

July 16, 2020

Autonomous Driving with Airflow

Where big-data meets high performance computing

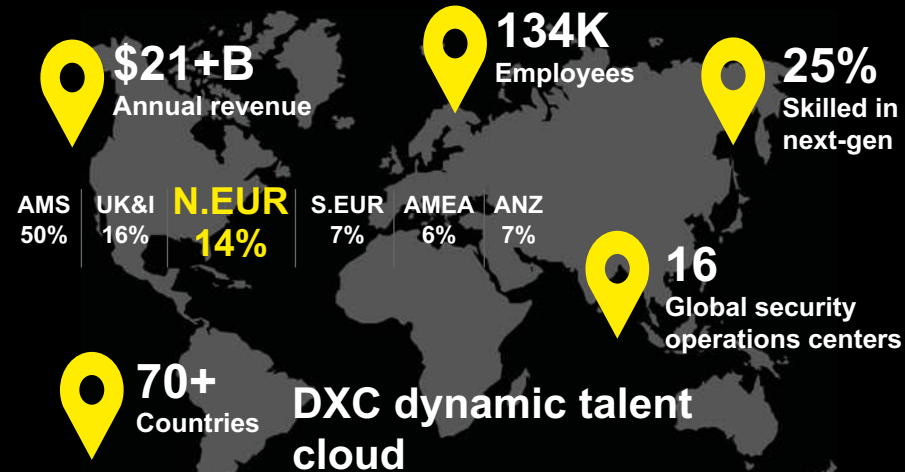
Amr Nouredin – Solution Architect
Michal Dura – Big Data Engineer



