



SPARK+AI
SUMMIT EUROPE

Lambda Architecture in the Cloud with Azure Databricks

Andrei Varanovich, InSpark

#SAISDev6



Selfie



INSPARK

Data & AI Lead



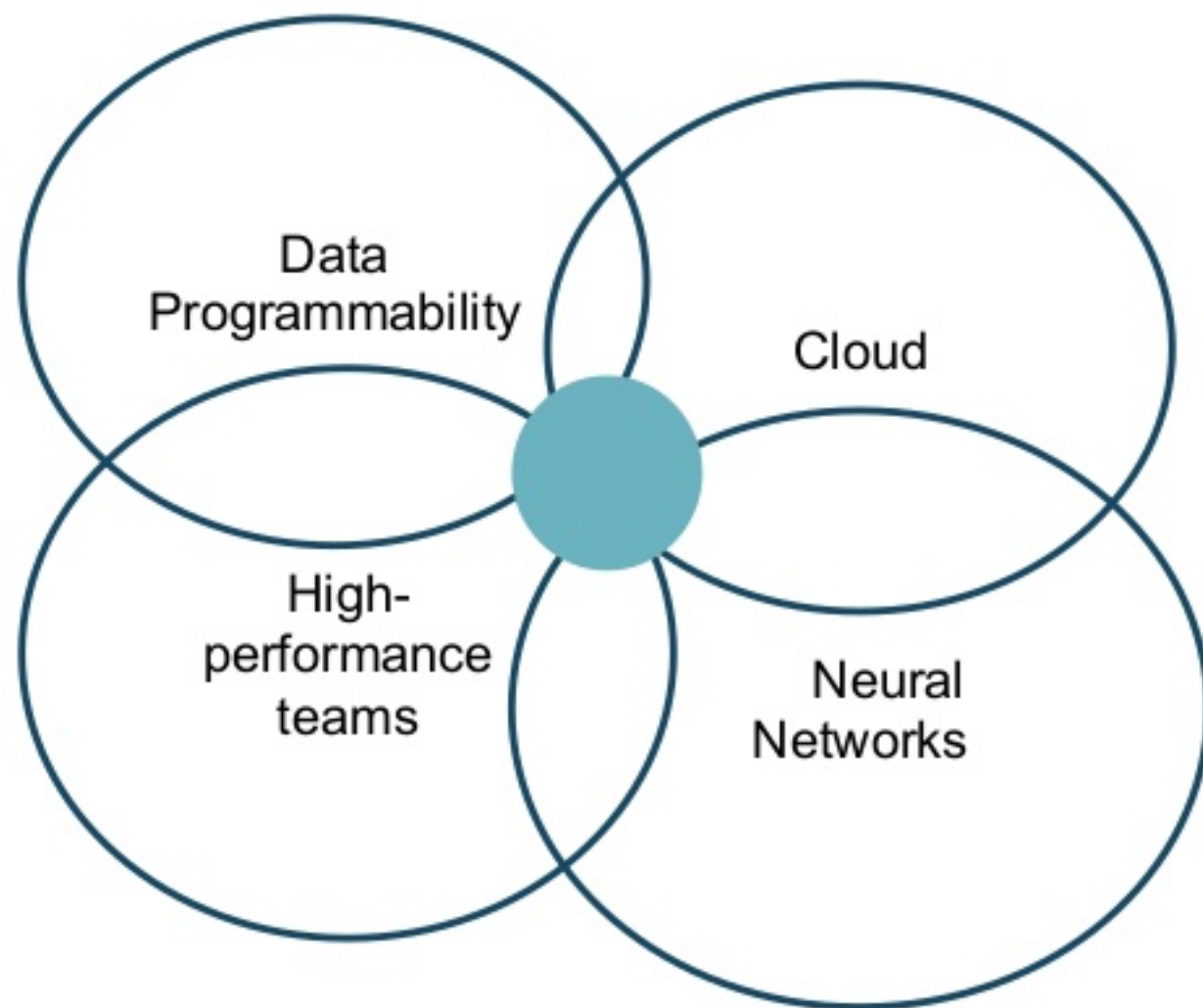
@DrGigabit



andrei.varanovich



andrei.varanovich@inspark.nl



**Big Data problem is many
small data problems**

