

Data Analytics Platform

or How to Make Data Science in a Box Possible

Krzysztof Adamski ING WBAA

Moscow, 11th October 2018

thinkforward

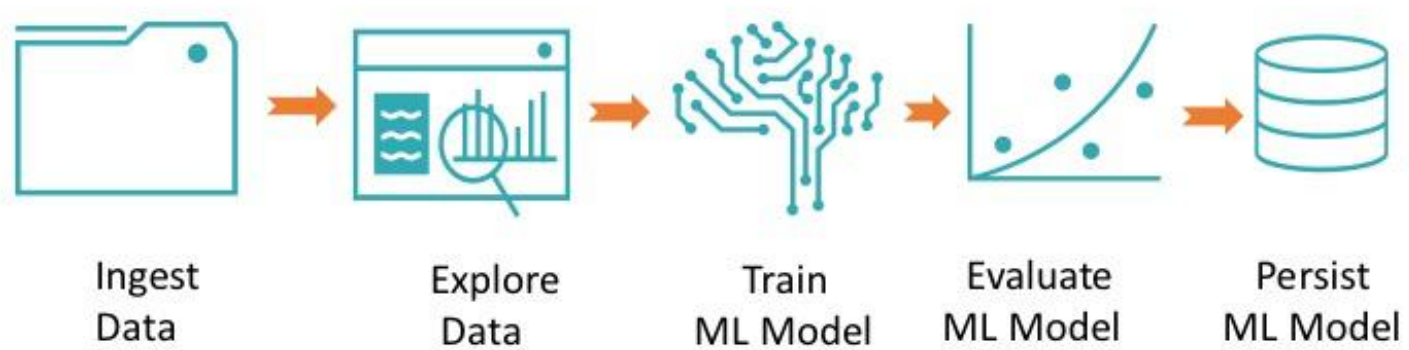


What am I doing here



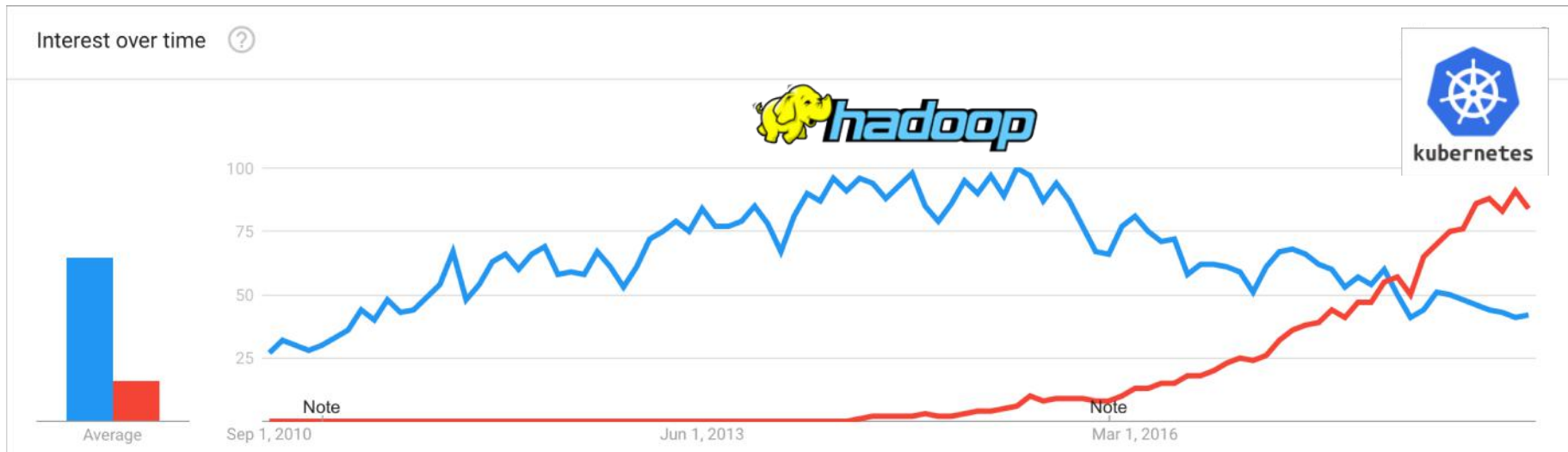
- Data Infrastructure Architect
- Focused on enabling data-driven decisions
- Main areas:
 - Large scale data storage and processing – this has to run somewhere
 - Data platforms security – we are a bank and trust is the biggest asset
 - Sharing the knowledge – as a nature is complex