© 2019 DXC Technology, Proprietary and Confidential

The world's leading independent, end-to-end IT services company

Scale & Skills



DXC Value for AD



Accelerate time to market

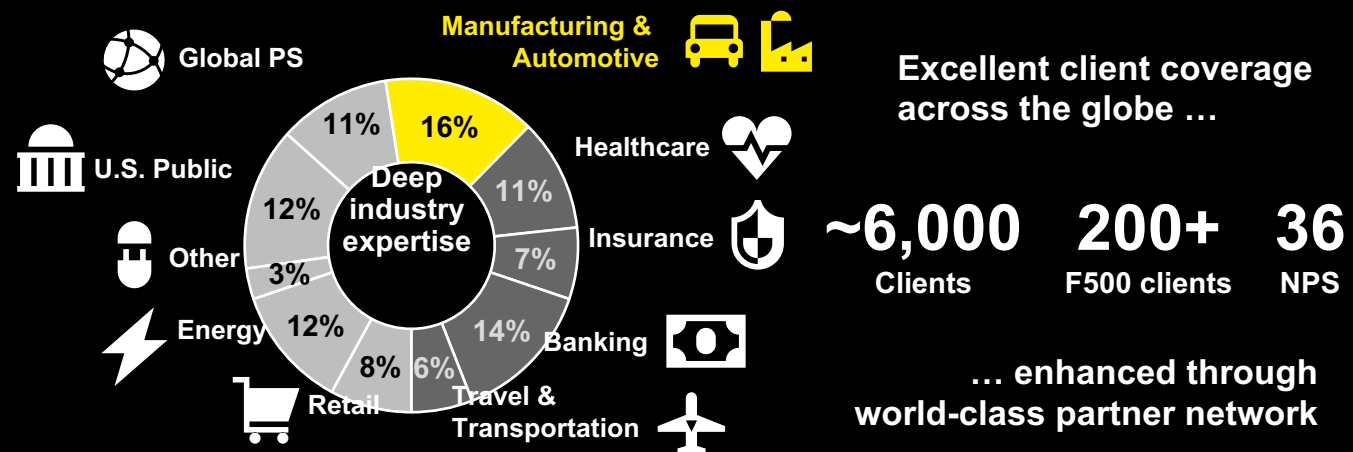
Reduce cost and risk

Improve market leadership



DXC.technology

Customer Intimacy



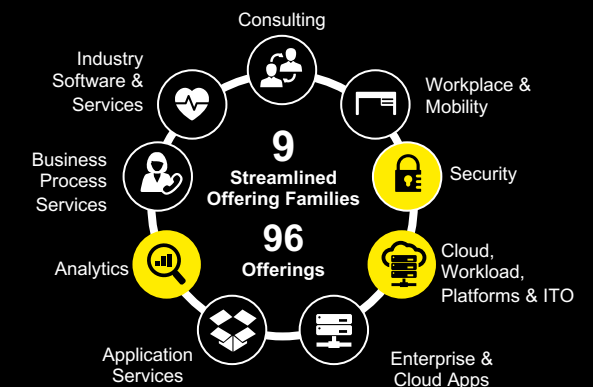
Technology-Driven Innovation

\$4B Digital revenue

250+ global partners

14 strategic co-investing partners

Streamlined offerings

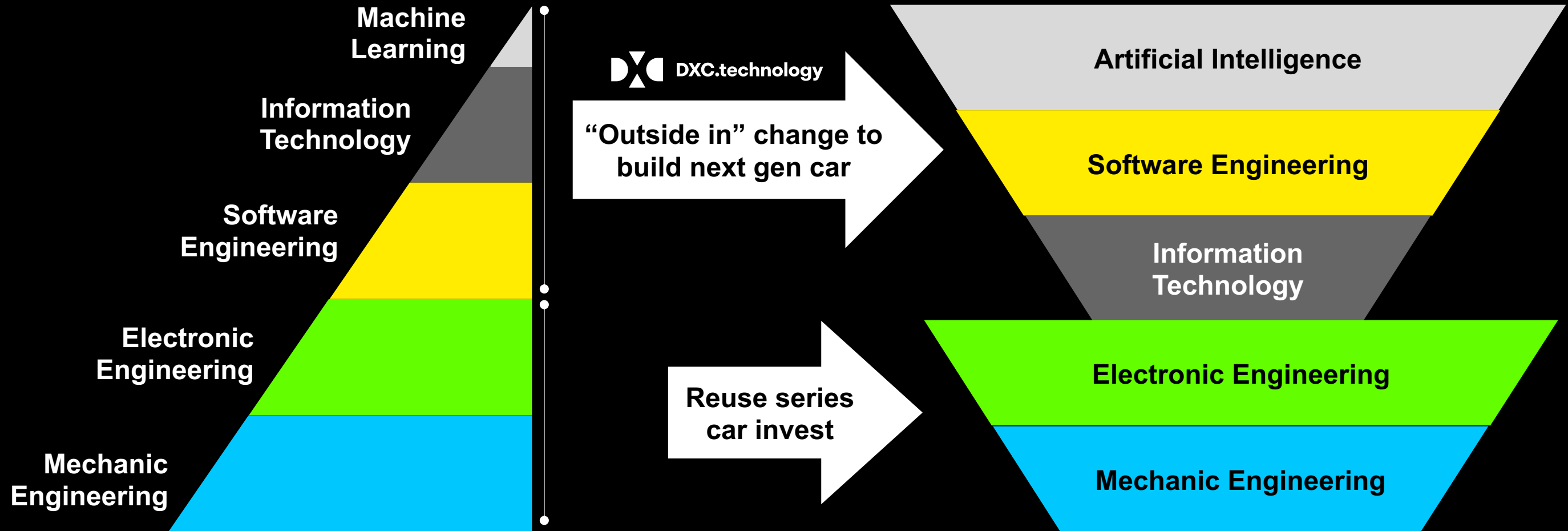


Autonomous Driving



Level - 0 DRIVER	Level - 1 FEET OFF	Level - 2 HANDS OFF	Level - 3 EYES OFF	Level - 4 MIND OFF	Level - 5 PASSENGER
No Assistance	Assisted	Partially Automated	Highly Automated	Fully Automated	Autonomous
Human	Transfer of responsibility				Machine

