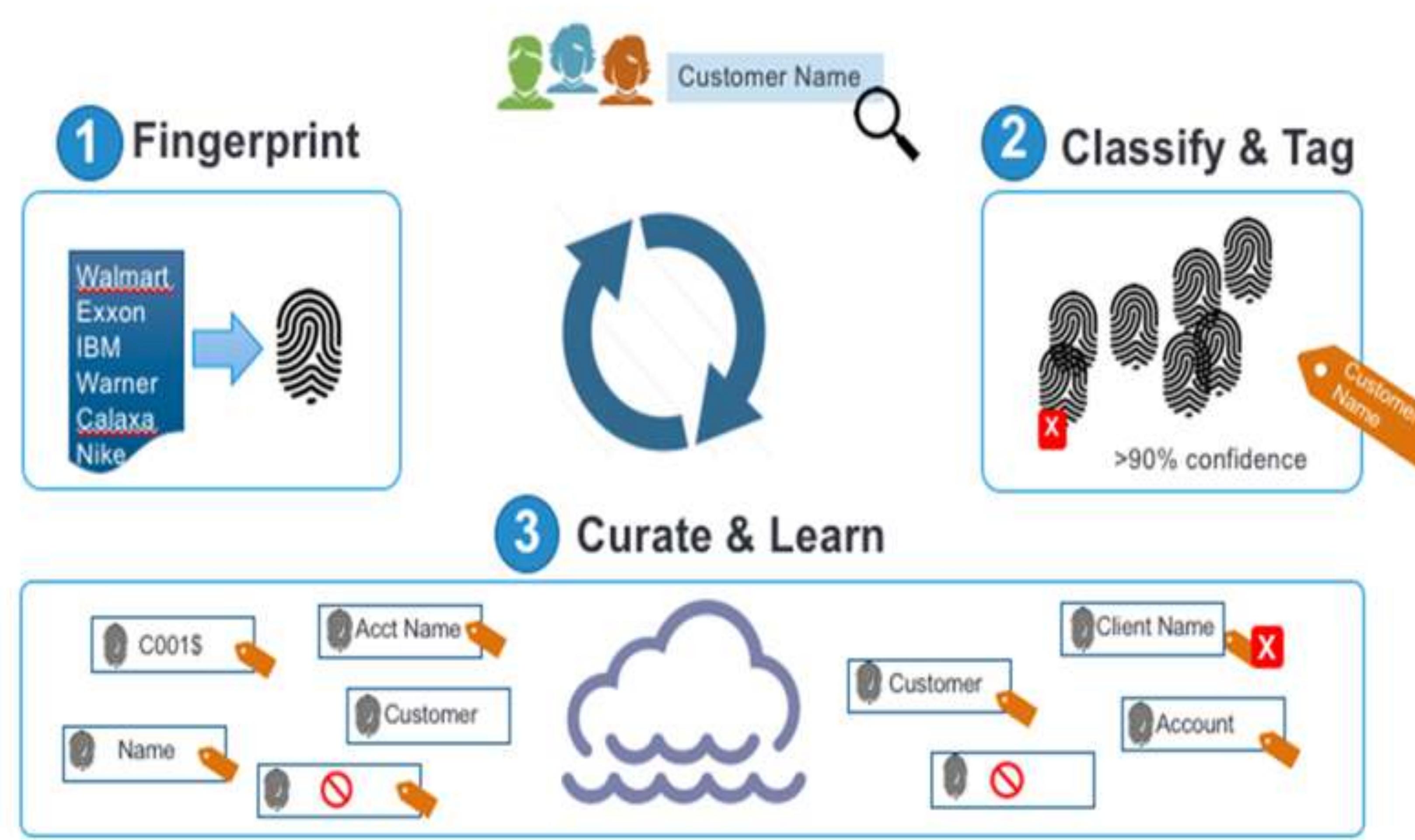
Curated + Popular content

Thumbnails for maximum context

Basic organization functionality

# Automated Data Detection

One of the most exciting parts of the market today is the availability of **automated data and metadata profiling tools**. They are AI/ML tools that can be used to help discover the structure and content of incoming data.



# Data Platform

¿What is the data management platform (**DMP**) used for collecting and managing data? ¿It allows to use big data and artificial intelligence algorithms to process and analyze large data sets? ¿Can it grow and acquire new functionalities as the data grows and the business needs evolve?



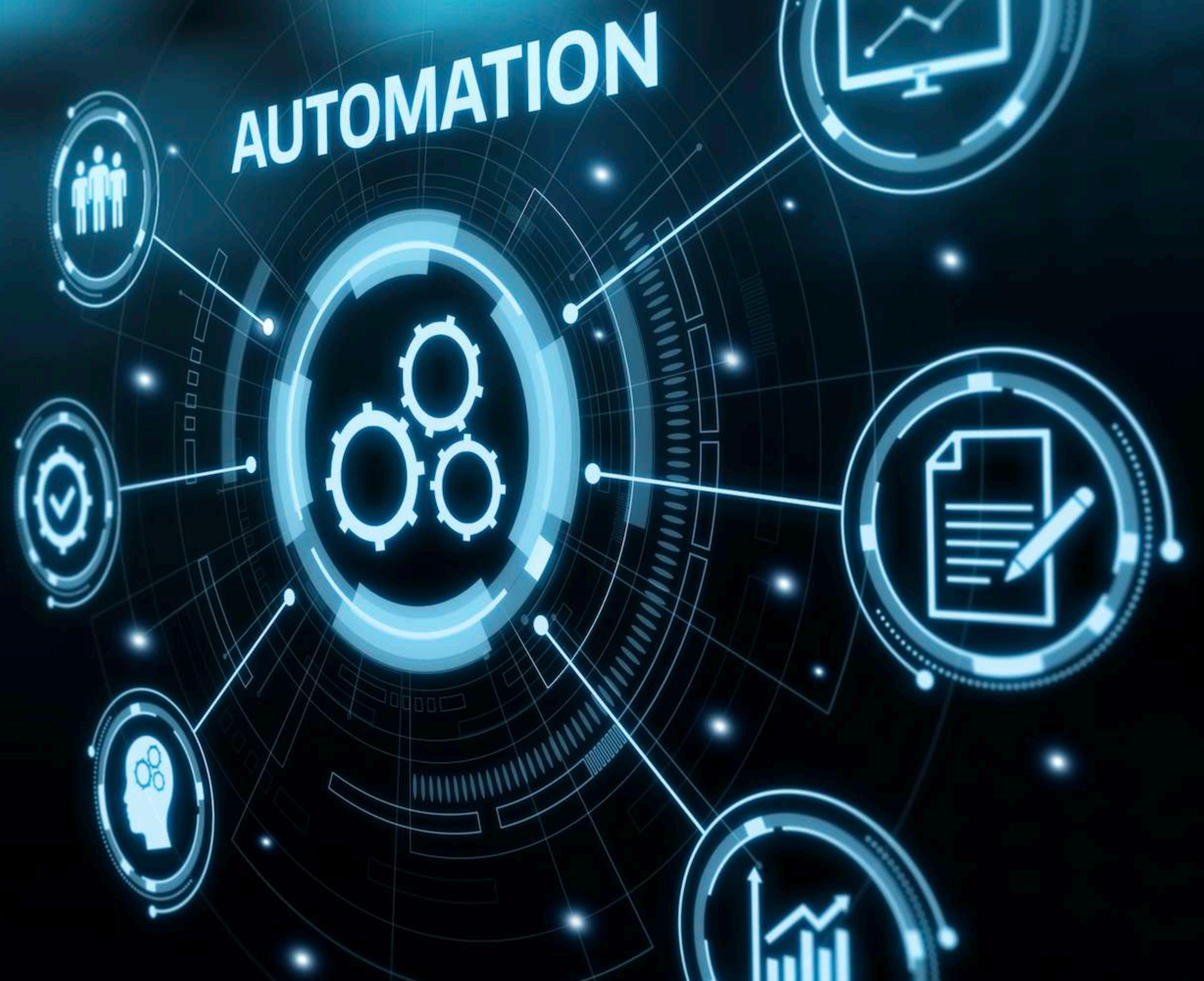
## Data Governance

¿What are the processes and policies in place that ensure that high **data quality** exists through its complete lifecycle? ¿What data **controls** there are that support business objectives? ¿It is the data **available**, **usable**, **consistent** and **secure**?



# 3rd Barrier: AI Ops

AUTOMATION



The **AI strategy** aims to obtain the maximum potential of this technology for our businesses by **solving all the barriers** to its implementation.

**1**

## Business Value

**Usage:** What are the main drivers of AI value in each business?

**AI ambition:** How much do we believe AI can transform our business?

**AI maturity:** How innovative in the use of AI does the company need to be?

**2**

## Data Access

**Types of sources:** What data will be needed at the company to develop AI?

**Data Platform:** What technological tools will we use to facilitate access to data at the company?

**Data Governance:** How are we going to ensure that they are of the right quality?

**3**

## AI Operations

**Organization:** Where should we locate the AI specialists at the company?

**Sourcing Model:** Who develops and maintains the AI models?

**Process industrialization:** How efficient and scalable is AI development, deployment and maintenance?

**4**

## AI Culture

**Skills training & hiring:** How can we train the company staff? How do we attract & retain talent?

**Communication:** how do we make it easier for everyone to understand what AI brings to the company?

**Government:** How do we guarantee the proper use of AI (security, privacy, ethics) at the company?

**5**

Are the legal/regulatory issues sufficiently resolved?

**6**

Is society ready to accept the transformation of the value proposition?

**7**

It may be legal, profitable and working - but is it ethical?



# Why do 87% of data science projects never make it into production?

Venturebeat.com

But if this is a universal understanding, that AI empirically provides a competitive edge, **why do only 13% of data science projects**, or just one out of every 10, actually **make it into production**?

“One of the biggest [reasons] is sometimes **people think**, all I need to do is **throw money at a problem or put a technology in**, and **success comes out the other end**, and **that just doesn’t happen**,” Chapo said. “And we’re not doing it because we **don’t have the right leadership support**, to make sure we **create the conditions for success**.“

The other key player in the whodunit is **data**, Leff adds, which is a double edged sword — it’s what makes all of these analytics and capabilities possible, but **most organizations are highly siloed**, with owners who are simply not collaborating and leaders who are not facilitating communication.

“I’ve had **data scientists** look me in the face and say **we could do that project**, but we **can’t get access to the data**,“ Leff says. “And I say, **your management allows that to go on?**“