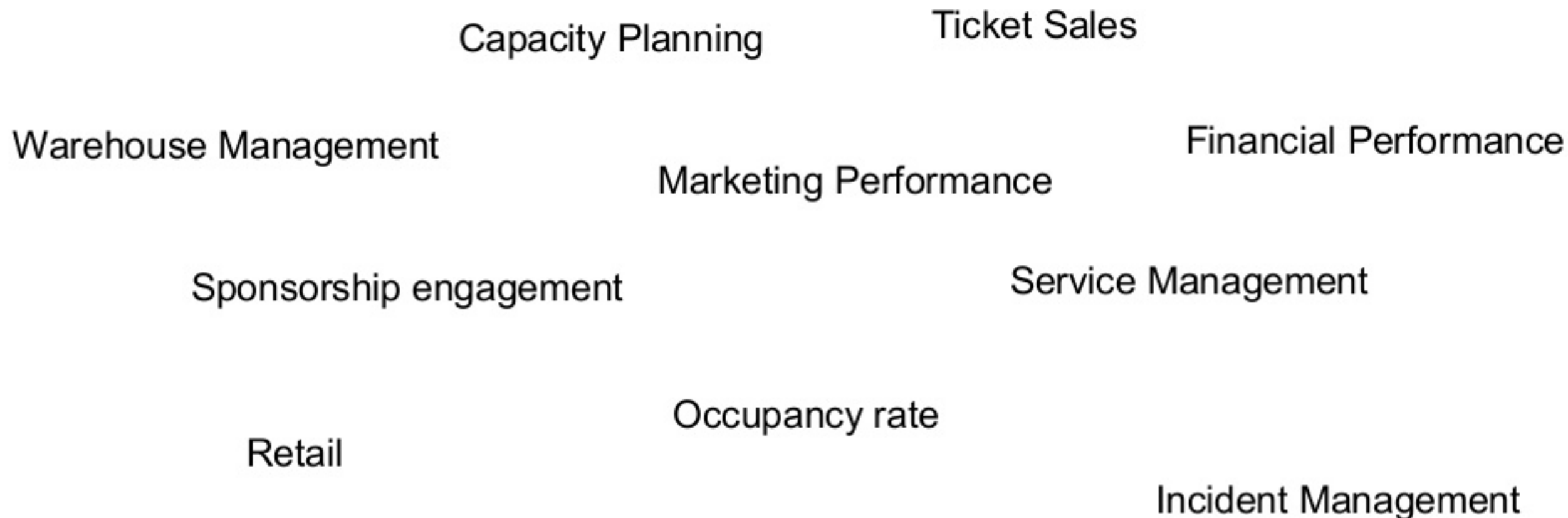
A wide-angle photograph of the Rijksmuseum in Amsterdam, showing its iconic red brick facade and multiple towers. The museum is situated behind a canal, with trees and a walkway in the foreground. The sky is blue with some clouds.