R&D in Automotive Industry – Capabilities are Changing



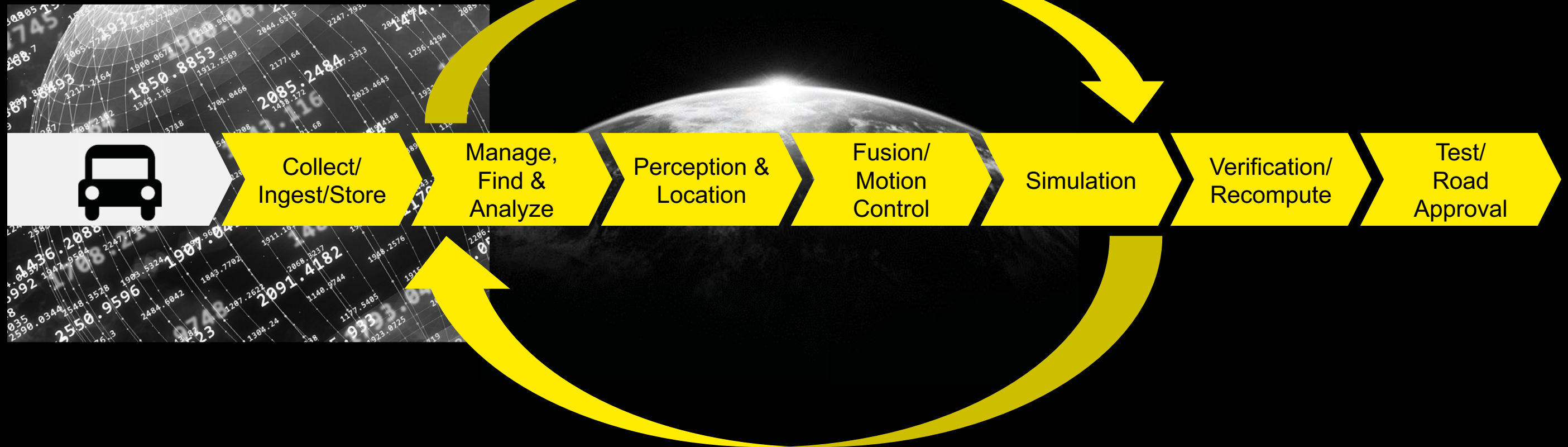
Extend & reuse current mechanic engineering capability combined with artificial intelligence, software building and IT is important to survive