Data Engineering Pipeline



The new post Hadoop world

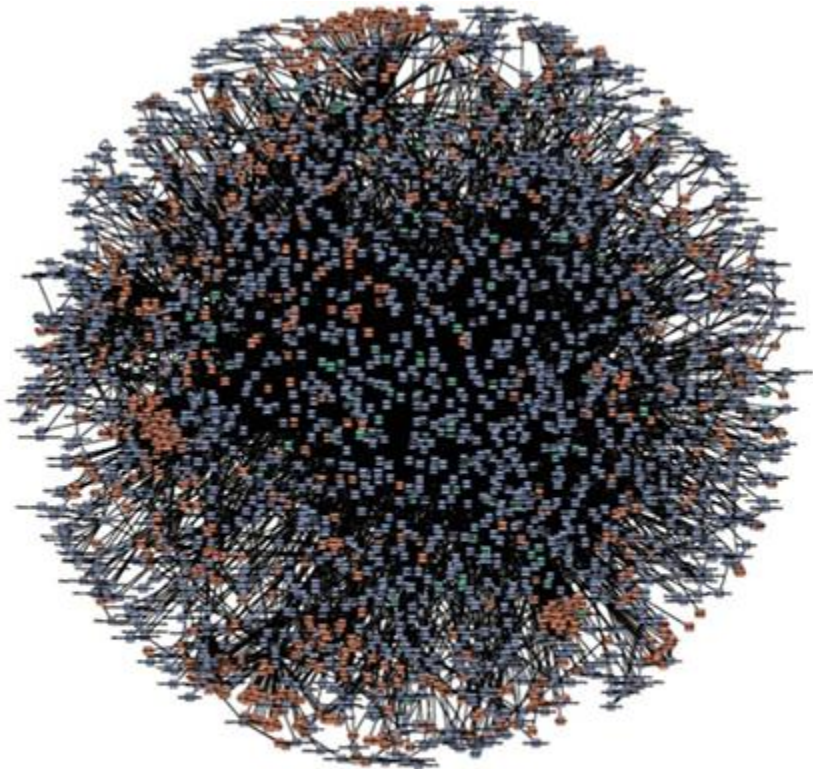
- Container based
- Need a more agile adoption of a growing set of technologies
- Stronger isolation is required
- Cloud-Readiness as an integral part of analytics frameworks