2.500.000 visitors per year

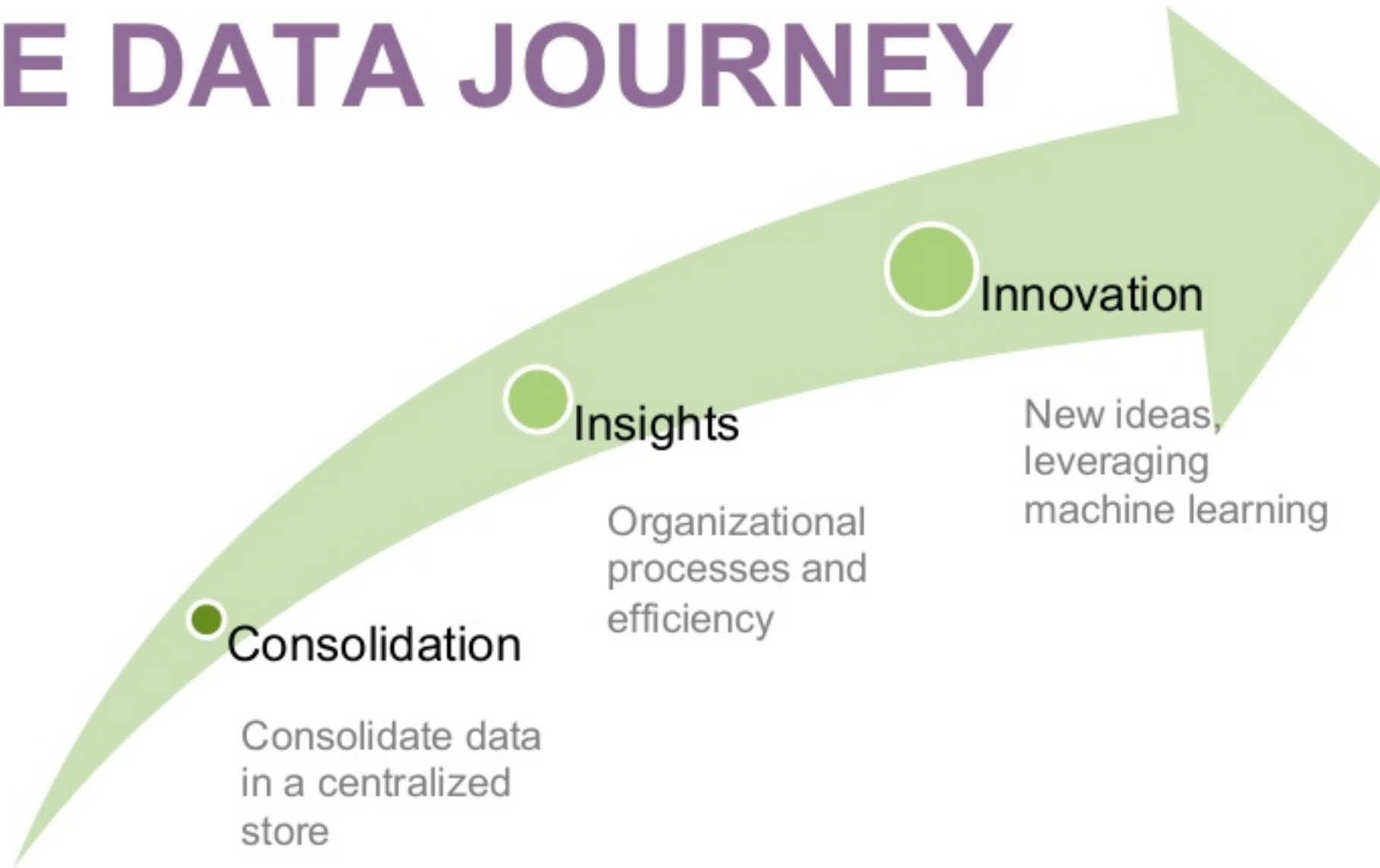
8.000 objects of art and history

1.000.000 objects stored from
the year 1200

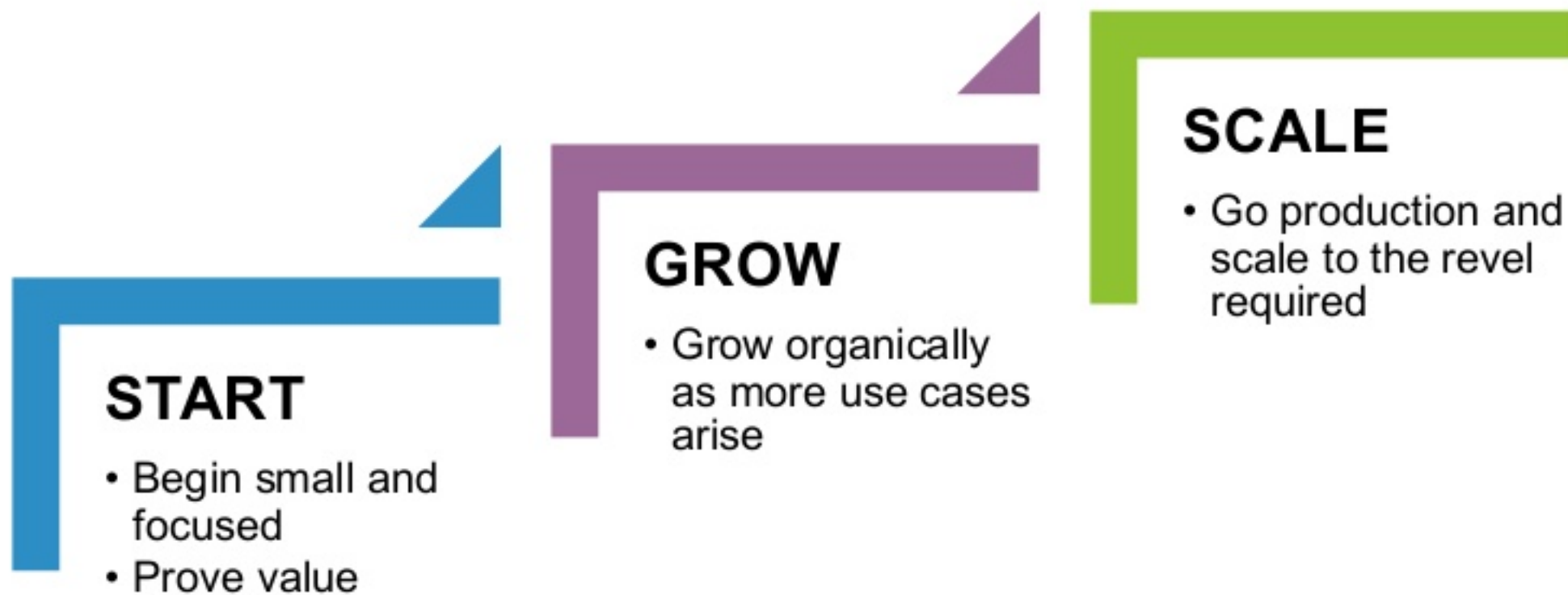