Requires an end-to-end data and AI capability ecosystem for AD development

Geographically distributed
R&D teams

AD data & models

Need for speed

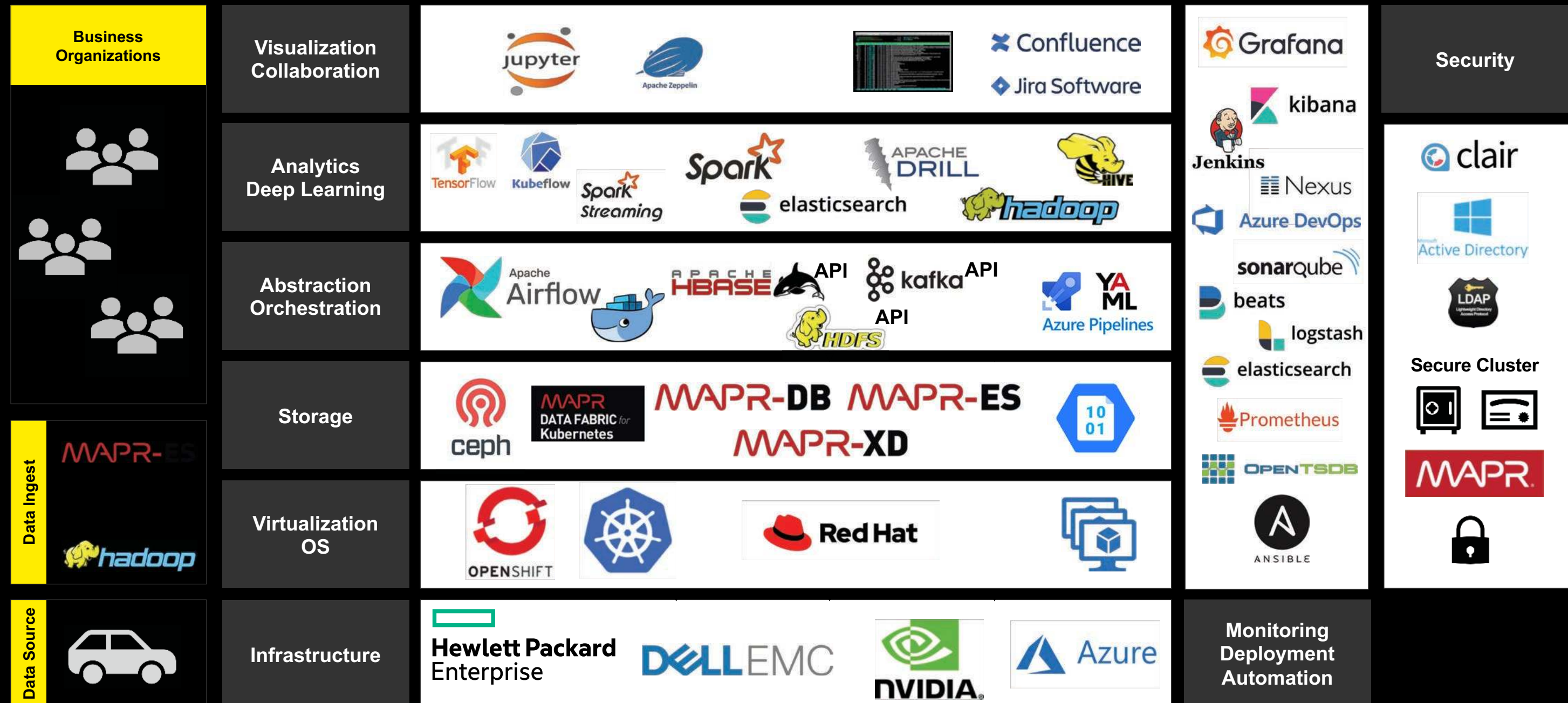


HIGH PERFORMANCE DATA DRIVEN DEVELOPMENT PLATFORM



FACTS & FIGURES

TECHNOLOGY STACK



Autonomous Driving with Airflow – In a Nutshell (1/3)

- **Airflow on OpenShift**
 - 1 scheduler instance – big risk
 - Multiple webserver – load balanced (Rest API Calls)
 - Numerous workers – multiple queues
- **Automated deployment via Helm charts**
 - Various configurations for different instances (also different Airflow versions)
 - History tracking via version control
- **On-demand Airflow instances (ex: development purposes, isolated testing – on a service level)**

Autonomous Driving with Airflow – In a Nutshell (2/3)

- **Integration with the Data and Storage Platform – MapR**
 - Loading DAGs from different locations
 - Loading job configurations, used by the different operators
- **Integration with the Compute Platform – OpenShift**
 - KubernetesPodOperator
 - KubernetesExecutor

Autonomous Driving with Airflow – In a Nutshell (3/3)

- **Metrics Collection & Monitoring**
 - StatsD → Prometheus → Grafana
- **Log collection and aggregation: ElasticSearch + Kibana**
- **Large scale orchestration: aiming at orchestrating jobs at the scale of 100,000's / month**
 - Ingestion, simulation, reprocesssing, machine learning, ...etc
 - Complex DAG dependencies

July 16, 2020

Apache Airflow - Robotic Drive orchestrator

Platform Orchestration Requirements



Open Source



Scalability



Easy to adapt / extend



Active community

What do we Orchestrate?