Hadoop as an isolated island



And the world nowadays looks like this



amazon.com®

NETFLIX

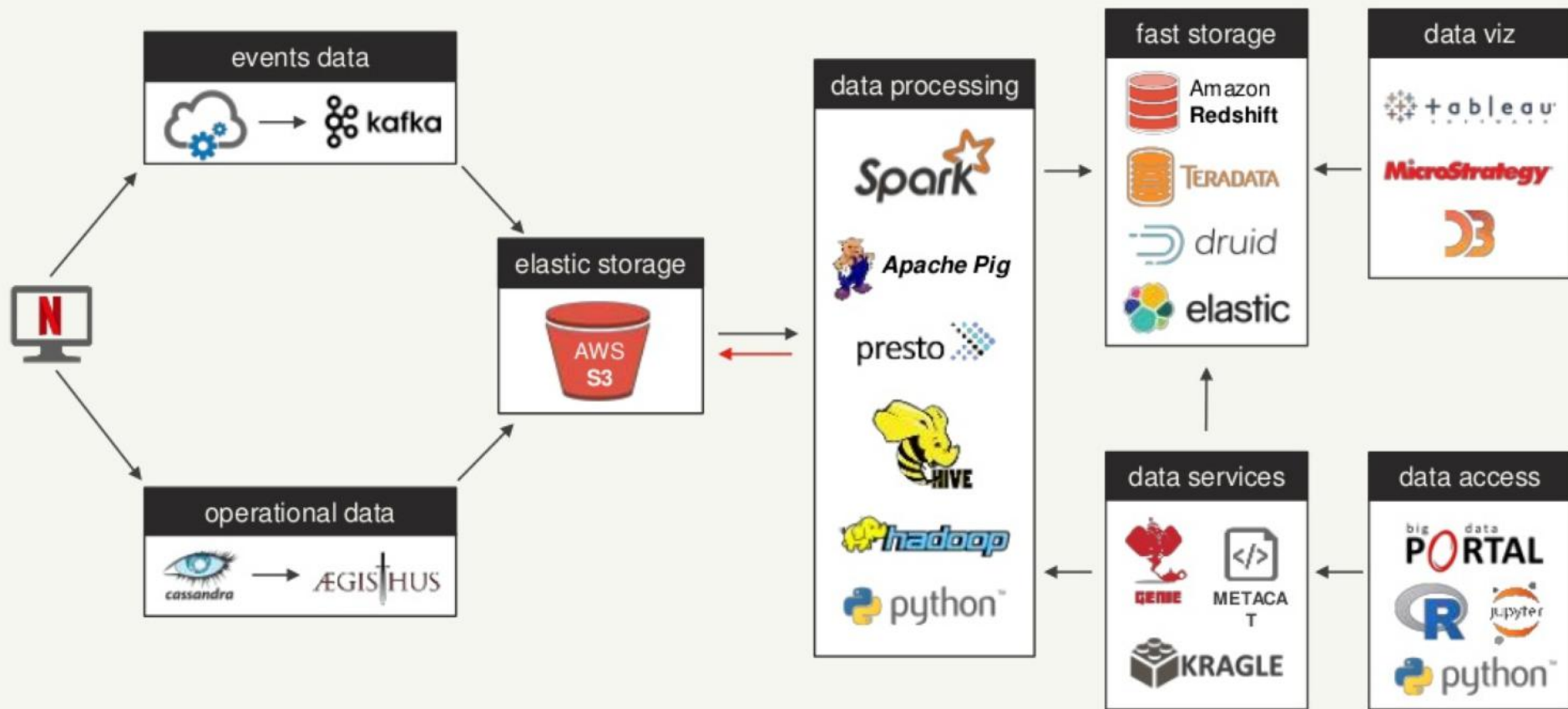
Netflix scale

Around **700 microservices** to control each of the many parts of what makes up the entire Netflix service



<https://medium.com/refraction-tech-everything/how-netflix-works-the-hugely-simplified-complex-stuff-that-happens-every-time-you-hit-play-3a40c9be254b>

Big Data Platform





Azure



cloudera

PLATFORM



MAPR

ORACLE
Data Cloud

TERADATA



dremio



databricks



IBM WATSON

I see platforms everywhere



Team & Expansion

On average we receive one application every 15 minutes

Madrid
Southern European Hub

100+ people

- Business Developers
 - User Experience Designers
 - Data Engineers
 - Data Scientists
 - Software Engineers
- 16 nationalities

Data Talents

- 65 NL
- 9 PL
- 5 (23) LDN

Where our people come from

Thought Leadership

Cancer research hackathon

Inspiration workshops and translator training

Experimentation days

Code contributions

Apache Airflow
Apache Spark

10 commandments

AI Master Thesis Supervisors

Data Ethics

DATA SCIENCE IN A BOX

Core Tech & Algorithms

ENTITY MATCHING

PRIVITAR

Safe sharing data with 3rd parties

PPI

Possible Private Individual Protection

PEER ANALYSIS

IADO

Realtime pricing in emerging market bonds (Katana)

Anchor Products

Katana Lens

Trading signals for asset management companies.

Katana Internal

Pricing analytics for bond traders.

Domino

Network analytics for payments and loans.

450+ users

New Products

Mandoline

R.I.P.

Celebrate failures!

KYC

MENTORING
Startupbootcamp

With FinTechs

WHOLESALE BANKING
advanced analytics

Building algorithmic data-driven products with a 10x impact.

Intellichain

Automated supply chain analytics.

Beyond Banking

AXYON.AI

Loan syndication prediction.

