



The background of the slide features a complex, abstract digital graphic. It consists of a dark, textured surface with numerous glowing blue and purple dots of varying sizes. These dots are interconnected by a network of thin, light-colored lines, creating a sense of a data grid or a neural network. In the center-left, there is a dense cluster of binary code digits ('0's and '1's) in green and yellow. To the right of this cluster, the words 'Build.', 'Unify.', and 'Scale.' are stacked vertically in large, bold, white sans-serif font. The 'U' in 'Unify.' is colored purple, matching the logo's color scheme.

Build.  
Unify.  
Scale.

WIFI SSID:Spark+AISummit | Password: UnifiedDataAnalytics



SPARK+AI  
SUMMIT 2019

# AI-Powered Retail Experience with Databricks

Akhil Dhingra, Zalando  
Saurav Verma, Zalando

#UnifiedDataAnalytics #SparkAISummit

# Zalando SE

- Founded in 2008 in Berlin.
- Europe's leading online fashion platform
- Connects customers, brands and partners.



# Zalando SE



**17**  
countries



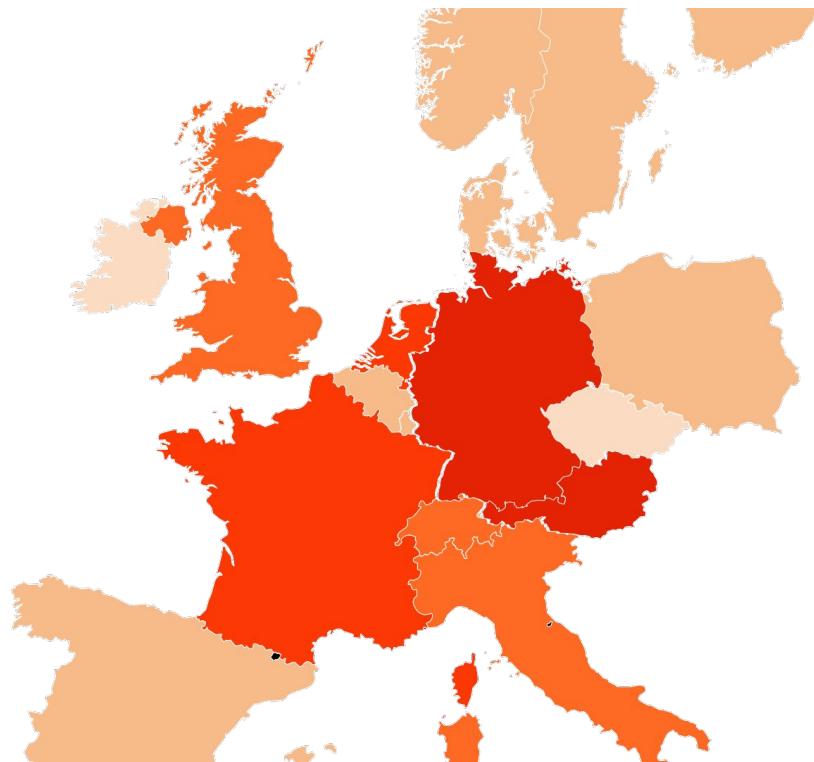
**>1,000m**  
visits last quarter



**>80%**  
mobile traffic



**>26m**  
active customers



# Big-Data Stack @ Zalando



AWS Step Functions



Google BigQuery



# About Us

## Akhil Dhingra

Product Manager, Data Solutions @Zalando  
Exp: 7+ Years, Ex-Groupon, Ex-Wingify | MBA



## Saurav Verma

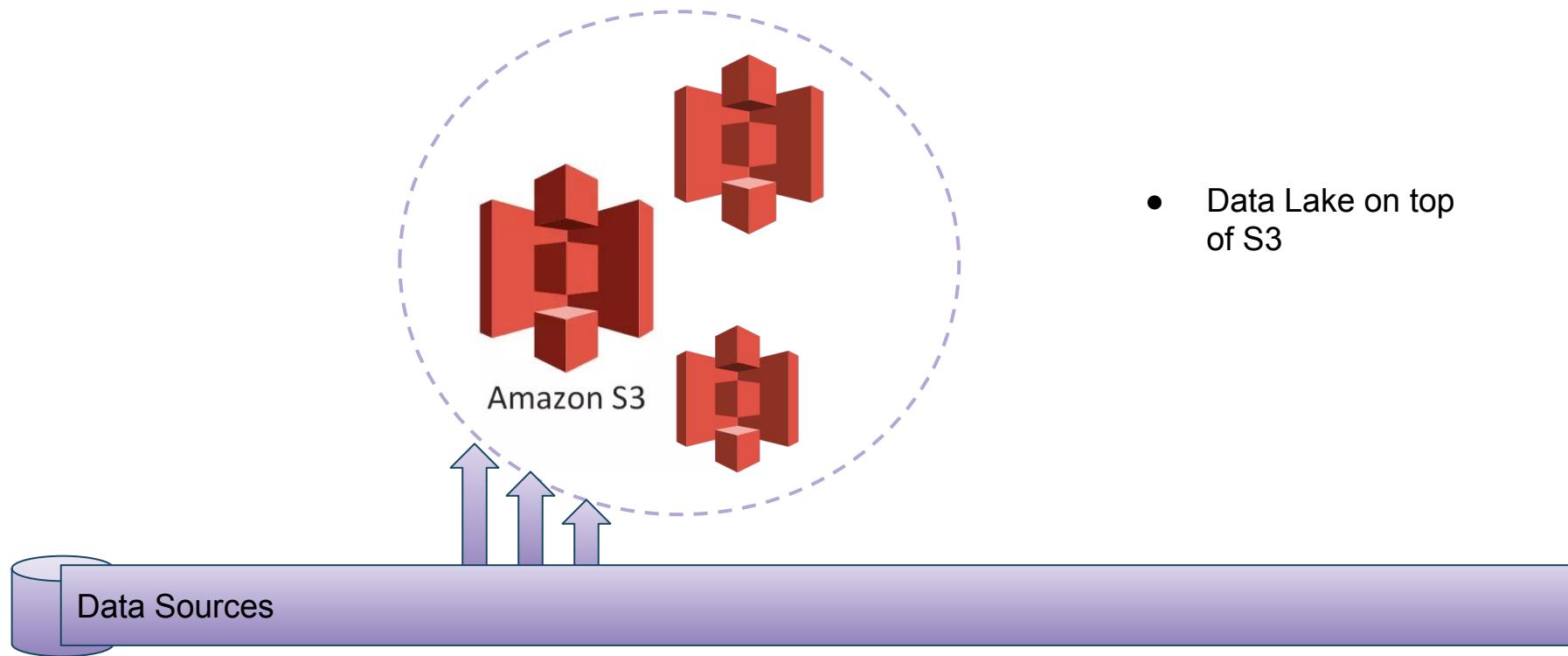
Senior Engineer, Data Lake @Zalando  
Exp: 9+ Years , Ex-Visa | Masters NUS