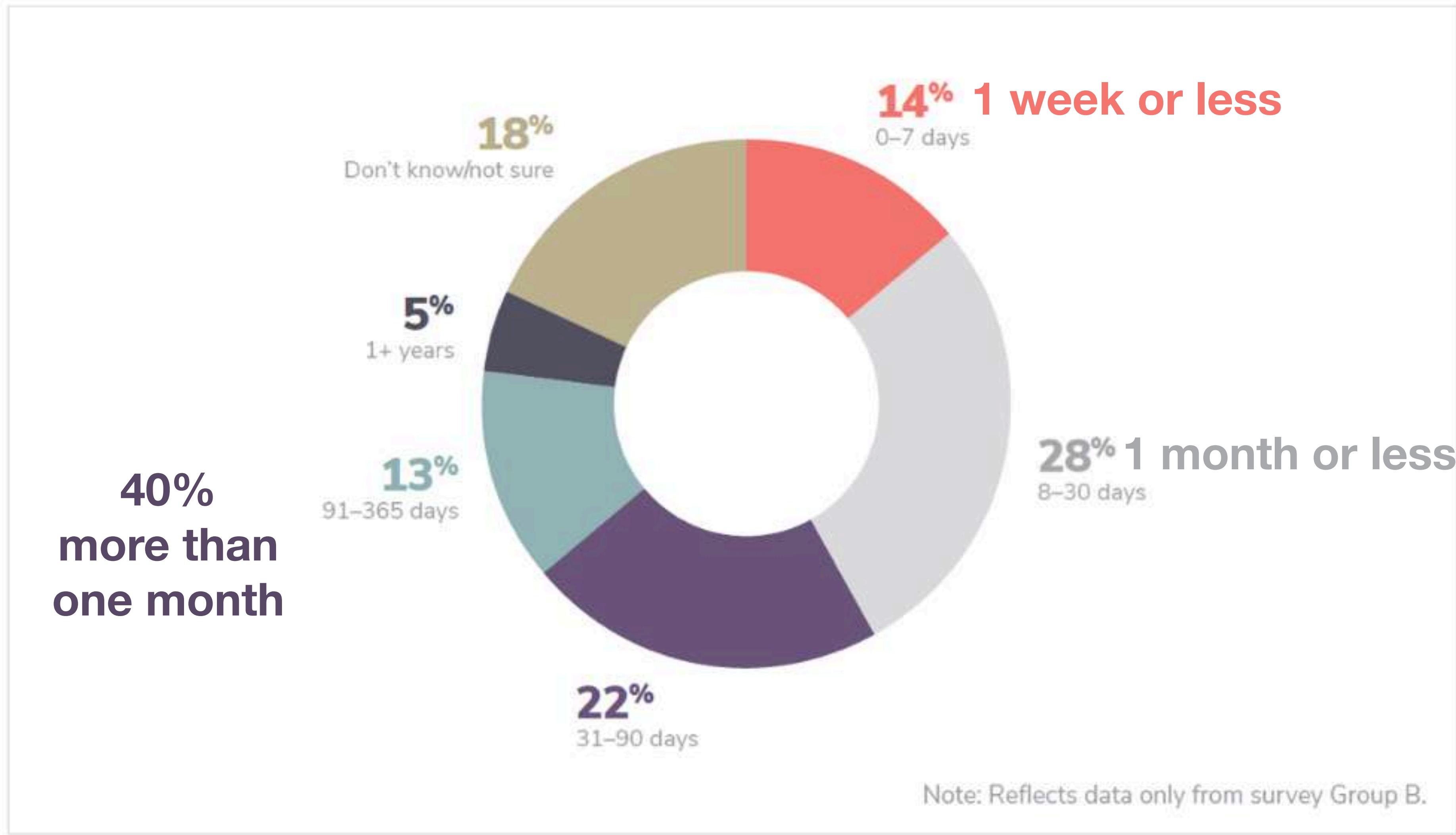
And the third issue, intimately connected to those silos, is the **lack of collaboration**. Data scientists have been around since the 1950s — and they were individuals sitting in a basement working behind a terminal. But **now that it’s a team sport**, and the importance of that work is now being embedded into the fabric of the company, it’s essential that every person on the team is able to collaborate with everyone else

# Add It Up: How Long Does a Machine Learning Deployment Take?

Lawrence E Hecht



## Machine learning model deployment timeline



*Source: Algorithmia's "2020 State of Enterprise ML". This question about how long it takes to deploy an ML model into production was only asked to a subset of respondents at a company that has an ML model production.*



**ginablaber**  
@ginablaber

Follow



The story of enterprise Machine Learning: “It took me 3 weeks to develop the model. It’s been >11 months, and it’s still not deployed.”

**@DineshNirmalIBM #StrataData #strataconf**

10:19 AM - 7 Mar 2018

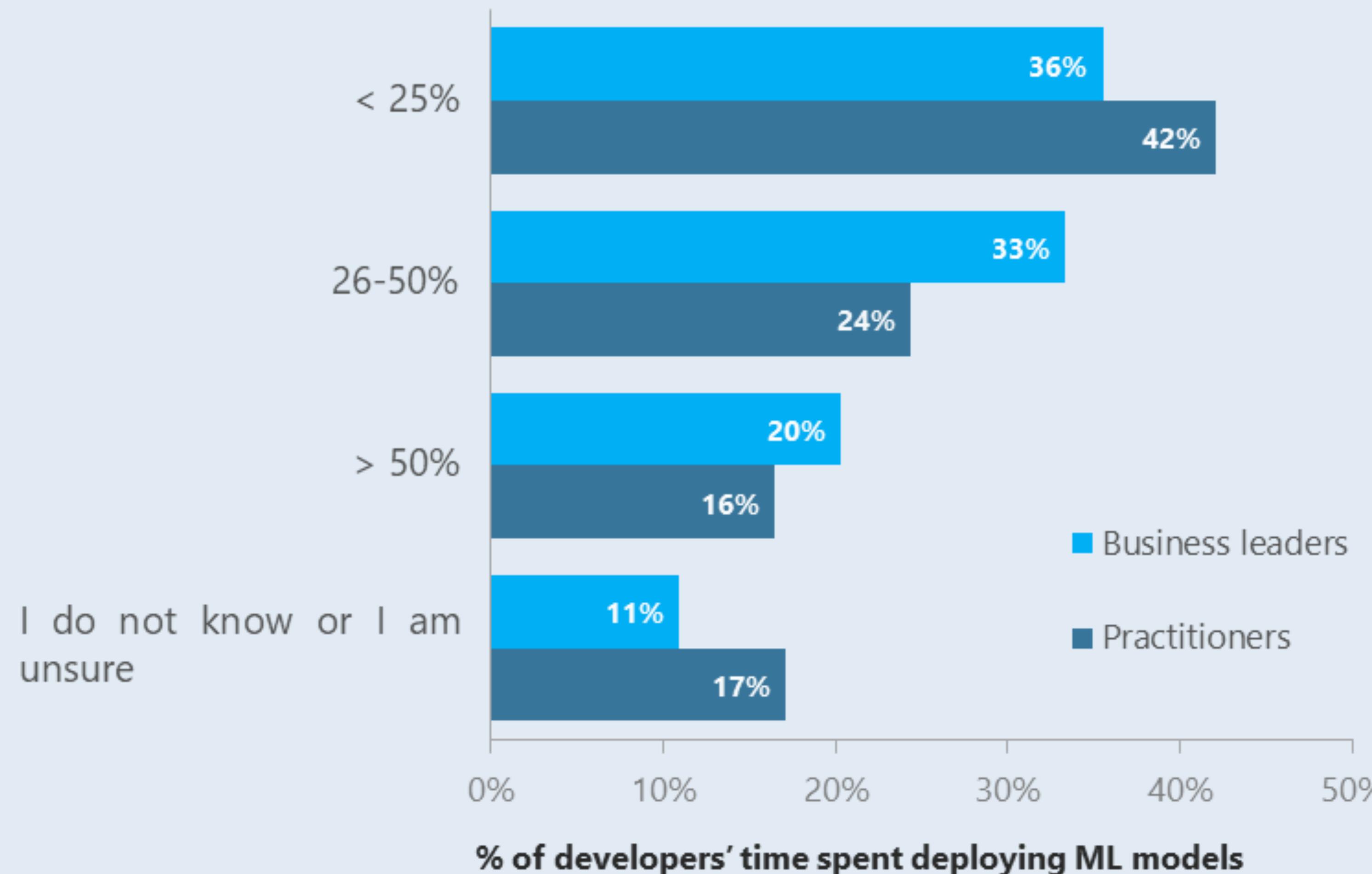
7 Retweets 19 Likes



7

19

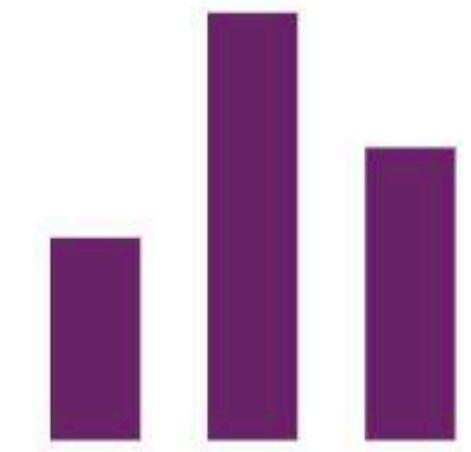
# 54% of Business Leaders Think Developers Spend >25% of Their Time Deploying ML Models



Source: Algorithmia's "2020 State of Enterprise Machine Learning". 745 respond people took the survey in October 2019. Q. What percentage of your developers' time is spent deploying ML models? Approximately 40% of were practitioners with job titles like data scientist and research engineer.

**THAT'S REALLY INTERESTING...**

**I HAVE NO IDEA WHAT YOU ARE TALKING ABOUT**

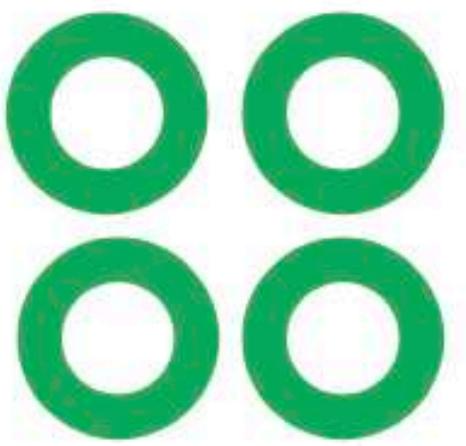


## Data

Schema

Sampling over Time

Volume



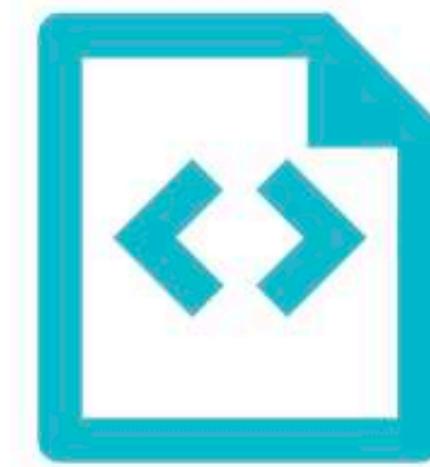
## Model

Algorithms

More Training

Experiments

Hyperparameters



## Code

Business Needs

Bug Fixes

Configuration

# Hidden Technical Debt in Machine Learning Systems

D. Sculley, Gary Holt, Daniel Golovin, Eugene Davydov, Todd Phillips  
{dsculley, gholt, dg, edavydov, toddphillips}@google.com  
Google, Inc.

Dietmar Ebner, Vinay Chaudhary, Michael Young, Jean-François Crespo, Dan Dennison  
{ebner, vchaudhary, mwyong, jfcrespo, dennison}@google.com  
Google, Inc.

## Abstract

Machine learning offers a fantastically powerful toolkit for building useful complex prediction systems quickly. This paper argues it is dangerous to think of these quick wins as coming for free. Using the software engineering framework of *technical debt*, we find it is common to incur massive ongoing maintenance costs in real-world ML systems. We explore several ML-specific risk factors to account for in system design. These include boundary erosion, entanglement, hidden feedback loops, undeclared consumers, data dependencies, configuration issues, changes in the external world, and a variety of system-level anti-patterns.