Data Ingestion




Machine Learning



Reprocessing



Simulation Jobs



July 16, 2020

Journey from PoC to Production

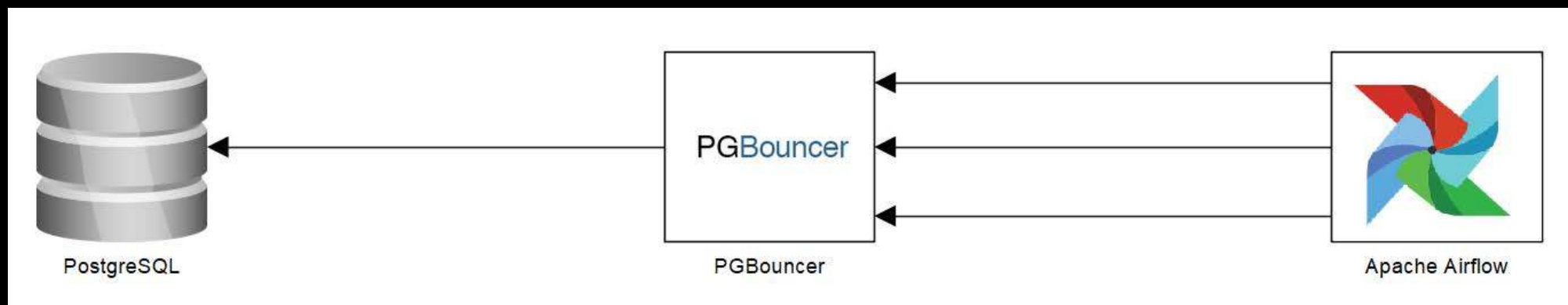
Airflow at Robotic Drive – the beginnings

- Initial work started at Q2 2019
- Airflow 1.10.2 with CeleryExecutor
- PostgreSQL 9.4
- RabbitMQ



Technical Challenges and Lessons Learned

- **Airflow stress and scalability tests**
- **Bottlenecks in the Architecture:**
 - PostgreSQL connection scalability: directly proportional relationship between number of running tasks and database connections
 - Scheduler configuration & performance



July 16, 2020

Tailor-made Solutions

Operators / Hooks Customizations (1/3)

Customization and standardization of SparkSubmitOperator

Included „properties_file” in operator constructor

- Spark application configuration can be provided via separate properties file
- It allows submitting jobs via Airflow, in the same fashion as submitting them standalone from the Hadoop cluster

Extend list of parameters where templating is supported

- Better DAGs reusability
- Reduced code duplication

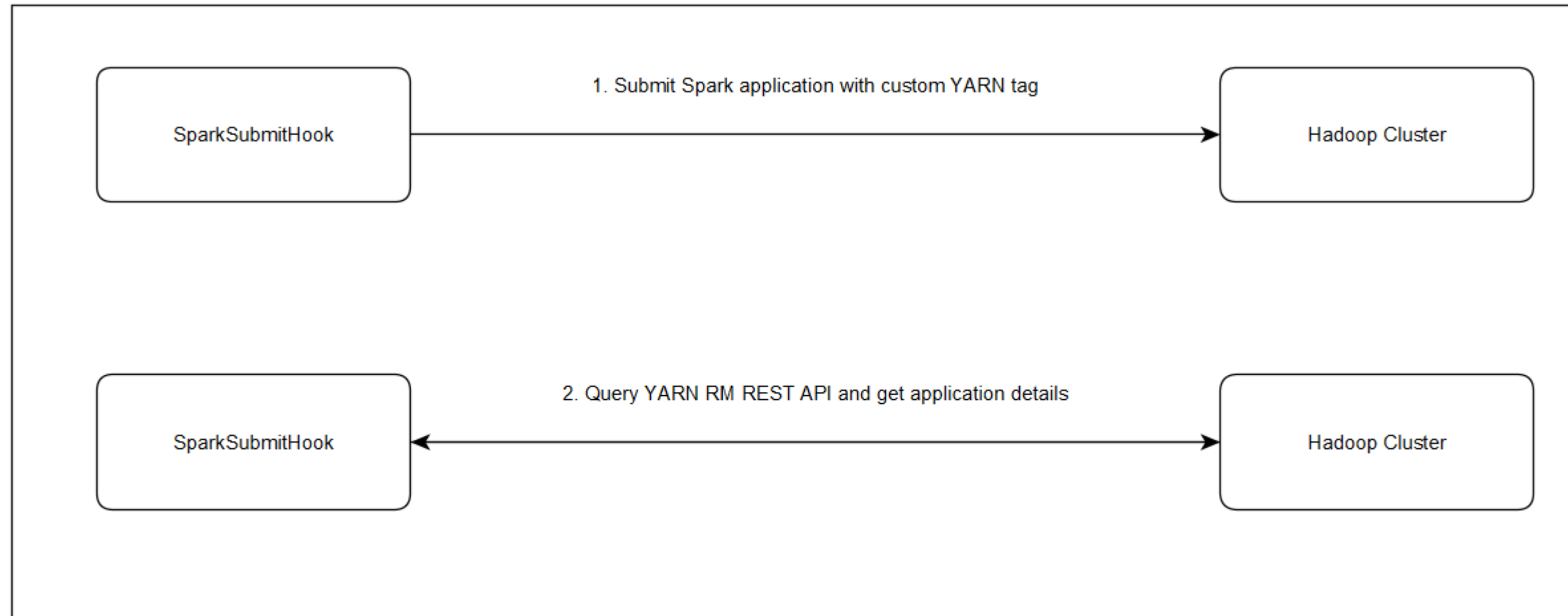
Operators / Hooks Customizations (2/3)