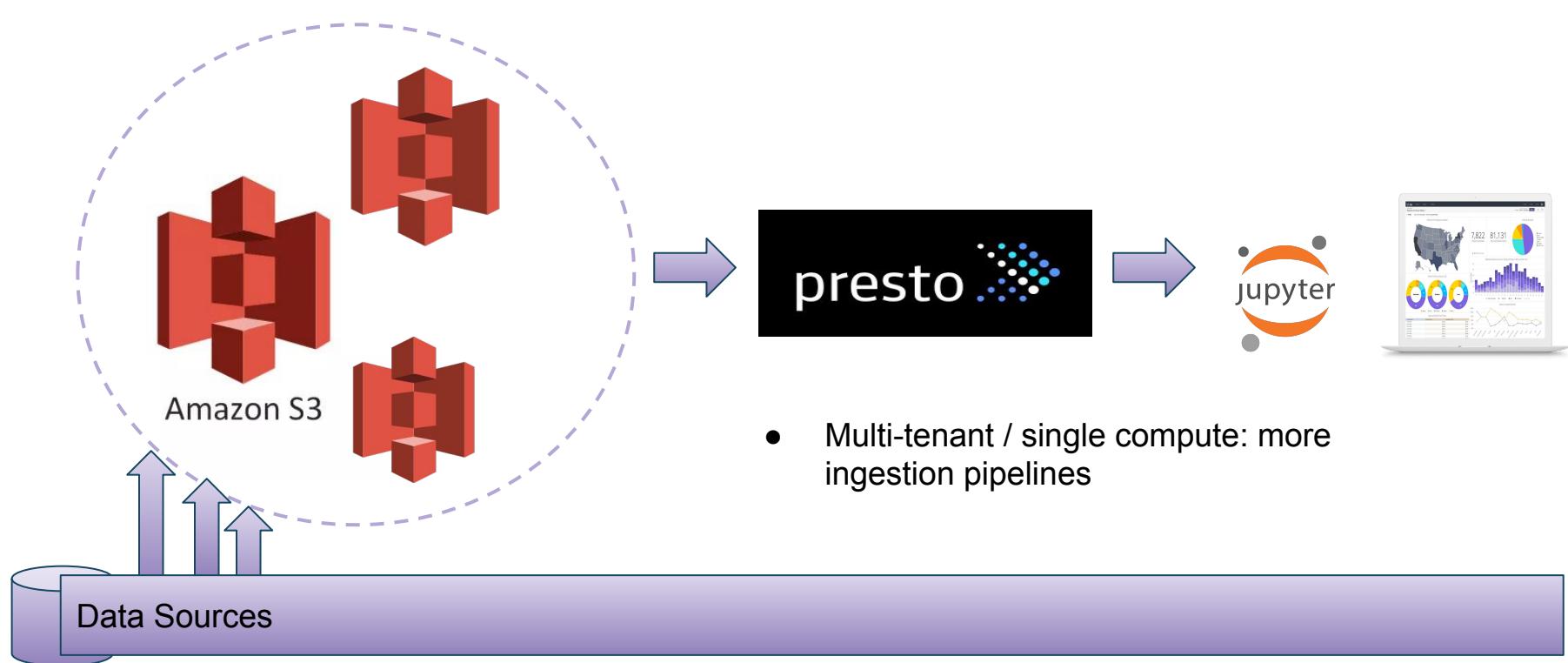
# Data Platform

Data Sources

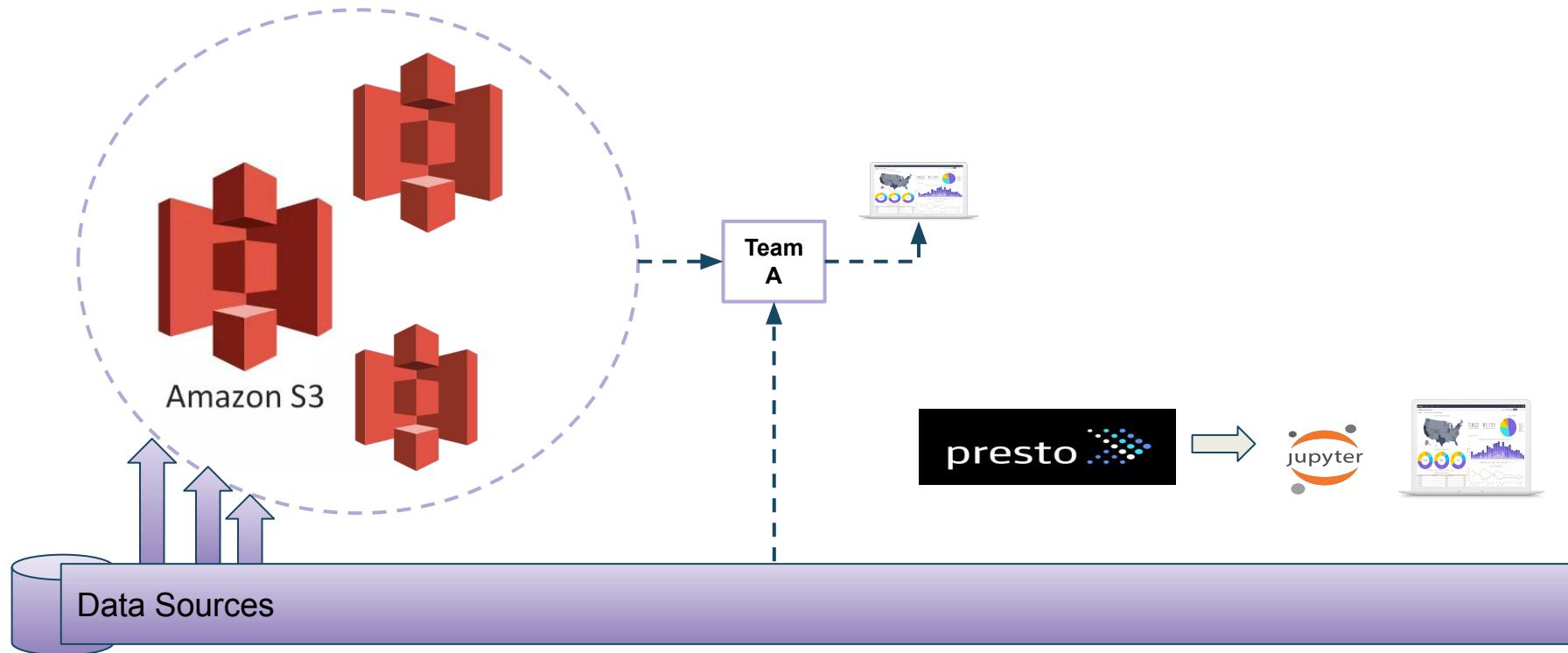
# Data Platform



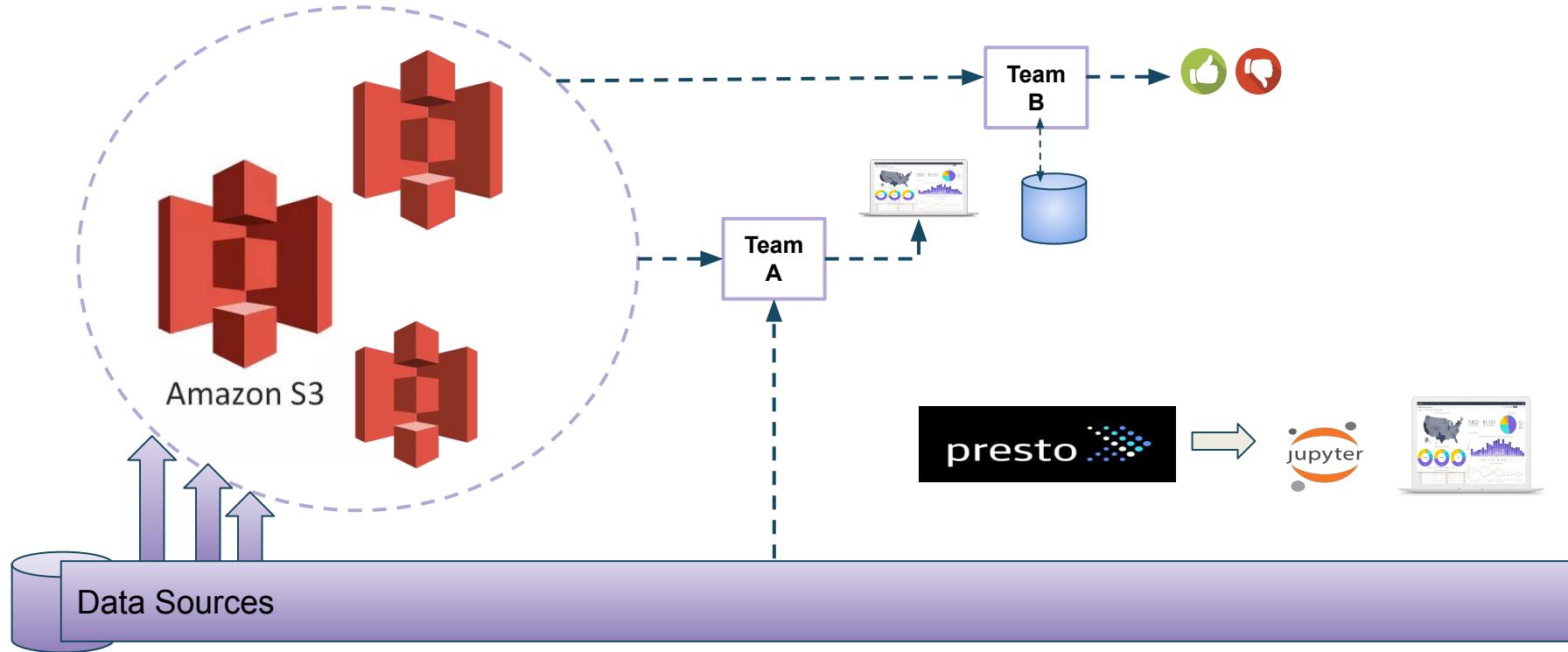
# Data Platform



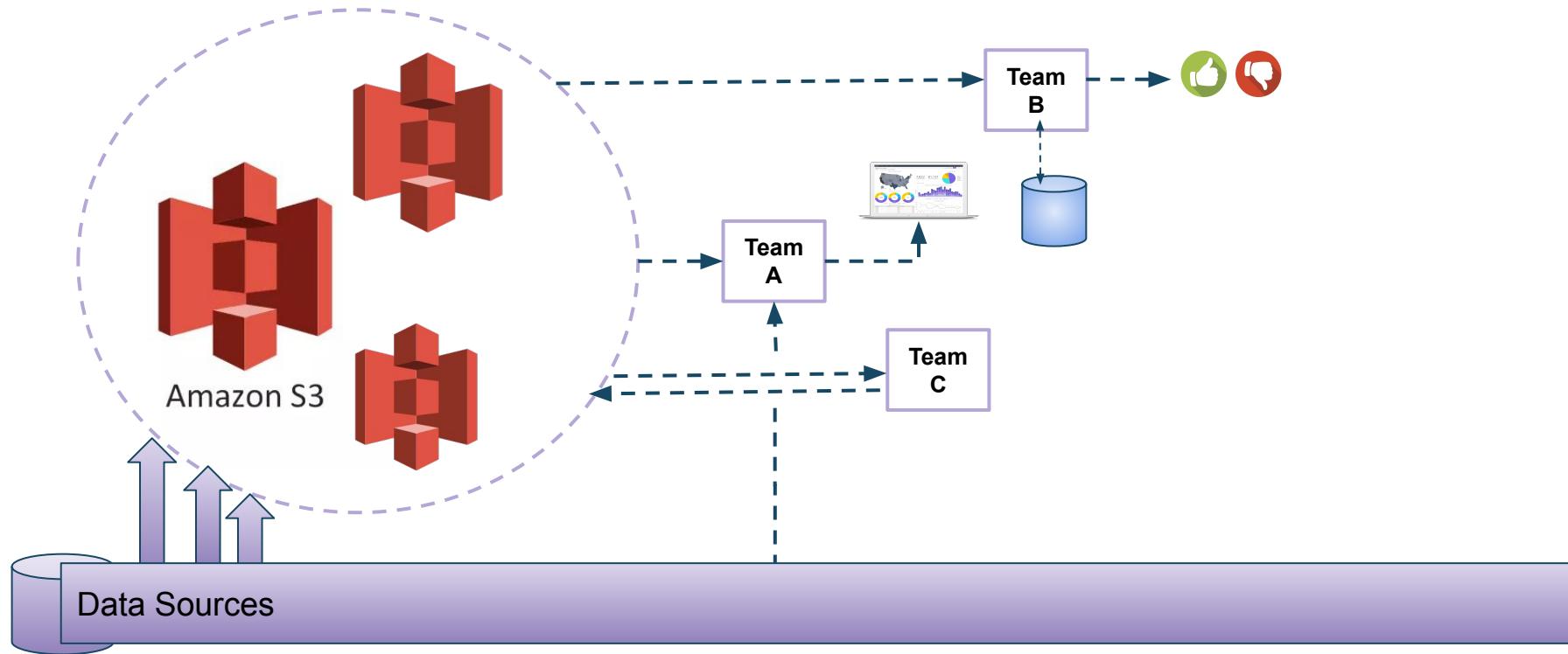
# Many Use Cases



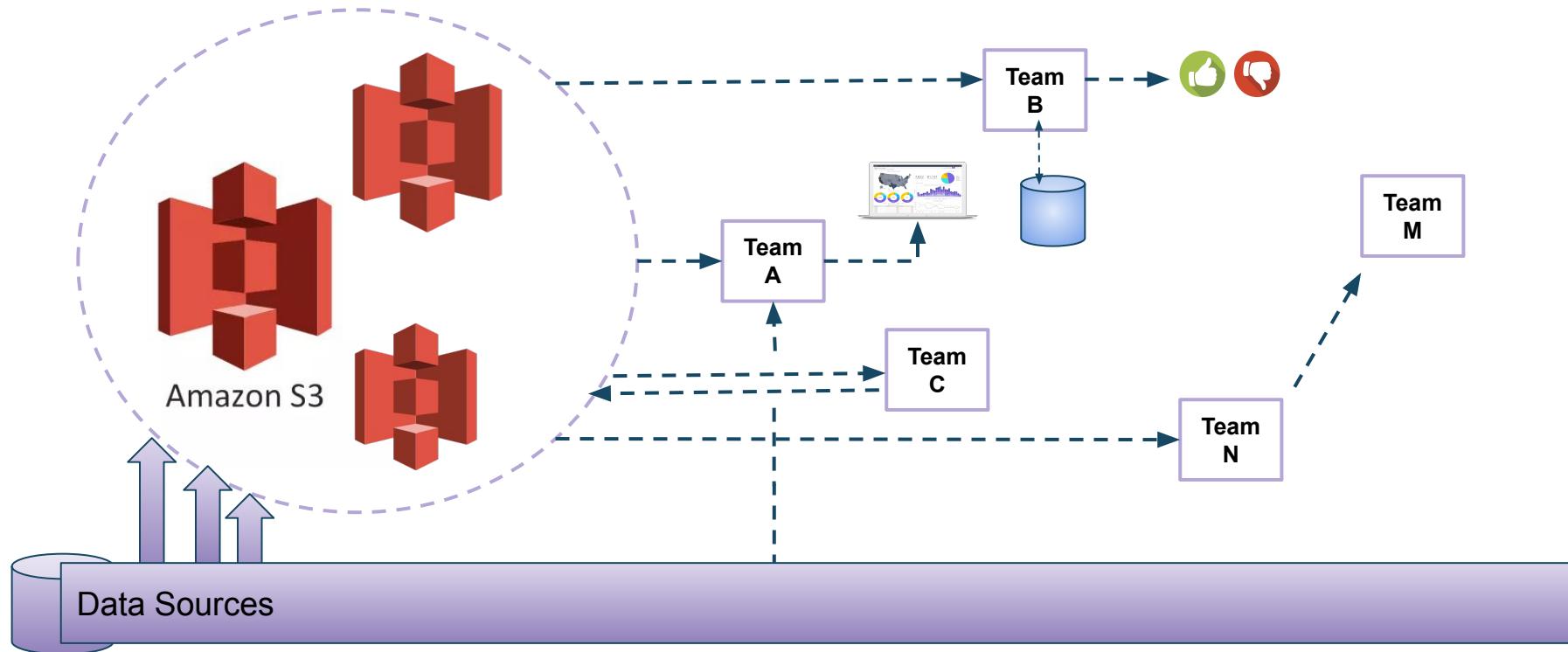