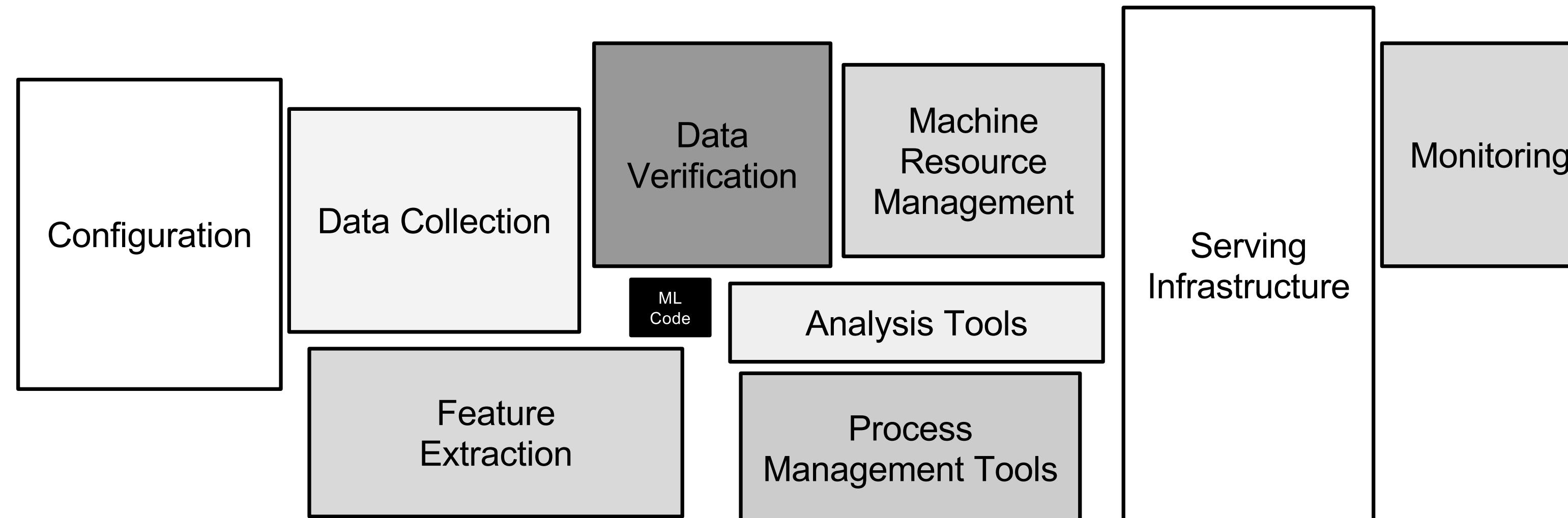
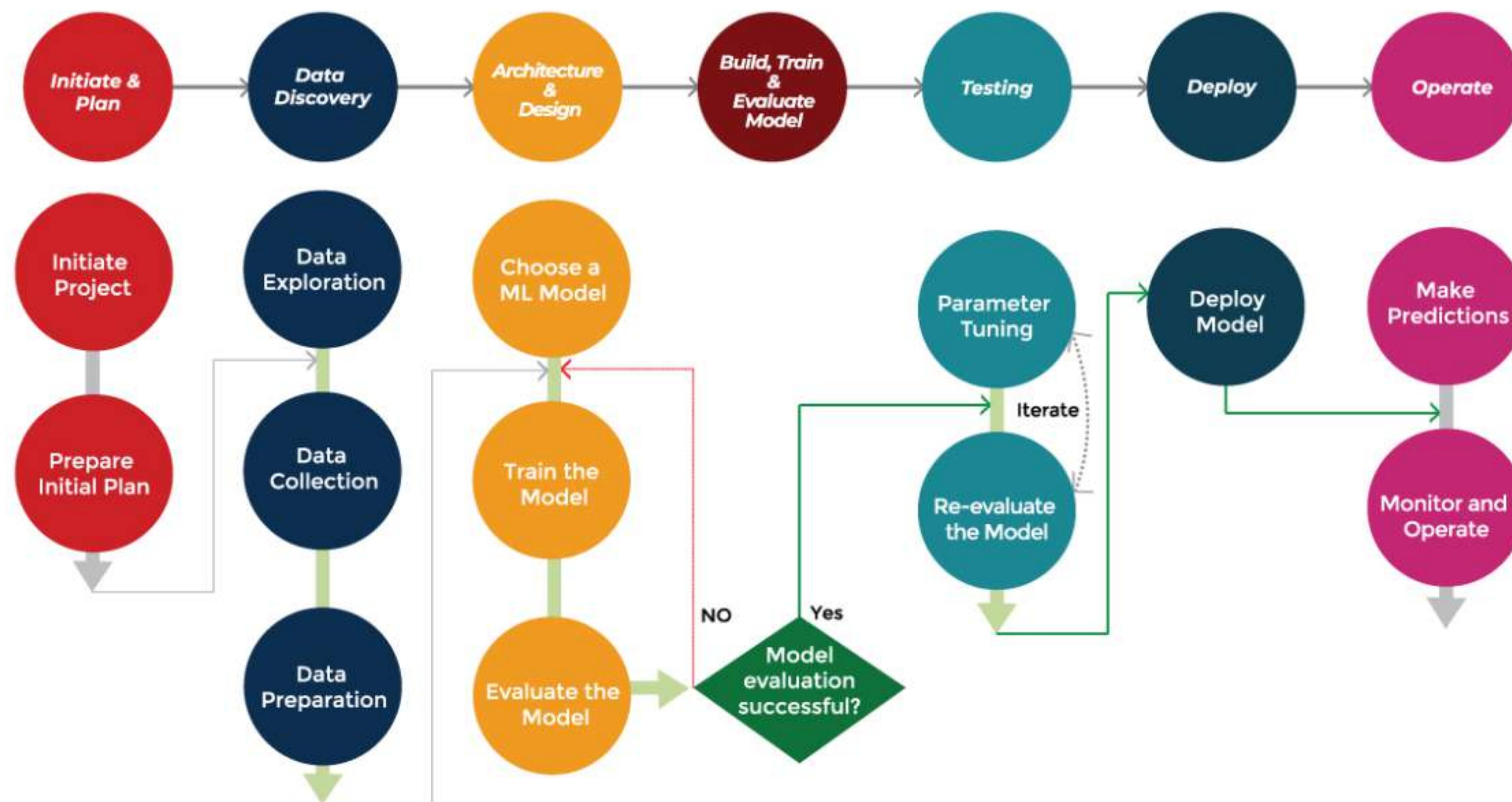


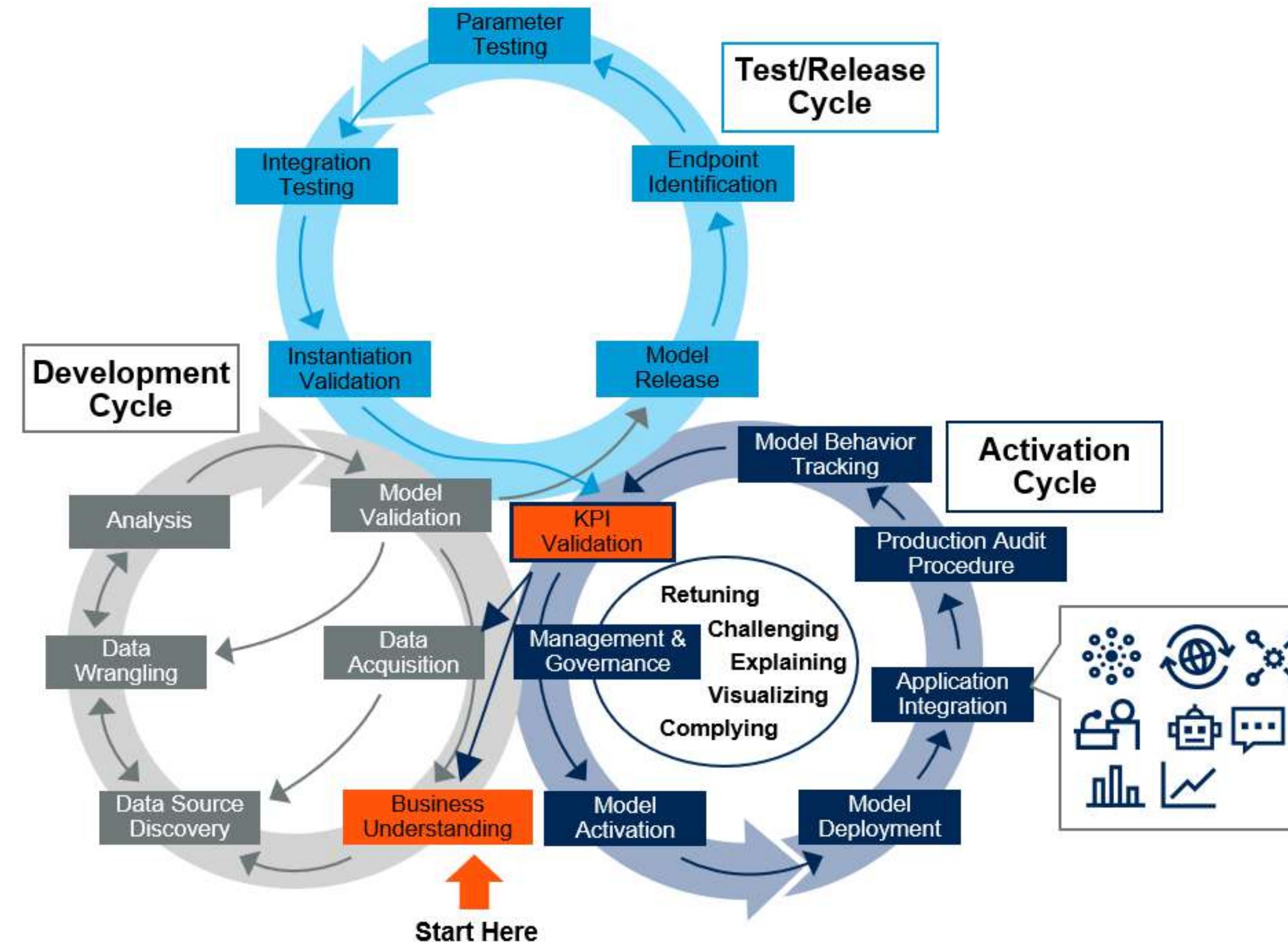
Figure 1: Only a small fraction of real-world ML systems is composed of the ML code, as shown by the small black box in the middle. The required surrounding infrastructure is vast and complex.