Under the hood



THE DATA JOURNEY

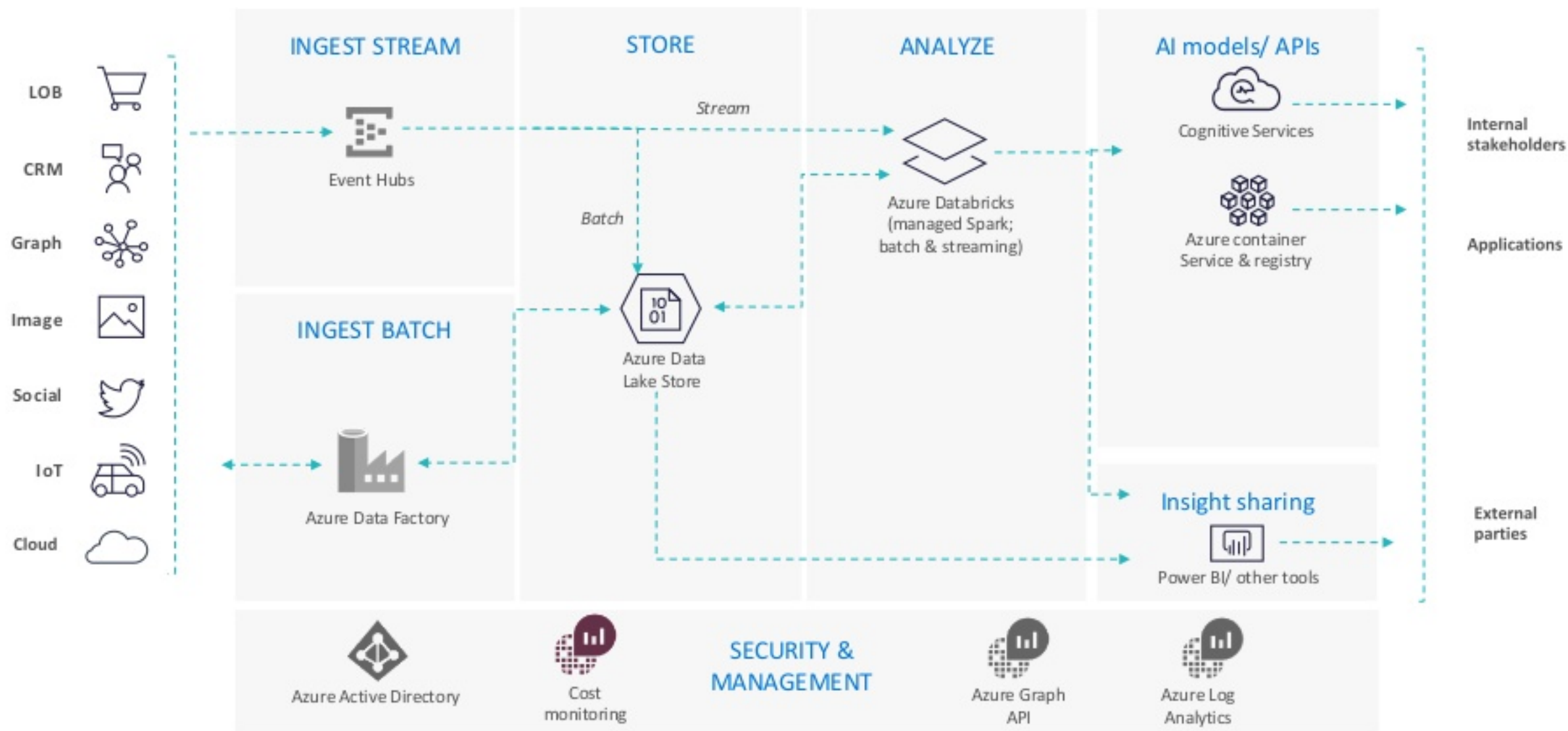


IN THE NEED FOR THE PLATFORM



We are in the need of the truly elastic data platform, to avoid any upfront planning, deployment and operations expenses, and put business value discovery first. The platform should support the [big]data projects in any stage, without the need to reengineer the whole solution.