# Many Use Cases



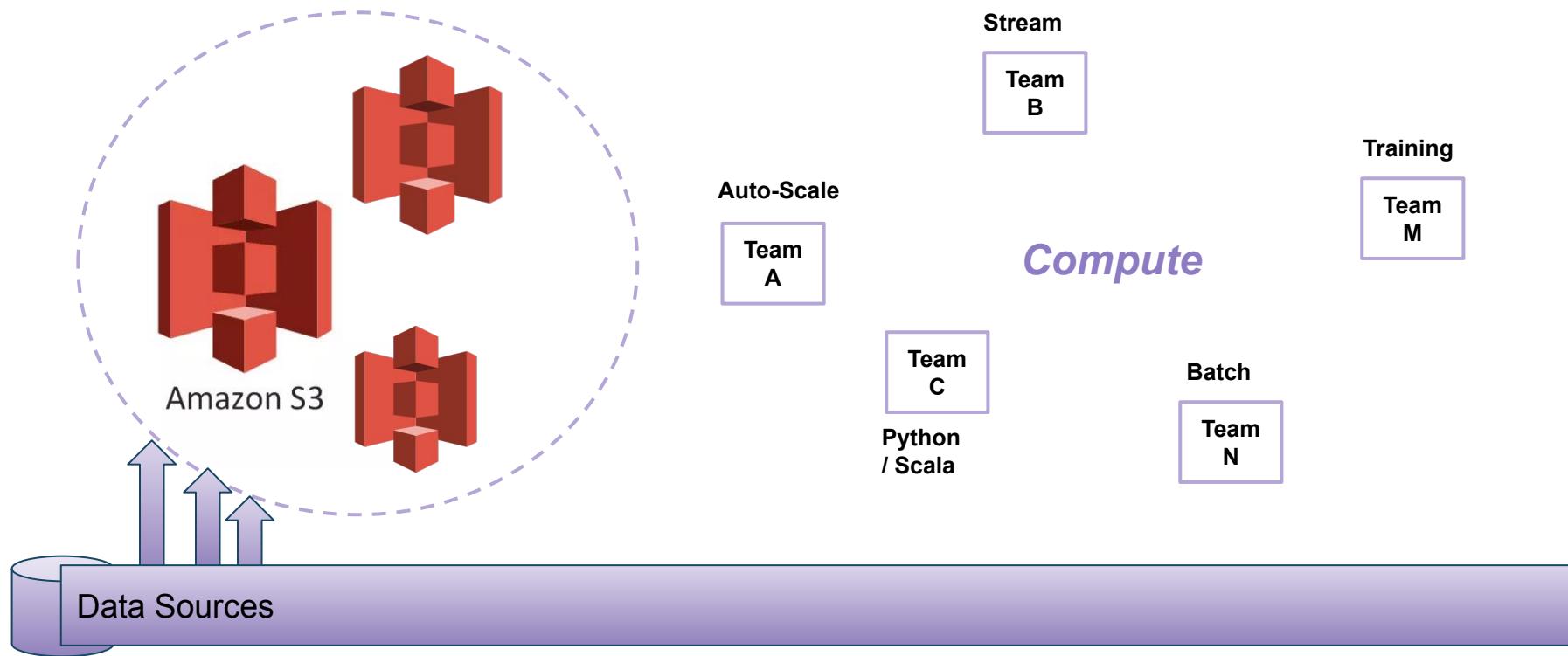
# Many Use Cases



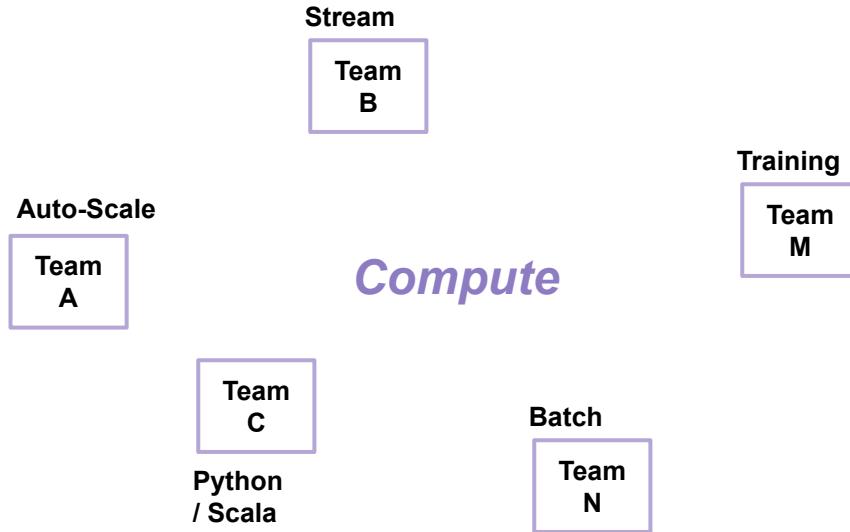
# Too Many Use Cases



# Too Many ... Compute

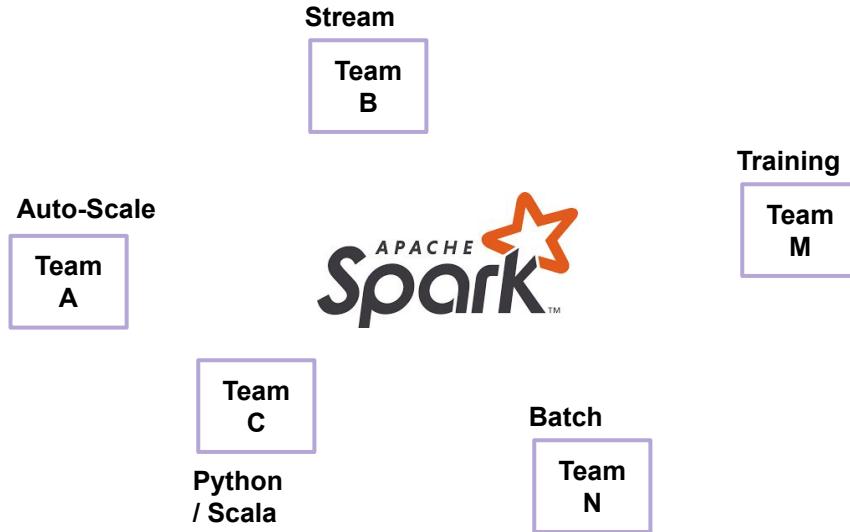