# MLOps Pipeline

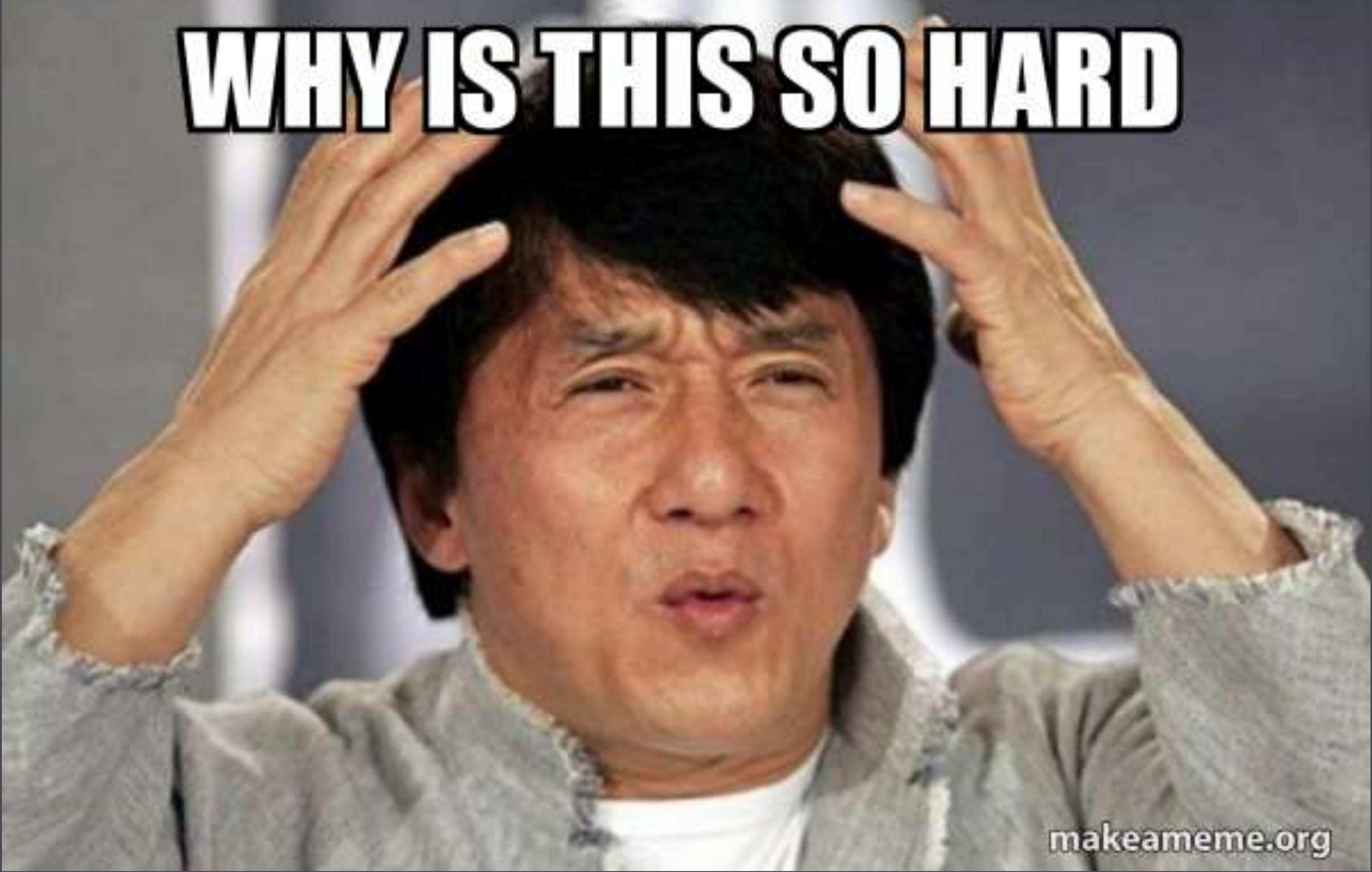


# MLOps Pipeline

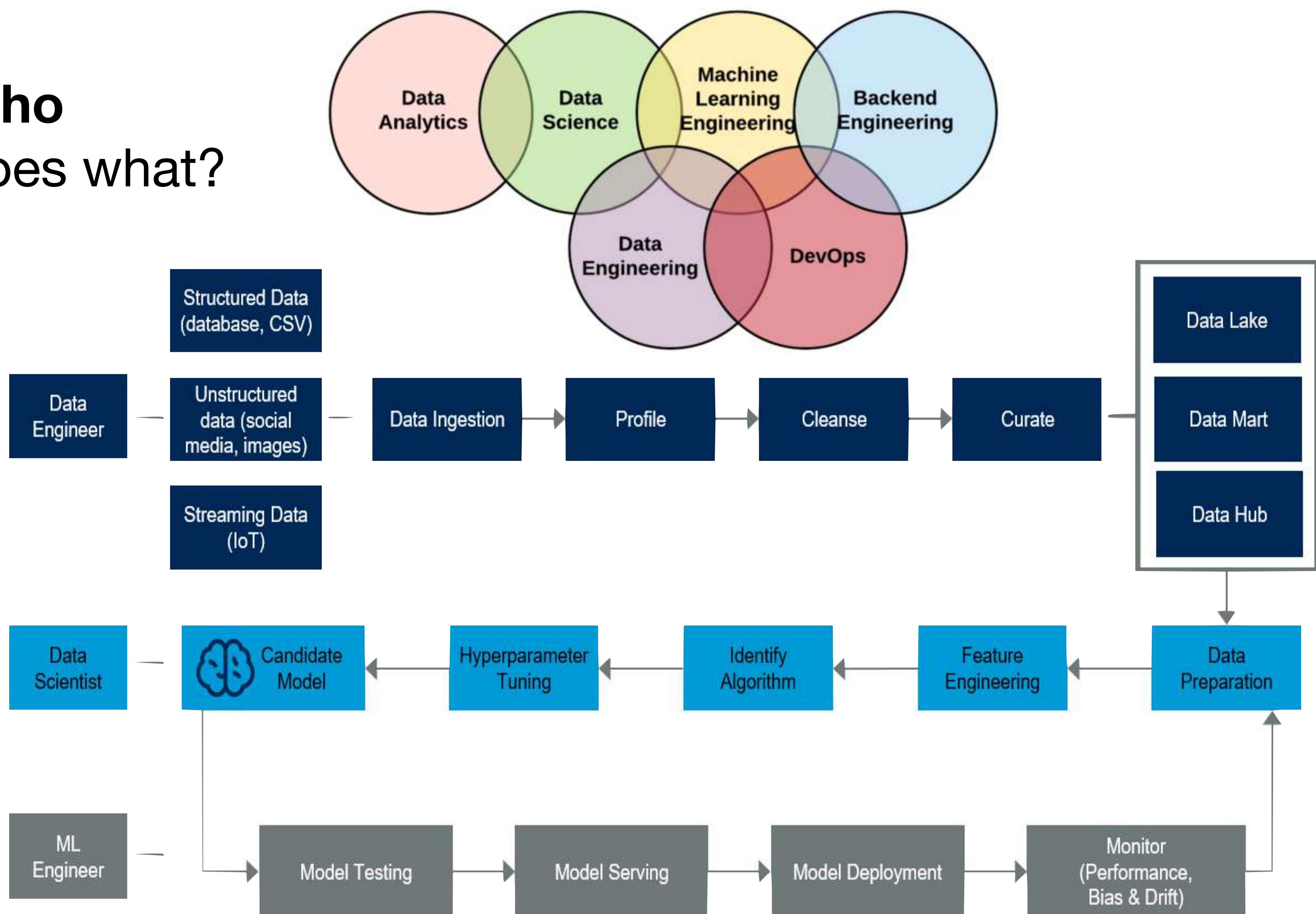
## Operationalizing ML Model Development for AI



# WHY IS THIS SO HARD



# Who does what?



Communication



Business  
Engagement

David G Pisano Retweeted



**Florian Leitner** @flowing

3 Jan 2019

Couldn't second this more; The hard part is finding engineers who can also build ML models. I believe pure data scientists who cannot engineer will only succeed in large companies, as they can more likely make use of employees who exclusively train models 24/7.

Emmanuel Ameisen @mlpowered

2018 has been a continuing flurry of exciting work in Machine Learning. If you are interested in being part of the field in 2019, I've written about how some of the most impactful trends of 2018 will impact this year! [medium.com/@emmanuelameisen/pic.twitter.com/6aLQ63XNEs](https://medium.com/@emmanuelameisen/pic.twitter.com/6aLQ63XNEs)



1

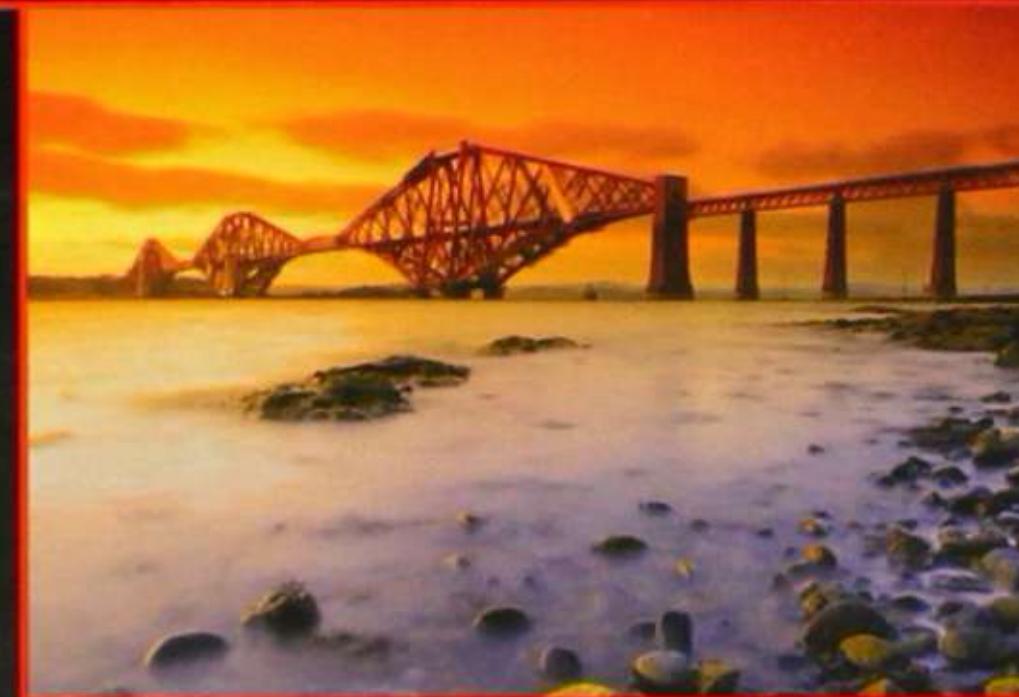
2

*The Addison-Wesley Signature Series*

# CONTINUOUS DELIVERY

RELIABLE SOFTWARE RELEASES THROUGH BUILD,  
TEST, AND DEPLOYMENT AUTOMATION

JEZ HUMBLE  
DAVID FARLEY



*Foreword by Martin Fowler*

MARTIN FOWLER SIGNATURE  
Book Martin Fowler

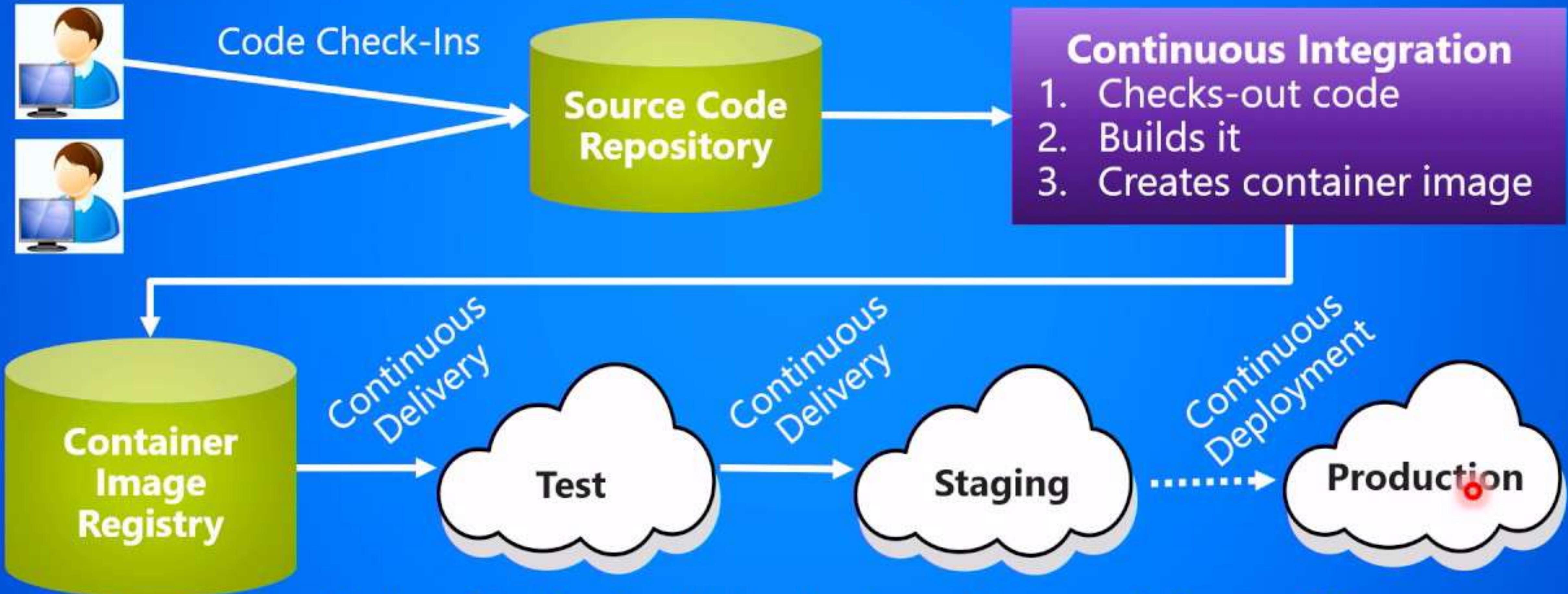


**“Continuous Delivery is the *ability* to get changes of all types — including new features, configuration changes, bug fixes, and experiments — *into production*, or *into the hands of users*, *safely and quickly* in a *sustainable way*”**

Jeff Humble and Dave Farley

# CI: Continuous Integration

# CD: Continuous Delivery, & Deployment

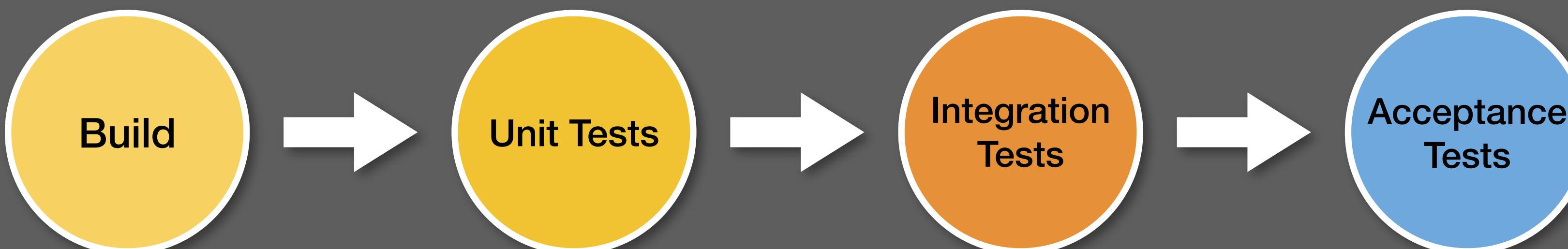


*Modern DevOps is all about automation; any failures are sent back to developers & stops further progress*

## Continuous Integration



## Continuous Delivery



## Continuous Deployment



# DATA Engineer

# The GAP

# DATA Scientist



WTF  
hyperparameters?

