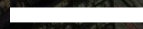


Towards convergence of Big Data and HPC considering hybrid edge-cloud infrastructures



Yiannis Georgiou - CTO Ryax Technologies



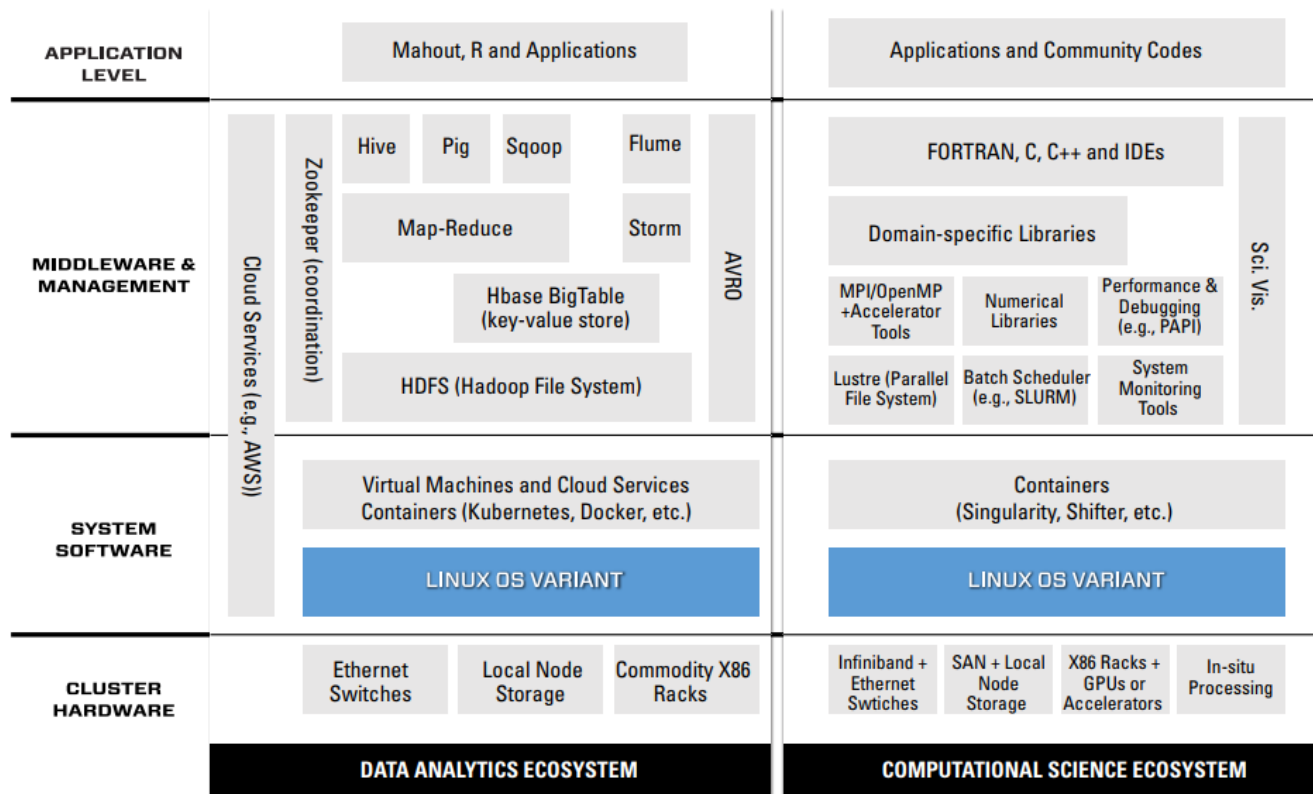
Big Data and Exascale Community Report

- The latest BDEC community study[1] reports on latest research, challenges and provides **future directions upon the convergence** of infrastructures. We focus on the following two aspects:
- The importance of **HPC and Big Data convergence** in terms of **resource management and scheduling**
- The importance of **edge computing and decentralized facilities for processing closer to sources**



[1]Tech Report: BIG DATA AND EXTREME-SCALE COMPUTING: PATHWAYS TO CONVERGENCE
Toward a Shaping Strategy for a Future
Software and Data Ecosystem for Scientific Inquiry

HPC and Big Data Stacks



[1]Tech Report: BIG DATA AND EXTREME-SCALE COMPUTING: PATHWAYS TO CONVERGENCE
 Toward a Shaping Strategy for a Future
 Software and Data Ecosystem for Scientific Inquiry

Resource Management Convergence

- More research effort is needed for converged resource and execution management with **radically improved** centralized intelligence.
- Proposed solution provide **higher level schedulers** to communicate with multiple types of resource managers and schedulers that are **specialized** for the particular hardware or ecosystem.