Enrich Airflow logs by adding YARN application details for Spark jobs

- Correlation between Spark application (triggered by Airflow) and submitted YARN job is challenging to discover when using YARN cluster mode
- Out of the box: YARN Application ID logged in the task logs only when a YARN job fails
- Extension: Extract and log the following for all Spark tasks:
 - 1) YARN Application ID
 - 2) YARN Tracking URL
 - 3) Diagnostics (Failure root cause)

Operators / Hooks Customizations (3/3)

YARN application details visible in Airflow logs for Spark jobs



Custom Authentication Methods

LDAP Secured REST API

- REST API usage allowed only for dedicated AD role
- Complete integration with LDAP
- Only one role specified for REST API – no separation between endpoints so far

July 16, 2020

Airflow in Production

Production-ready Airflow instances

Robotic Drive supports by default 3 main instances used in the platform:

- **Development** (used mostly to test new Airflow deployments / features)
- **Staging** (testing new DAGs)
- **Production** (for full production usage)

Current stable setup is created using **Airflow 1.10.10** and **Celery Executor**

User-based Airflow instances

- Deployment automated via **Helm Charts**
- Airflow created as a Kubernetes project

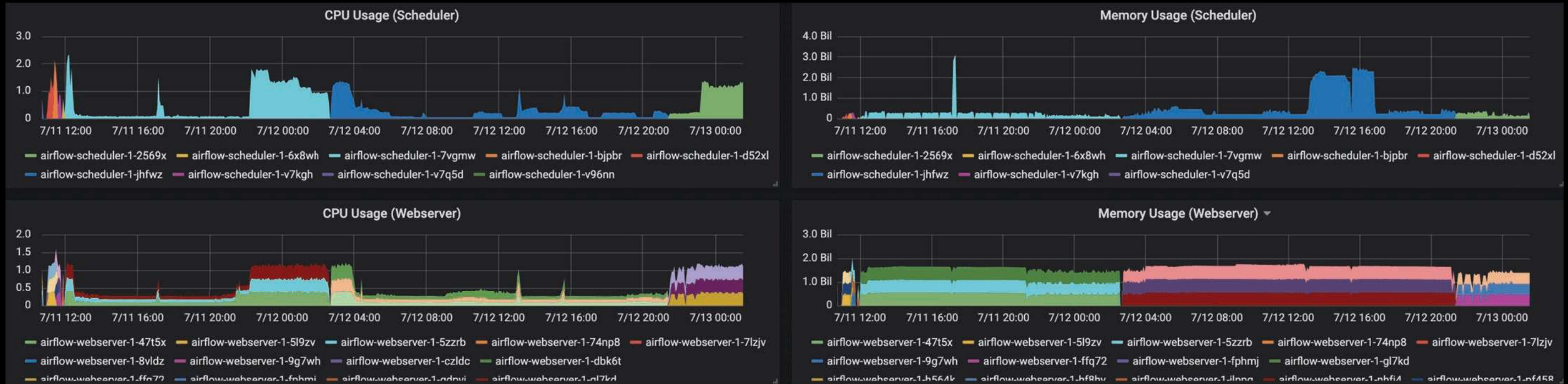
These 2 points makes it possible to fully parametrize Airflow deployment and create many Airflow instances on-demand. In the Robotic Drive Platform each user is able to create their own, separate Airflow instance.