WTF pull  
request?



?

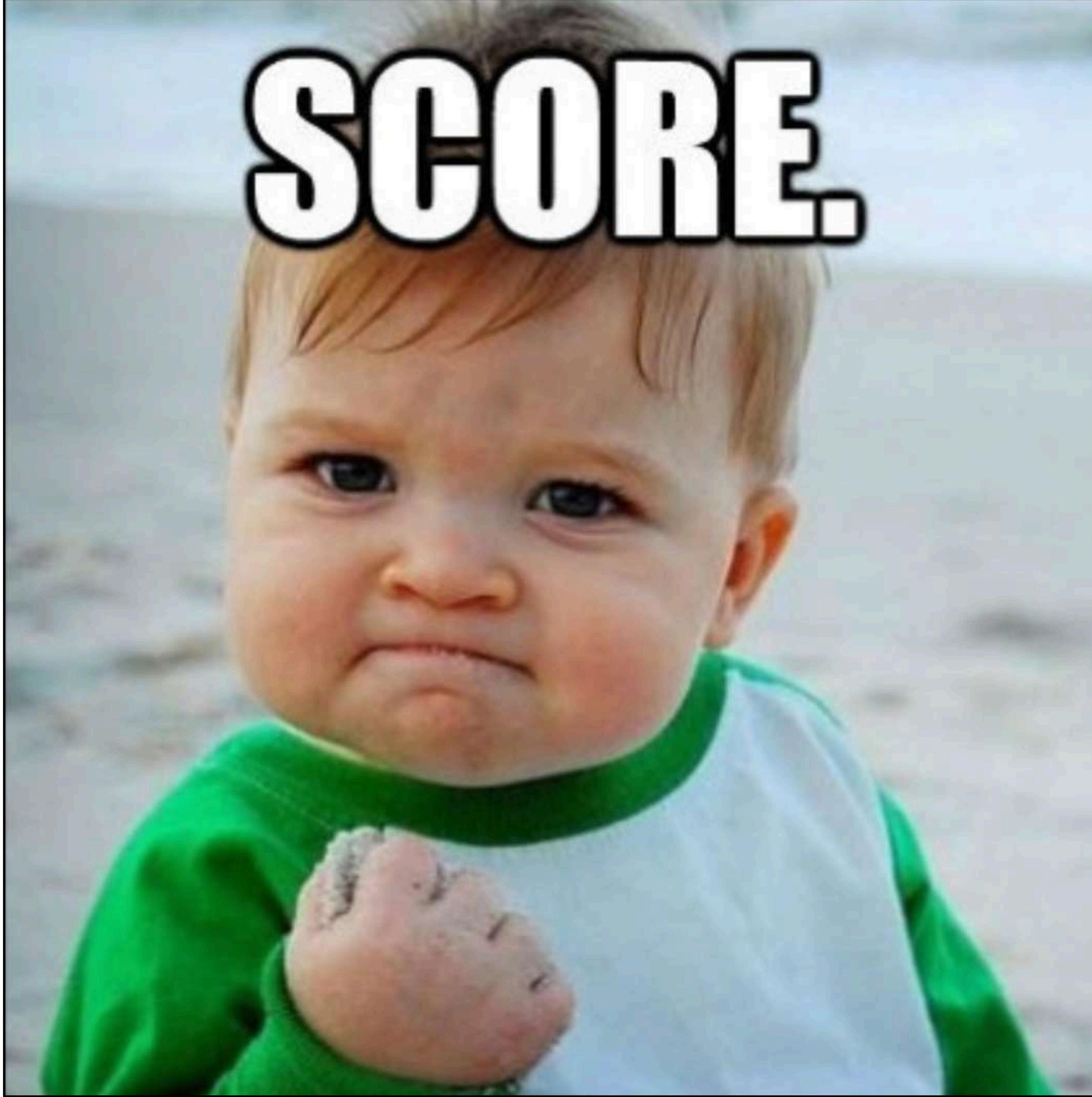


# Data Scientists just want to Data Science

The image shows a complex data science workspace with several open windows:

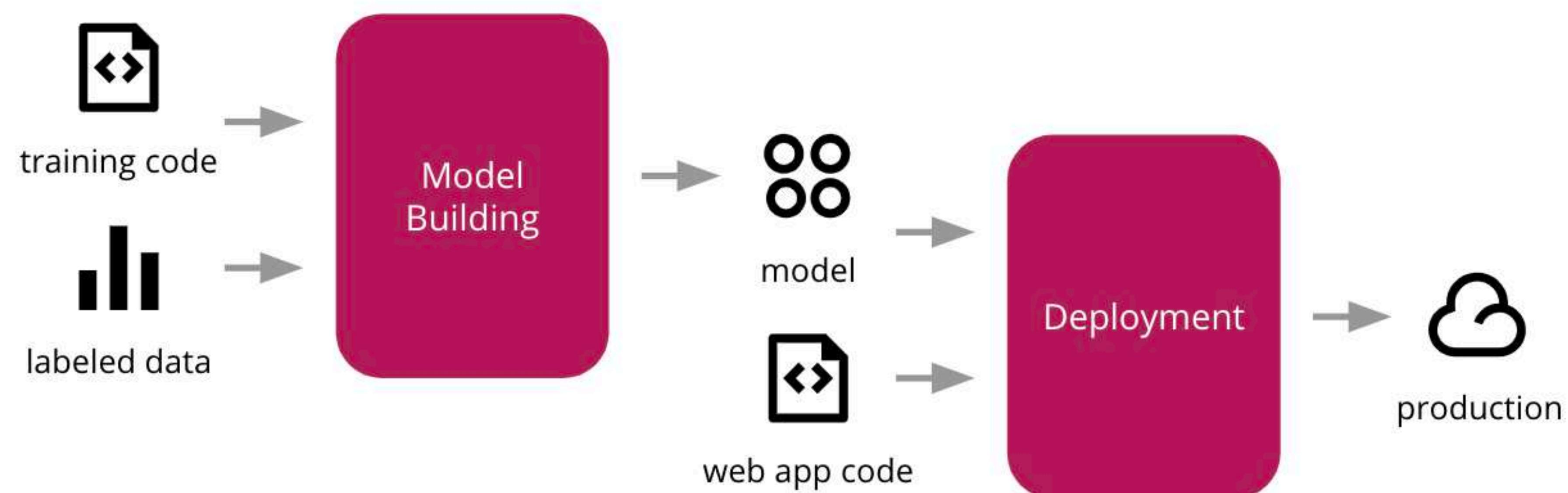
- Launcher:** Shows icons for Python 3, C++11, C++14, C++17, Julia 1.1.0, phylogenetics (Python 3.7), R, and Databricks 3.
- Output View:** Displays a scatter plot titled "Seattle Weather: 2012-2015" showing Maximum Daily Temperature (C) over time, with a bar chart below it.
- Altair.ipynb:** A Jupyter notebook titled "In Depth: Linear Regression". It contains a slide titled "Simple" with a scatter plot of x vs y, and code cells for generating the plot and printing the kernel spec.
- R.ipynb:** A Jupyter notebook titled "R". It shows a scatter plot of Sepal.Length vs Sepal.Width for the iris dataset.
- Lorenz.ipynb:** A Jupyter notebook titled "Julia". It displays a scatter plot of Sepal.Length vs Sepal.Width for the iris dataset and some Julia code.
- python notebook:** A Jupyter notebook titled "python notebook". It shows a Lorenz system differential equations and some ipywidgets code.
- Julia:** A Jupyter notebook titled "Julia". It shows a scatter plot of Sepal.Length vs Sepal.Width for the iris dataset and some Julia code.
- Linear Regression.ipynb:** A Jupyter notebook titled "In Depth: Linear Regression". It contains a slide titled "Simple", code for generating a scatter plot, and a code cell for printing the kernel spec.
- Python 3 | Idle:** A Python 3 idle window showing a scatter plot of x vs y.