INTELMATCH

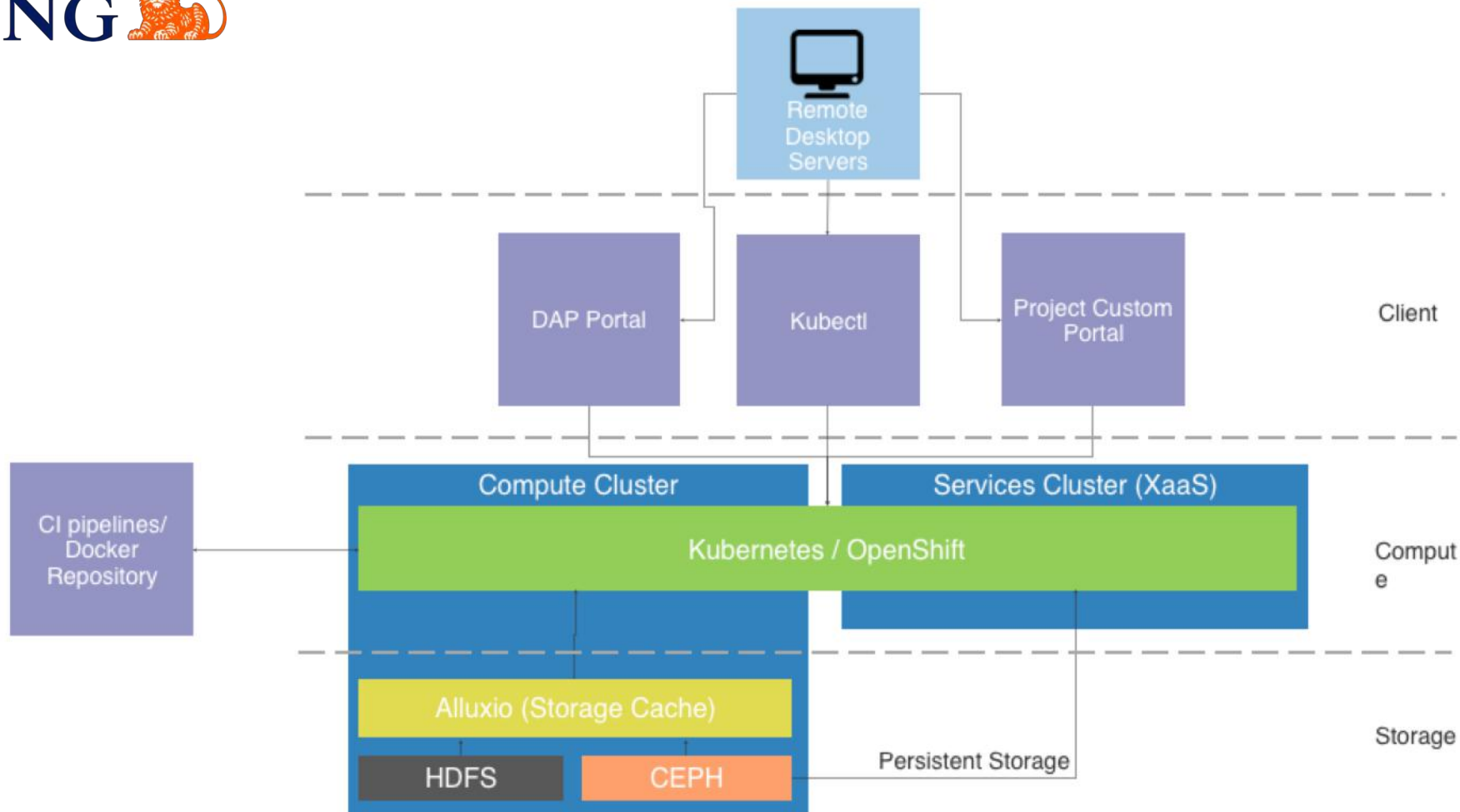
Smart knowledge management system.