# Too Many ... Compute



- Cost control problem at Scale
- More Time To Production
- No Best Practices
- Duplication of work / Data
- Dependencies
- Inconsistent Environment
- No Community Knowledge
- Accidental Complexity

# Spark as a Service



- Foundational piece of Zalando's Big Data Infrastructure
- GitOps Management, Decentralized Clusters
- Security / Compliance / CI-CD
- XX clusters/Jobs
- ~20 teams in production
- Thriving **#Databricks** community in Zalando

# Spark as a Service



**Migration Projects**  
ETLs | Data Preparation in  
Spark-S3

# Spark as a Service



## Others:

Structured Streams |  
Traceability



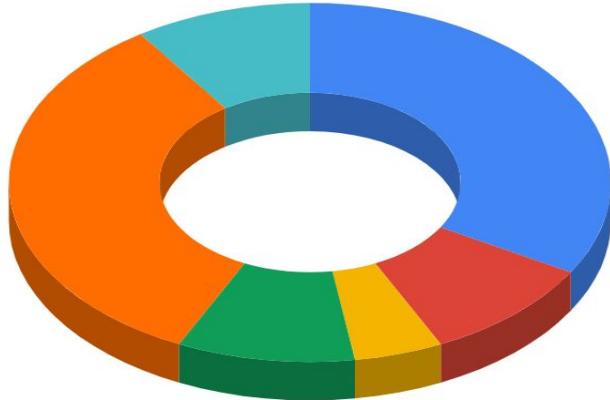
databricks



DELTA LAKE

# Spectrum of use cases

- Reporting
- Marketing
- Marketplace
- Search
- Data Products / Others
- Payments