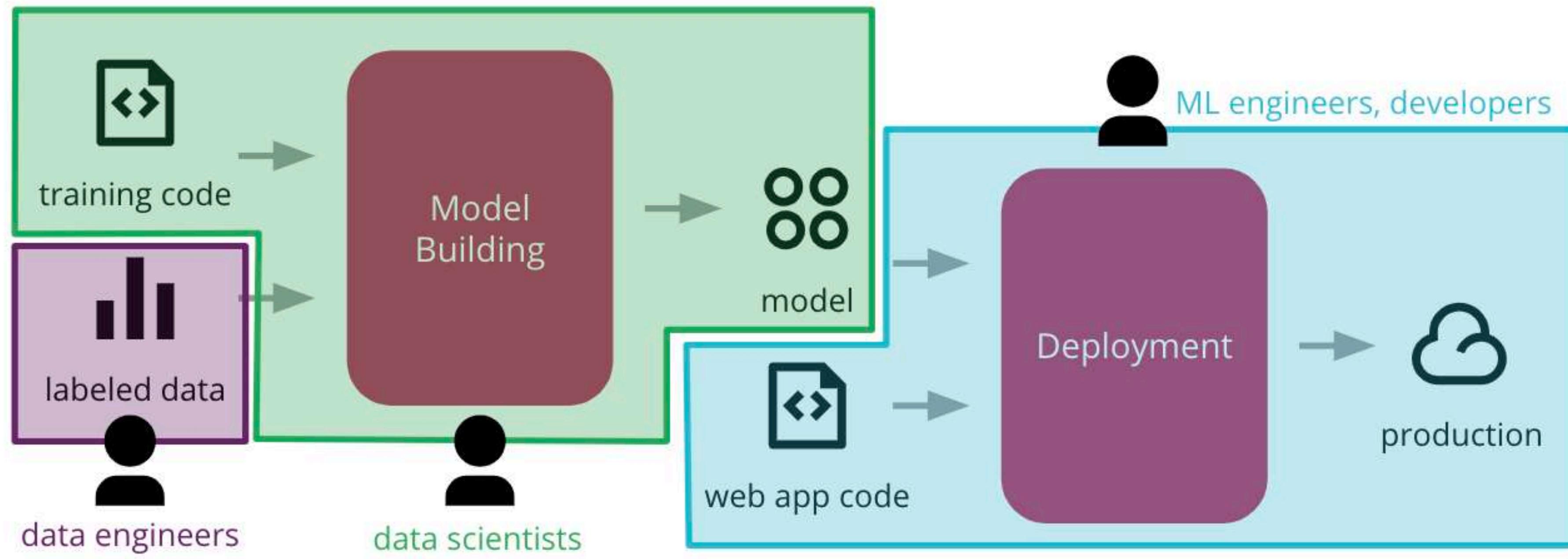
Typical data scientist work environment

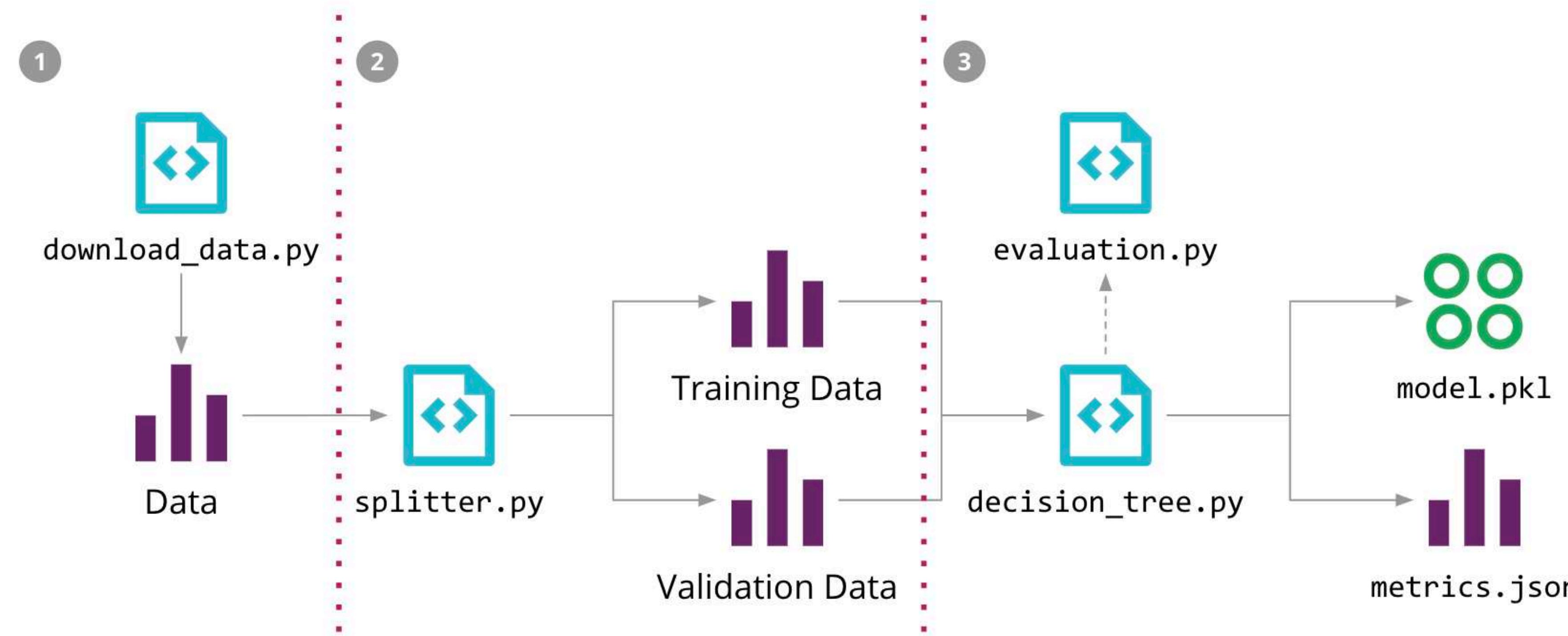


**SCORE.**

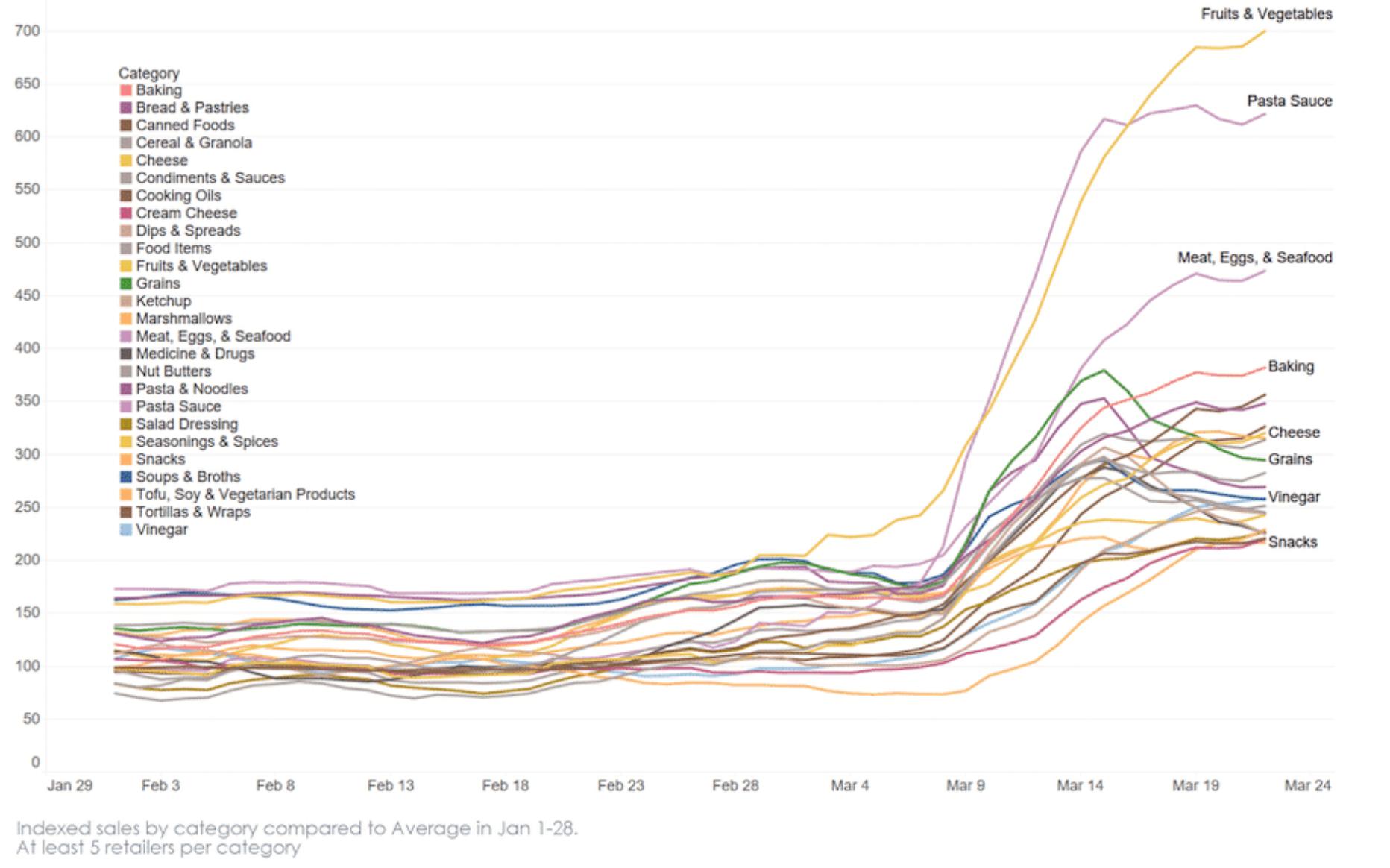
We've got  
the notebook  
into source  
control!



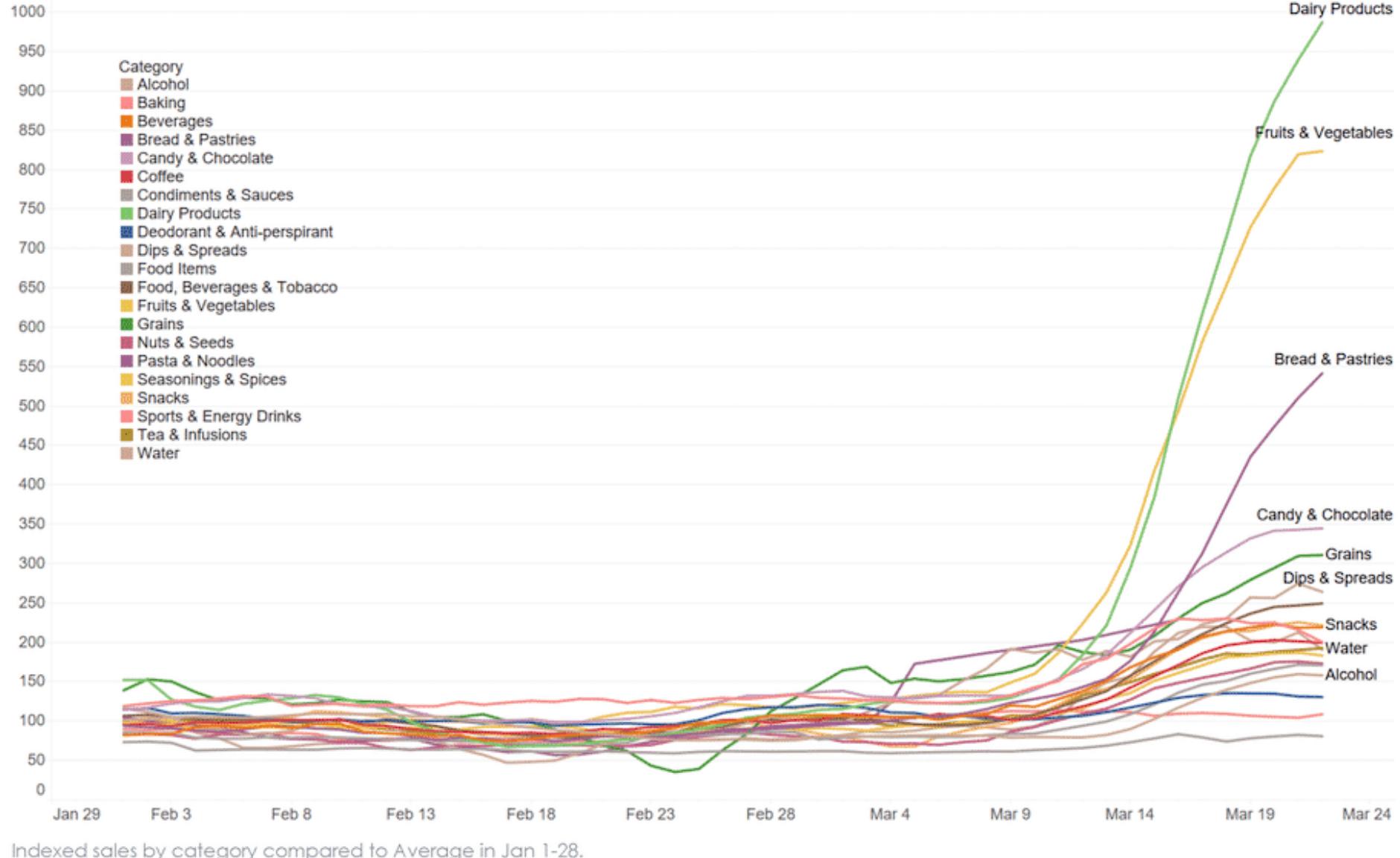




## Consumers Now Using Online Providers For All Grocery Needs



## Sales of Dairy Products and Fruits & Vegetables Rise in March



# Why is necessary to monitor and retrain?

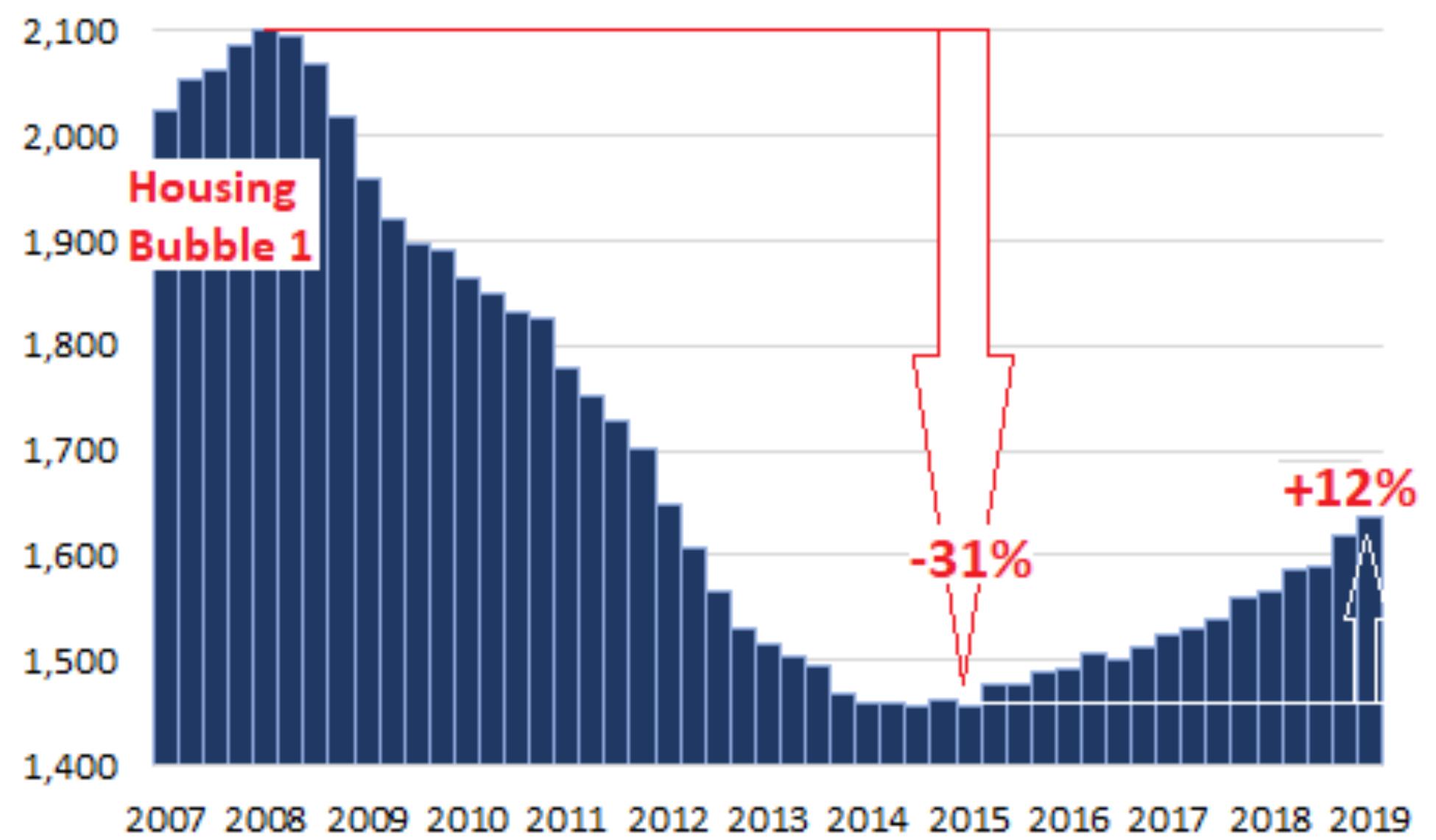
↑ Posted by [REDACTED] 3 years ago ↗  
**12.7k** Facebook's list of "suggested friends" is quite literally a list of people I've been avoiding my entire life.