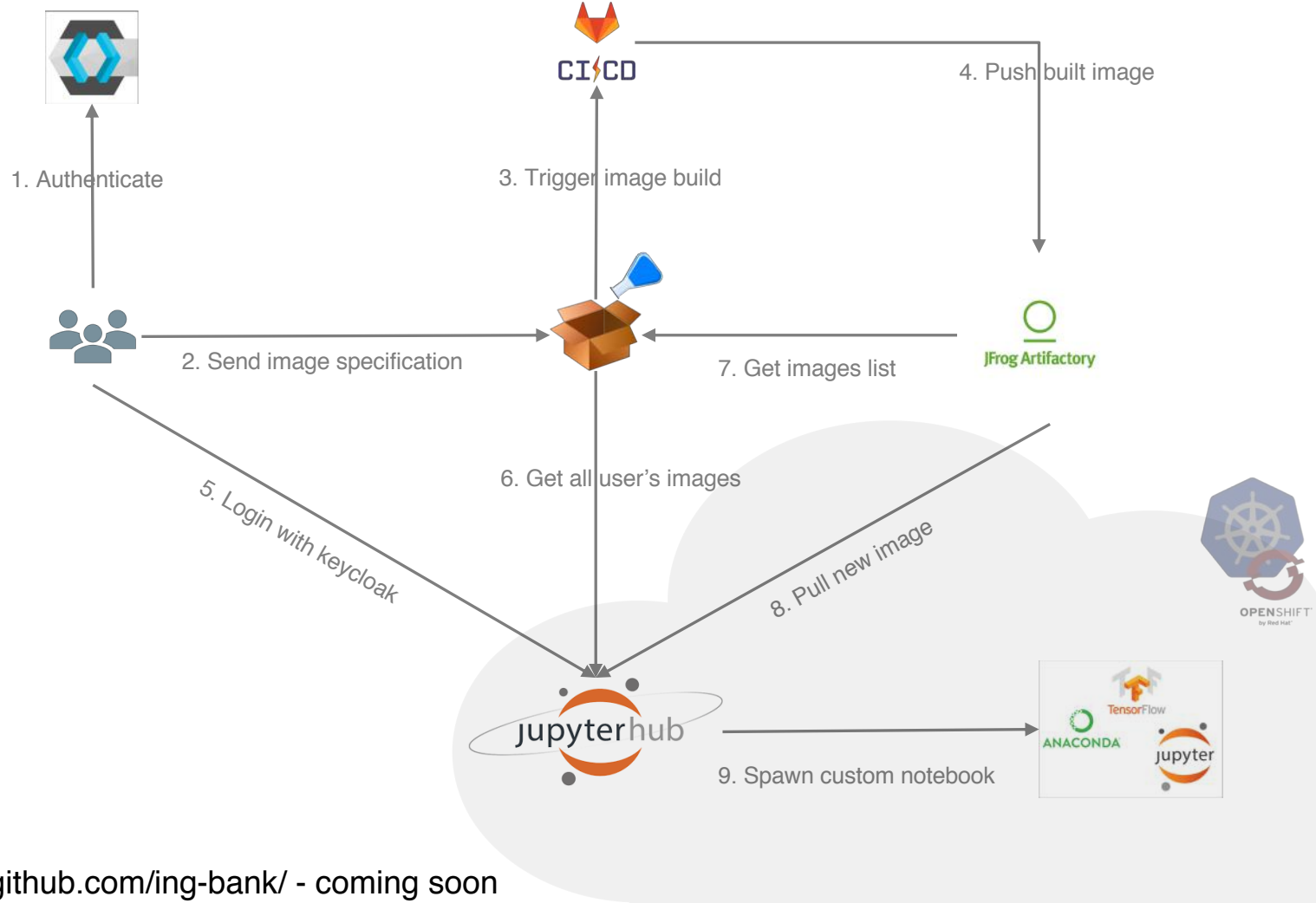
Owlin

News analytics.