# GDPR and Antitrust

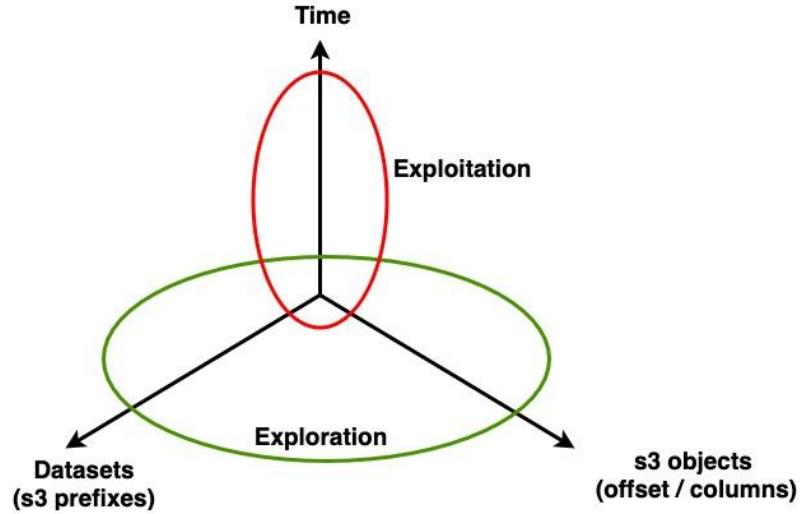
Compliance with GDPR and antitrust laws



# GDPR and Antitrust

## Probe (pilot)

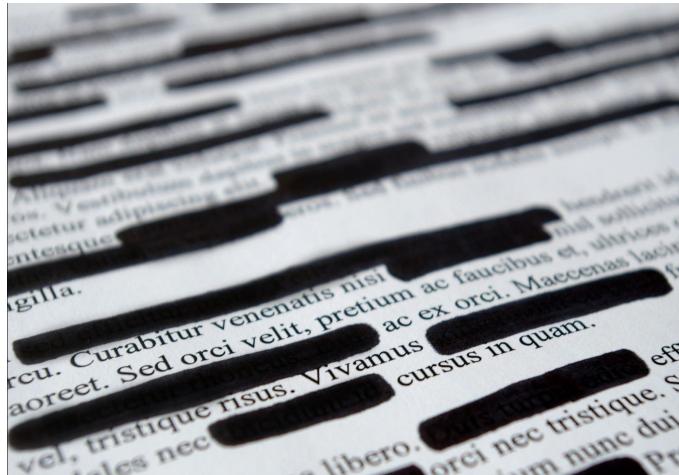
- Use marker event to create heat map of the data path.
- List of all datasets within the heat map.



# GDPR and Antitrust

## Pseudonymize/Remove

- Identifier based, on-demand, in-place record updater with field precision
- Great for semi-structured formats like JSON
- Use S3 Inventory + Streaming

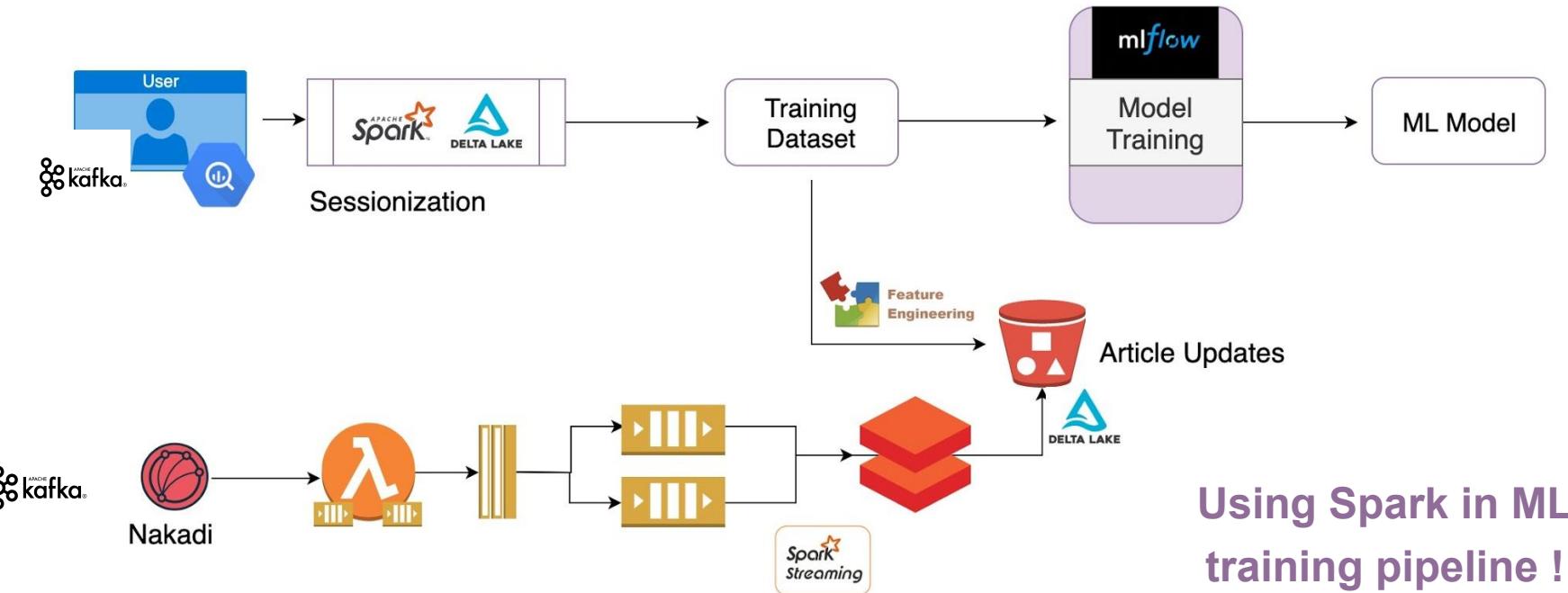


# Search & Ranking

Personalized article ranking for relevance and user engagement.