523 Comments Share Save Hide Report 95% Upvoted

criteo.



## Spain Median Home Price, €/Sq. Meter Quarterly



criteo.

WOLFSTREET.com

A close-up photograph of a green frog's head and upper body. The frog has bright green skin with some yellowish-green patches on its front legs. Its large, dark eyes are looking slightly to the left. Water droplets are visible on the surface of the frog's skin and the background, suggesting a rainy environment.

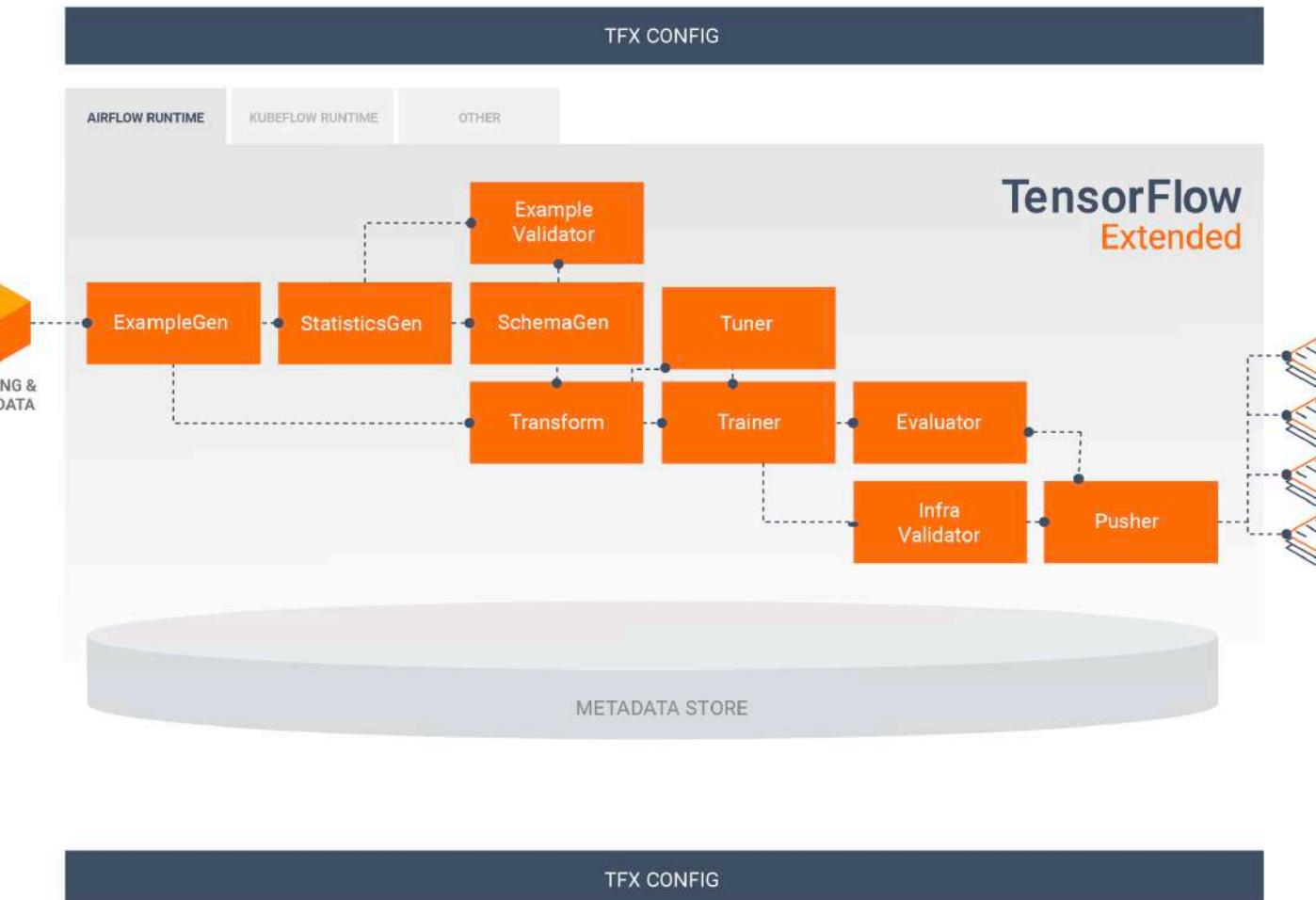
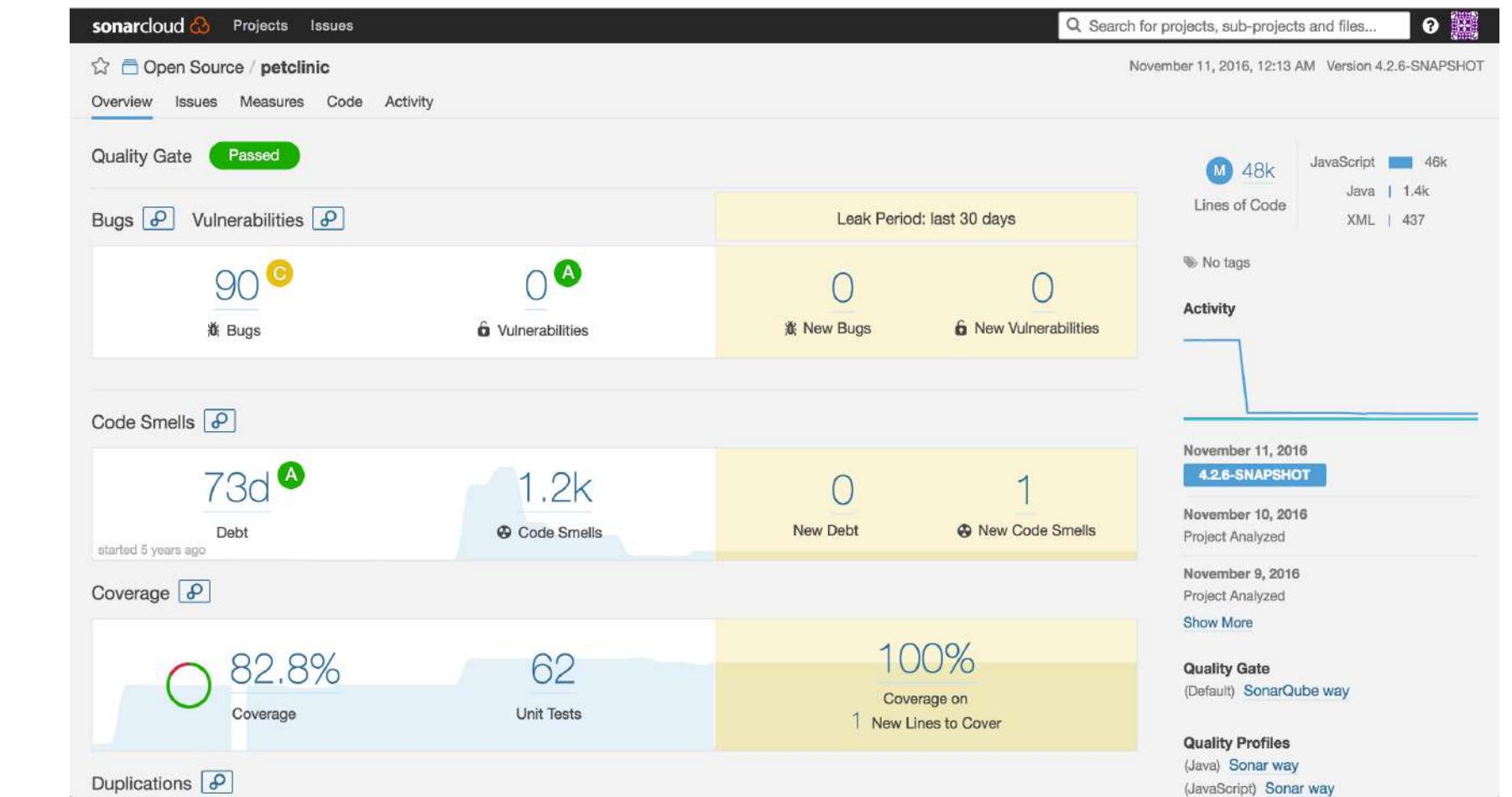
**I HAVE TO ACCEPT THAT**