Lambda Architecture on Azure



IoT / streaming data

Cloud storage

Data warehouses

Hadoop storage

Azure Databricks

Collaborative workspace

DATA ENGINEER

DATA SCIENTIST

BUSINESS ANALYST

Production jobs & workflows

MULTI-STAGE PIPELINES



JOB SCHEDULER



NOTIFICATION & LOGS

Optimized Databricks Runtime Engine



DATABRICKS I/O



APACHE SPARK



SERVERLESS



Rest APIs



Machine learning models



BI tools



Data exports



Data warehouses

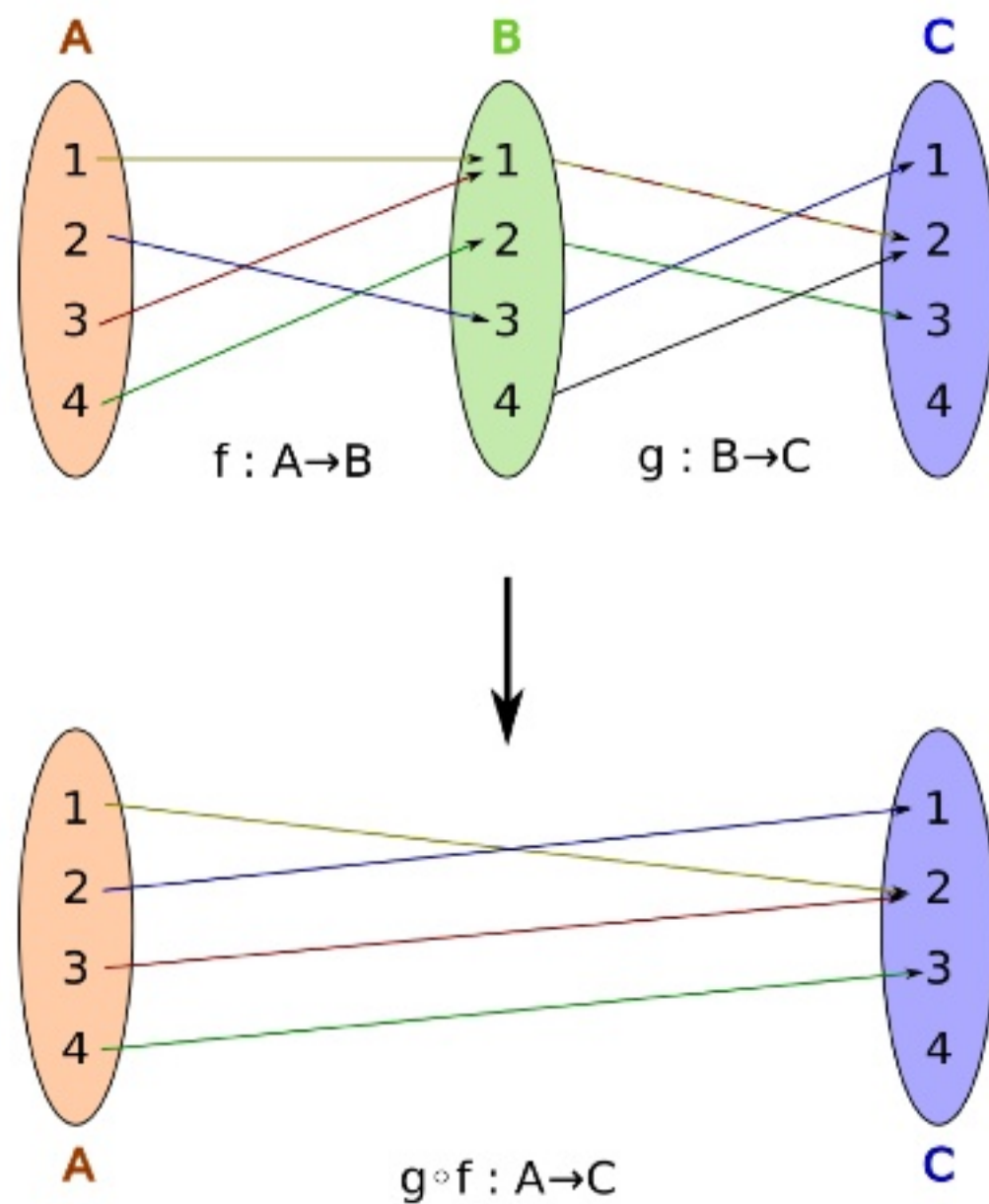
Simplicity is the ultimate sophistication