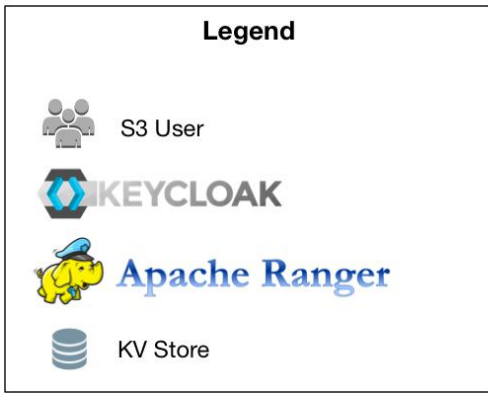
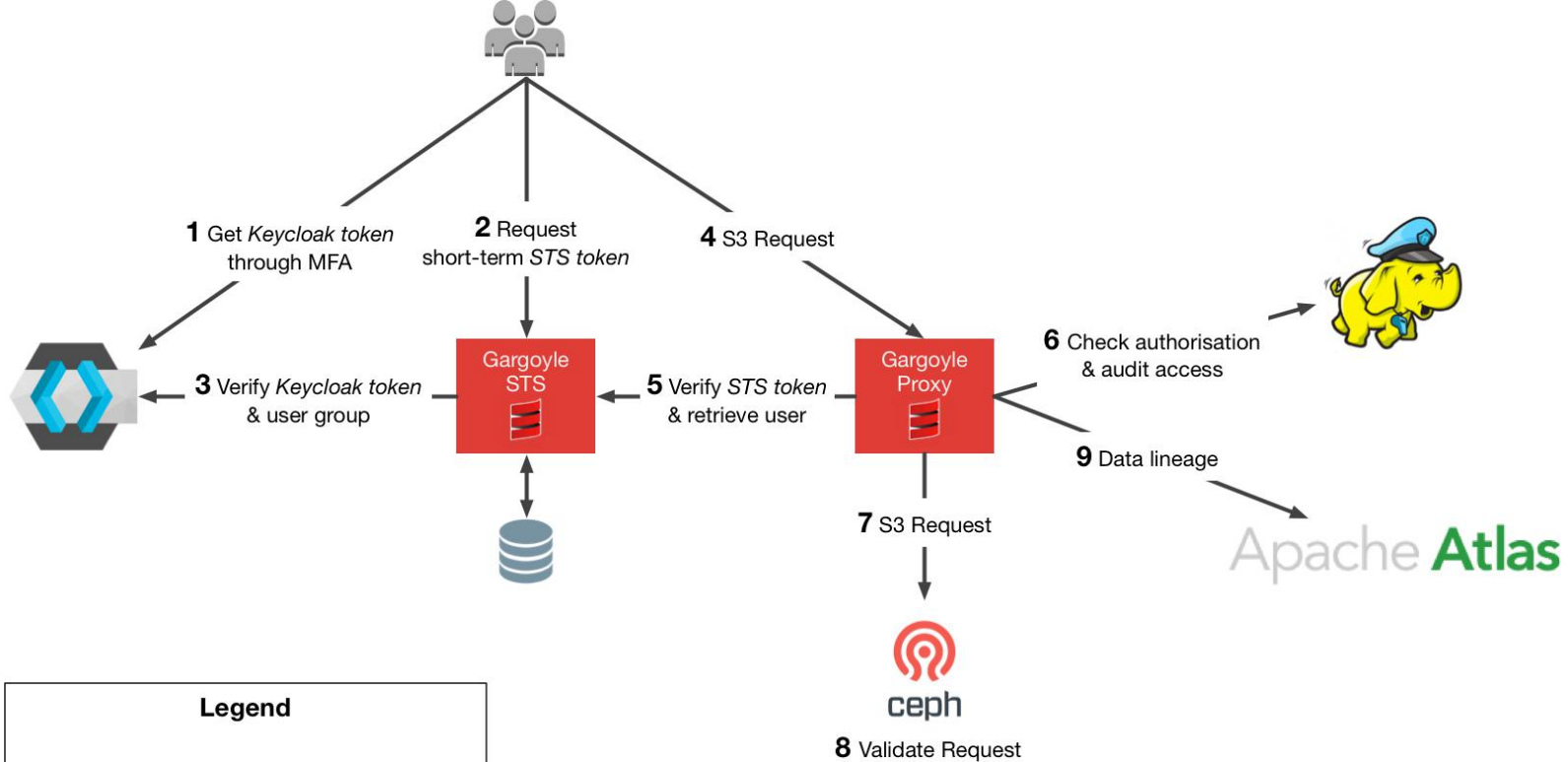
SPAZIODATI

Knowledge graph for credit packs.





<https://github.com/ing-bank/> - coming soon



<https://github.com/ing-bank/airlock>

Heads up for Spark 2.4 and beyond

- Pyspark
- R
- Client mode (finally you can have your notebooks running)
- Kerberos support is almost there
- Dynamic resource allocation and external shuffle service – still pending
- Spark streaming still missing driver resillience
- Highly recommended series of articles

<https://banzaicloud.github.io/tags/spark/>



**KEEP
CALM
IT IS
DEMO
TIME**

...and keep fingers crossed