**I'M FOREVER ALONE**

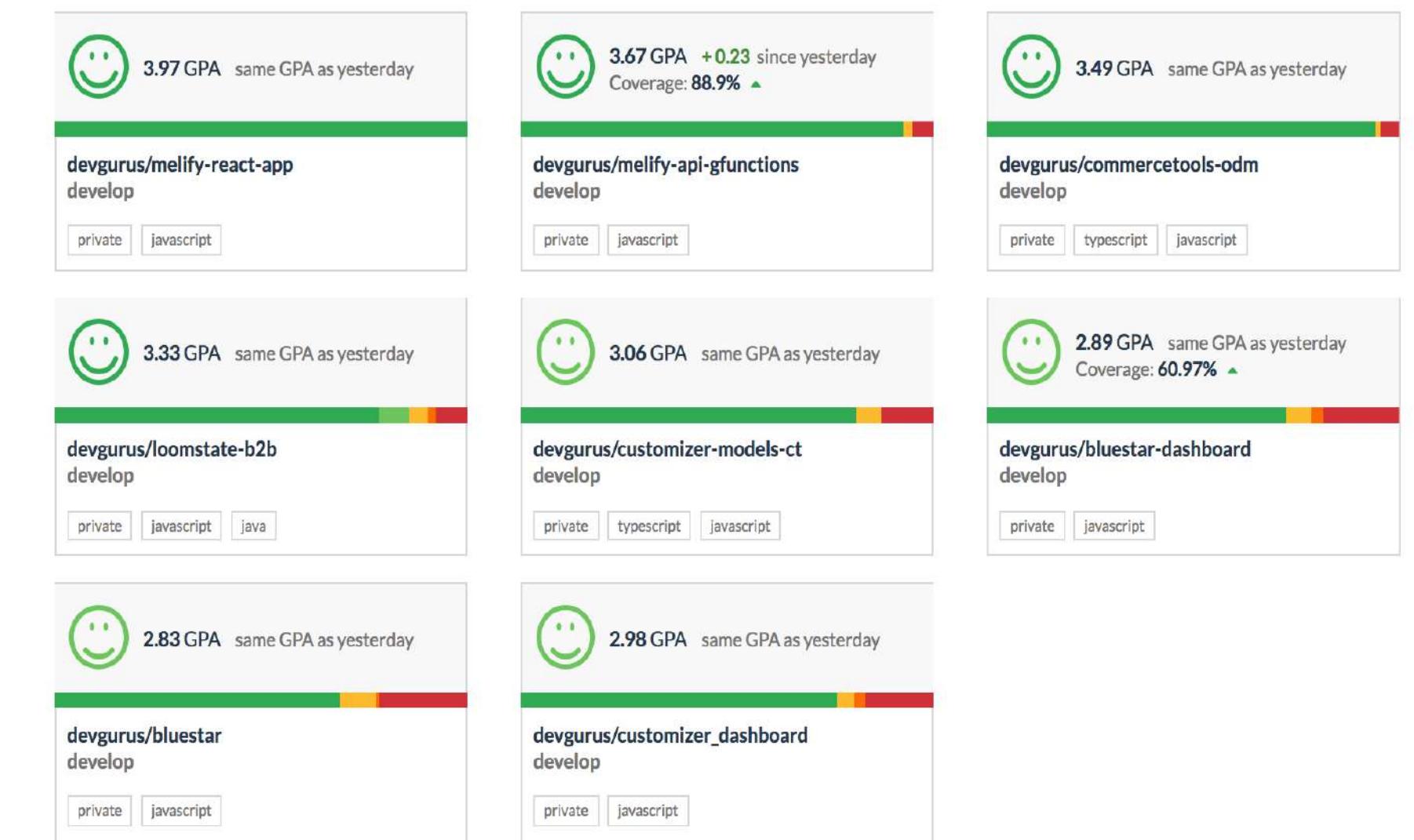
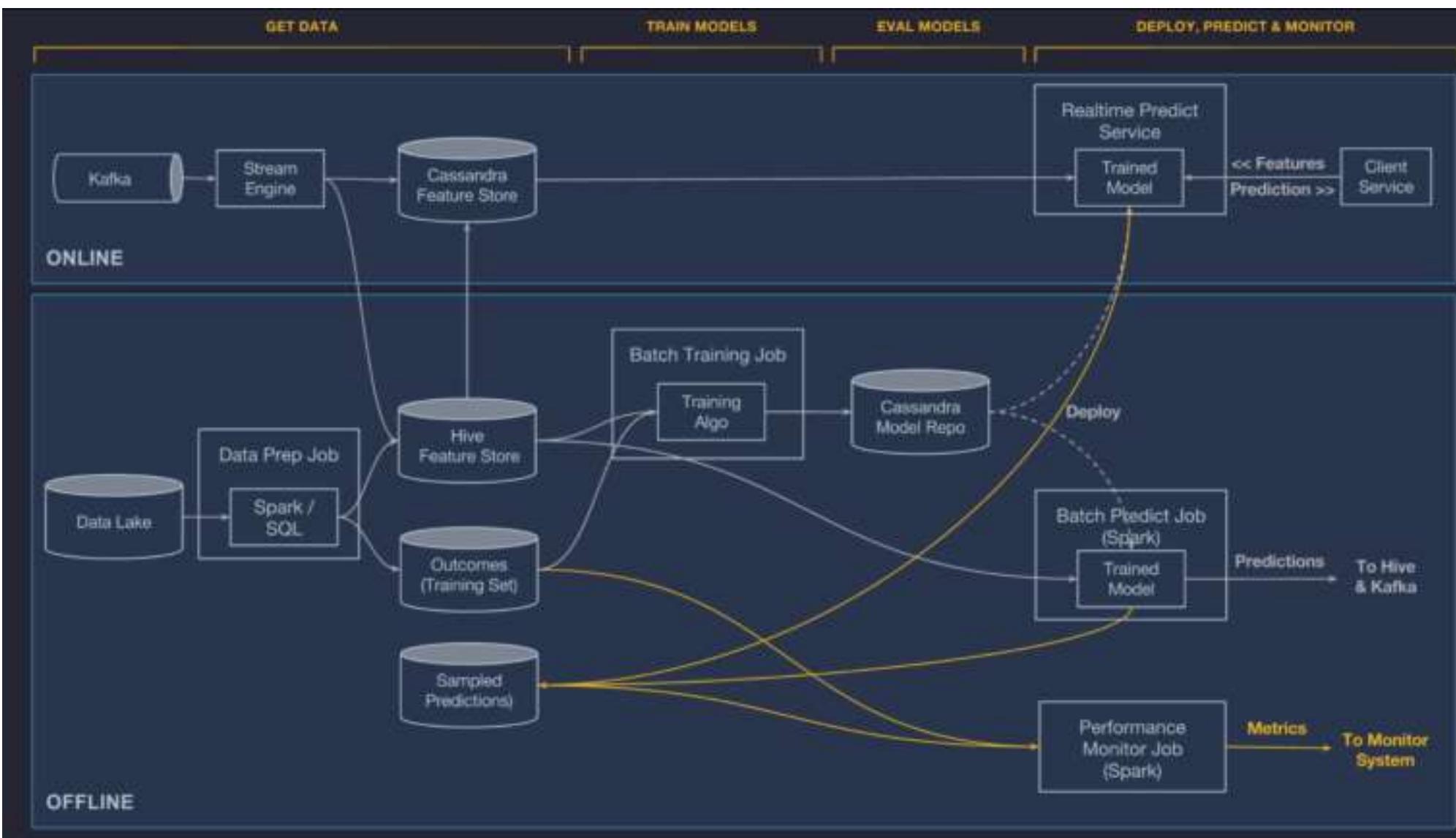
But I don't  
work at a big  
company with  
thousand of  
ML engineers!



**sonarqube**

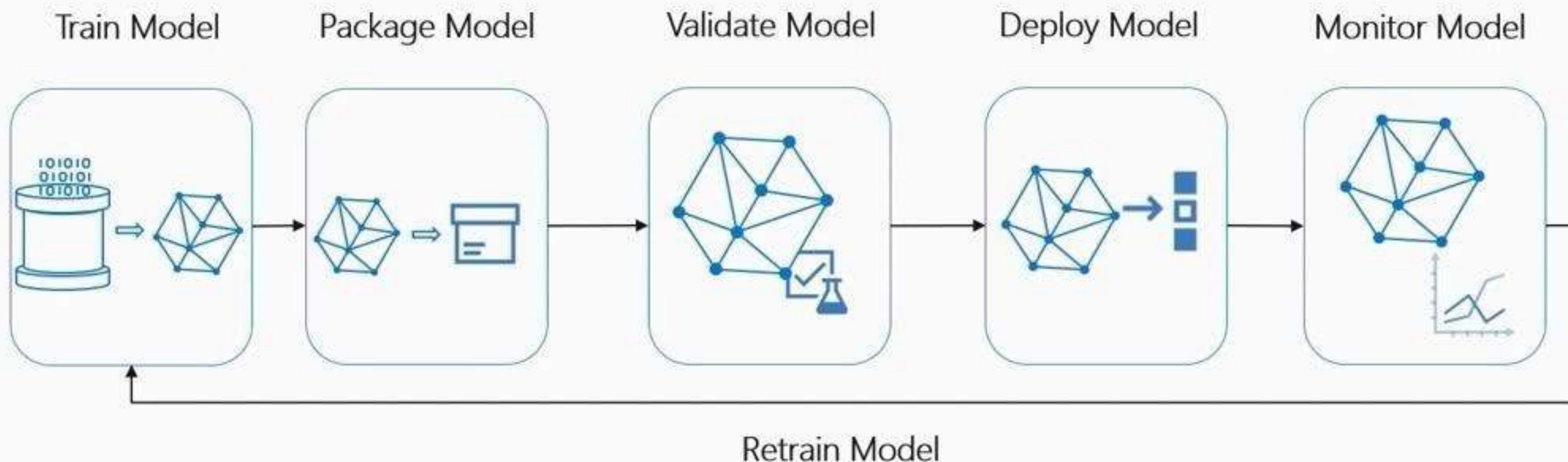


## Michelangelo: Uber's Machine Learning Platform



# What is the E2E ML lifecycle?

- **Develop & train model** with reusable ML pipelines
- **Package model** using containers to capture runtime dependencies for inference
- **Validate model behavior** – functionally, in terms of responsiveness, in terms of regulatory compliance
- **Deploy model** - to cloud & edge, for use in real-time / streaming / batch processing
- **Monitor model** behavior & business value, know **when to replace / deprecate a stale model**



# MLOps Benefits

## 1. Reproducibility + Auditability

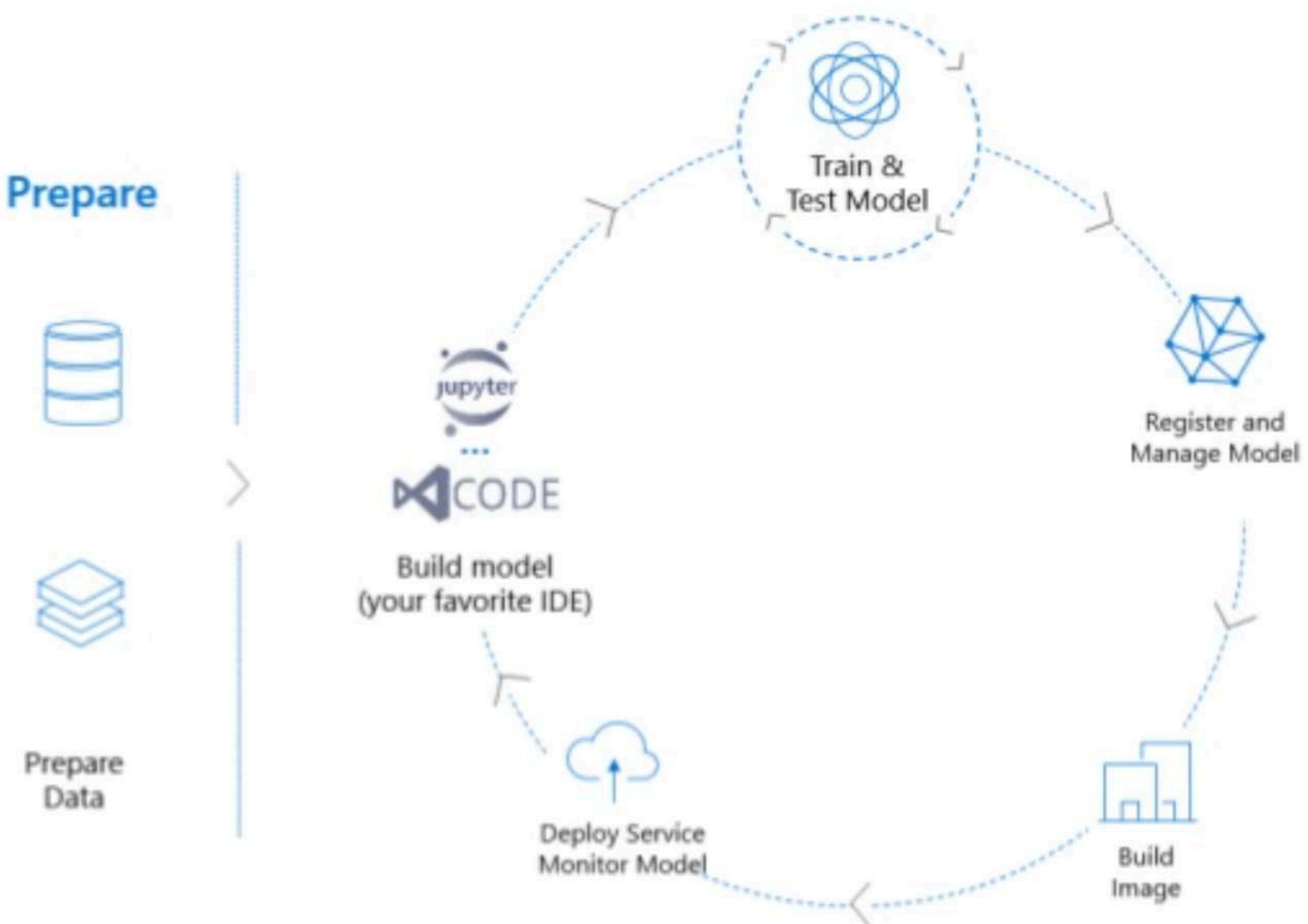
- Code drives generation and deployments
- Pipelines are reproducible and verifiable
- All artifacts can be tagged and audited

## 2. Validation

- SWE best practices for quality control
- Offline comparisons of model quality
- Minimize bias and enable explainability

## 3. Automation + Observability

- Controlled rollout capabilities.
- Live comparison of predicted vs. expected performance.
- Results fed back to watch for drift and improve model.



# What's next?

## Asynchronous Session 12

**TRADITIONAL CORPORATION BUSINESS CASE**  
*by Juanjo Casado*



dgonzalezp@faculty.ie.edu