Leonardo da Vinci



LAMBDA TO THE RESCUE



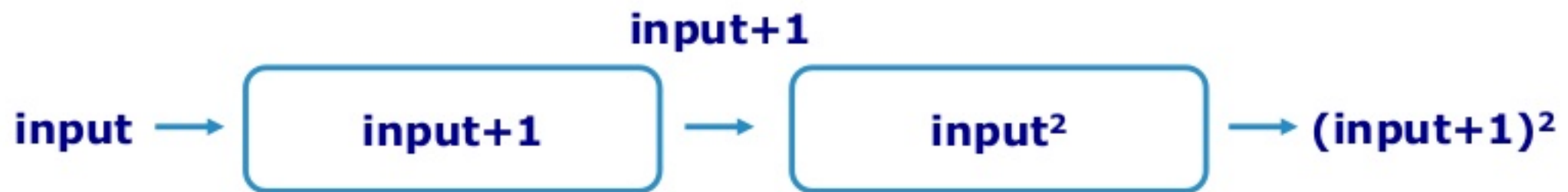


Composition of functions is applying one function to the result of another

$$f(x) = x+1$$

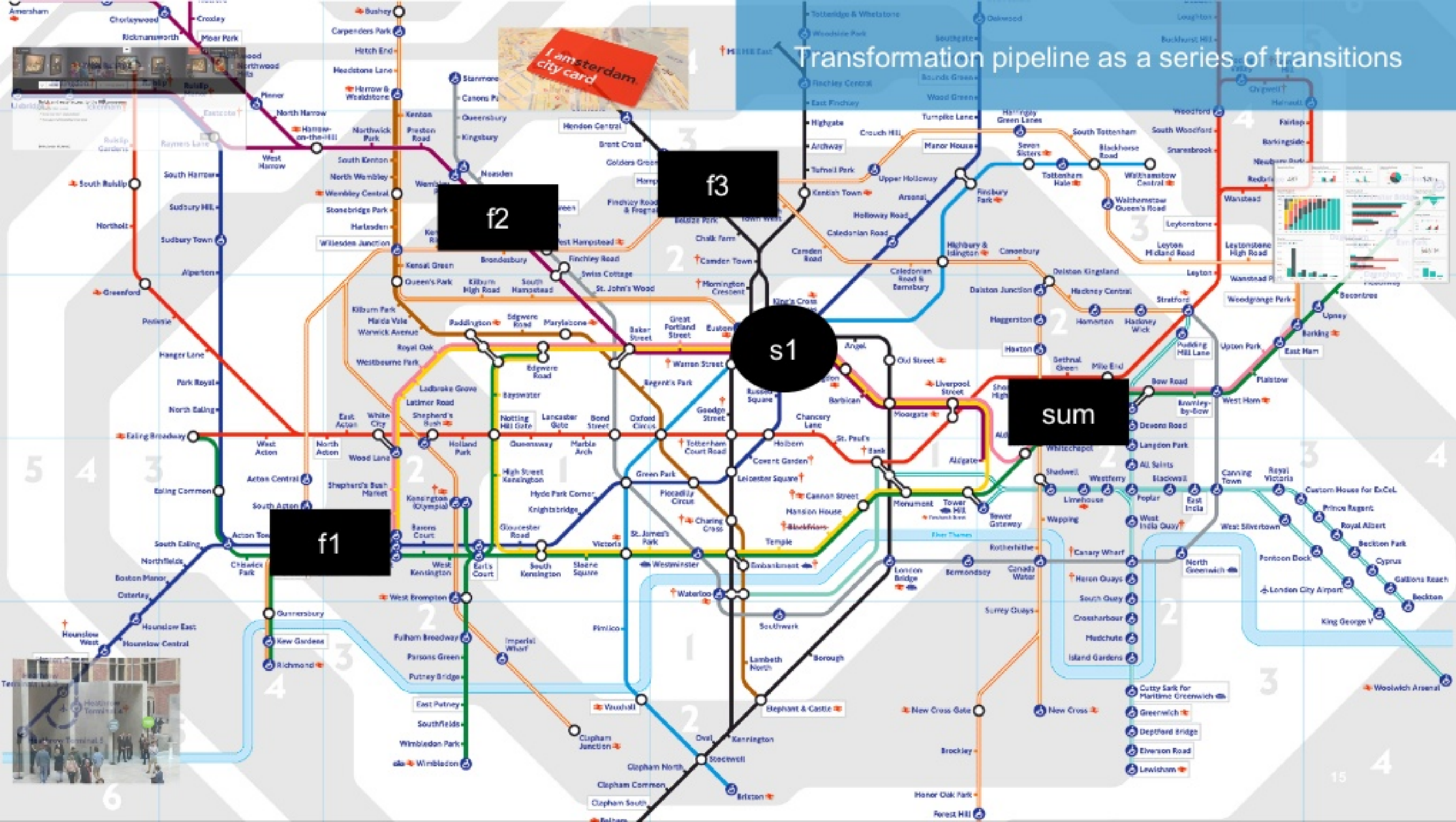
$$g(x) = x^2$$

$$(g \circ f)(x) = g(f(x))$$



$$(g \circ f)(x) = (x+1)^2$$

Transformation pipeline as a series of transitions



IoT / streaming data

Cloud storage

Data warehouses

Hadoop storage

Azure Databricks

Collaborative workspace

DATA ENGINEER

DATA SCIENTIST

BUSINESS ANALYST

Production jobs & workflows

MULTI-STAGE PIPELINES

JOB SCHEDULER

NOTIFICATION & LOGS

Optimized Databricks Runtime Engine

DATABRICKS I/O

APACHE SPARK

SERVERLESS

Rest APIs

Machine learning models

BI tools

Data exports

Data warehouses

Conclusions

... with proper design, the features come cheaply. This approach is arduous, but continues to succeed.

—Dennis Ritchie

- Standardization on Apache Spark allows us to move forward without introducing extra complexity.
- 100% PaaS offering is important – no need to maintain the infrastructure. All components we use offered as PaaS on Azure.
- Data pipelines as function composition allows us to ensure end-to-end consistency and spot the errors quickly.
- Saving intermediate states allows to quickly inspect the data sets.

Thank you!

Questions?