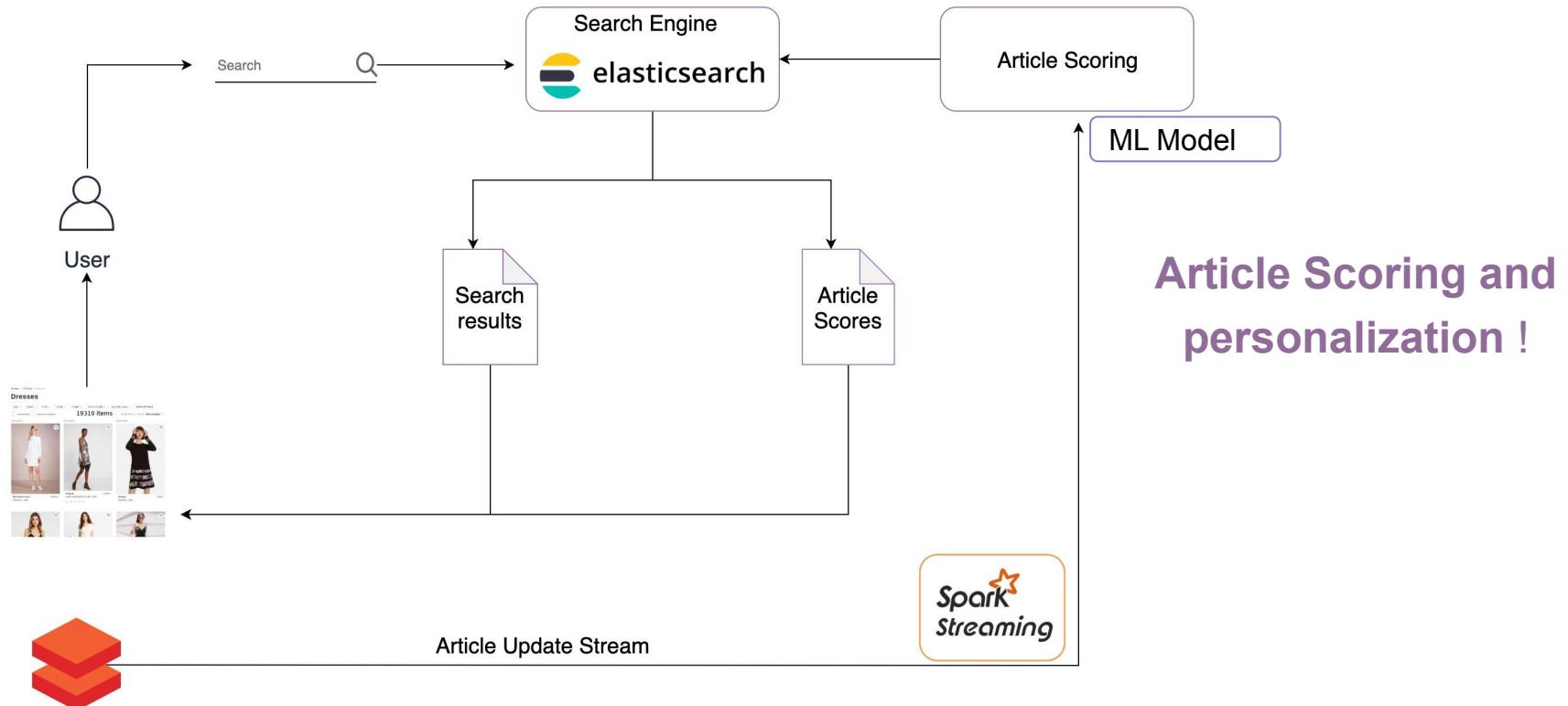


# Search & Ranking



Using Spark in ML  
training pipeline !

# Search & Ranking



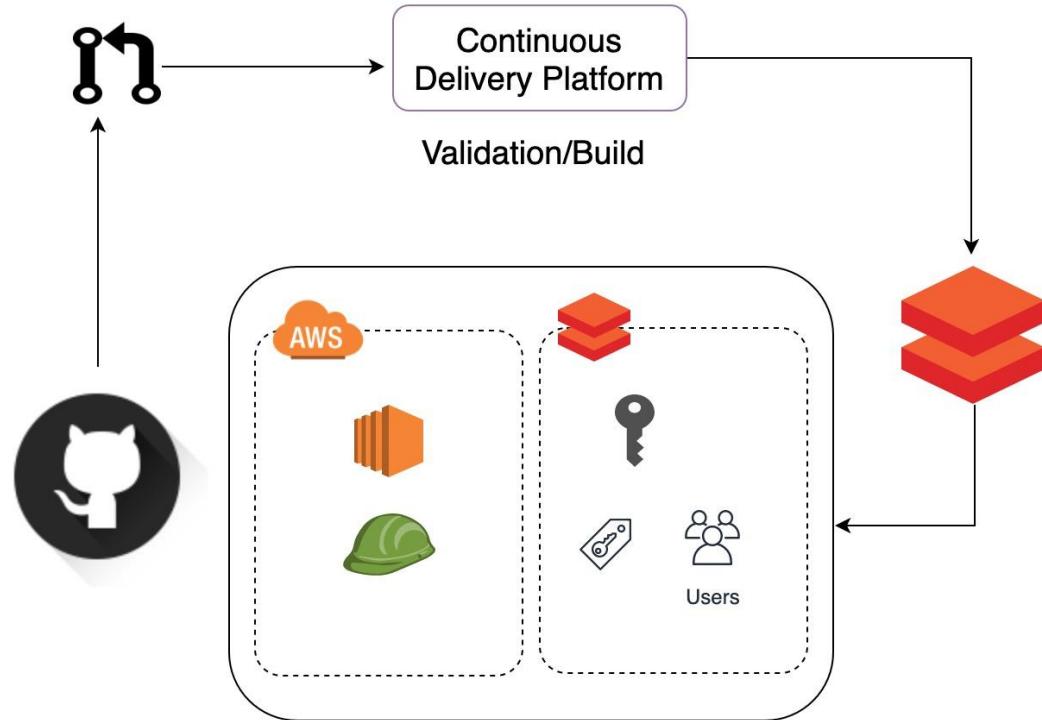
# Others

- Sizing: Reducing return rates due to size and fit issues.
- Experimentation @Scale
- Merchant Analytics
- Marketing Services



# First Impressions

- GitOps | Self Service



# First Impressions

- Multi-Tiered support system
- Delta Adoption | But few readers outside Databricks ecosystem
- Communicating pricing downstream
- Exploding Usage is Good
- Fits all Size?



# Thank you.

## AI-Powered Retail Experience with Databricks

Akhil Dhingra  
Saurav Verma

[www.zalando.com](http://www.zalando.com)  
[www.jobs.zalando.com/tech](http://www.jobs.zalando.com/tech)





SPARK+AI  
SUMMIT 2019

DON'T FORGET TO RATE  
AND REVIEW THE SESSIONS

SEARCH SPARK + AI SUMMIT



SPARK+AI  
SUMMIT 2019