It helps to eliminate problems related to testing new DAGs, Operators and other features and changes that might impact other users, especially when multiple developers are working on the same component.

July 16, 2020

Monitoring Airflow

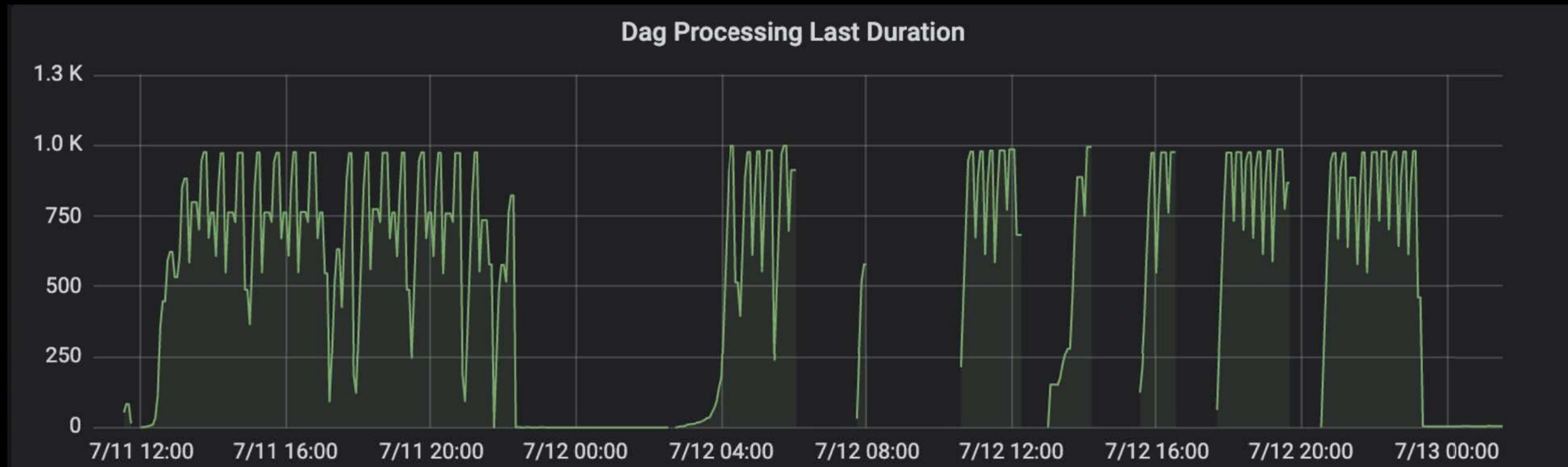
Monitoring Airflow



Monitoring Airflow



Monitoring Airflow



July 16, 2020

What's next?

Looking forward to...

- **Airflow 2.0: HA Scheduler + Performance Optimizations**
- **Advanced Authentication + Authorization**
- **Extend and stabilize monitoring metrics**
- **Stable API vs Experimental API**

Q&A

July 16, 2020

Amr Noureldin <amr.noureldin@dxccom>
Michal Dura <mdura2@dxccom>