Resource Management Convergence

- Research and tools to deploy mixed Big Data and HPC workloads:
 - **Different clusters** dedicated to each workload. But, data transfer and load balancing are challenging.
 - Let the user deploy the Big Data workload inside HPC batch jobs using a **set of scripts** [1].
 - Use **pilot-based abstraction**[2] to deploy and manage Hadoop or stream-processing (Spark, Flink) frameworks upon HPC infrastructure.
 - A **lightweight and less interfering** approach [3] where execution of Big Data applications is done by the HPC scheduler as HPC “best-effort” jobs .
 - INDIGO-Datacloud project[4] uses the **udocker runtime tool** along with Kubernetes **orchestration** to enable the deployment of both HPC and Big Data workloads.

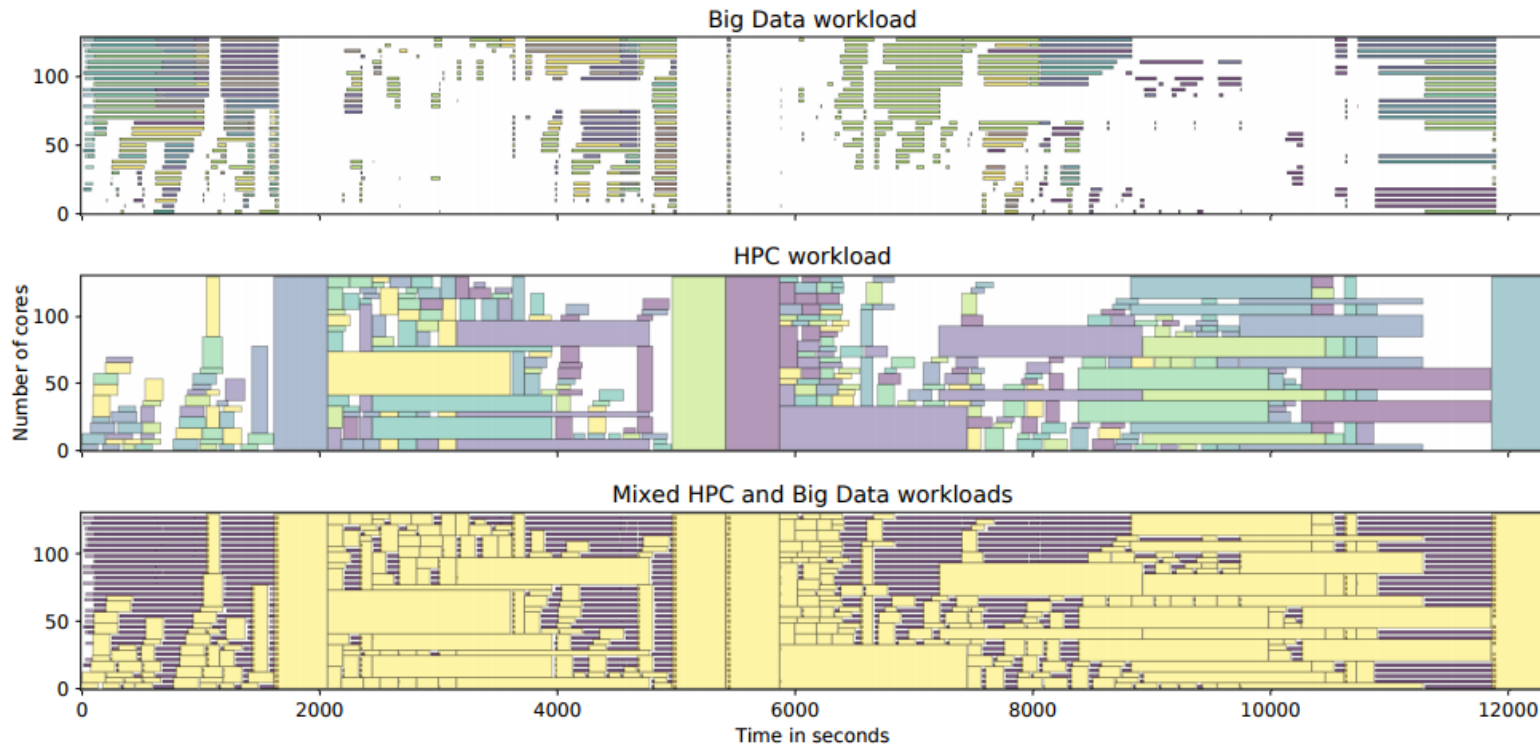
[1]<https://github.com/LLNL/magpie>

[2]André Luckow, George Chantzialexiou, Shantenu Jha: Pilot-Streaming: A Stream Processing Framework for High-Performance Computing. CoRR abs/1801.08648 (2018)

[3]Michael Mercier, David Glesser, Yiannis Georgiou, Olivier Richard: Big data and HPC collocation: Using HPC idle resources for Big Data analytics. BigData 2017: 347-352

[4]D. Salomoni et al. "INDIGO-DataCloud: Project Achievements", arXiv:1711.01981

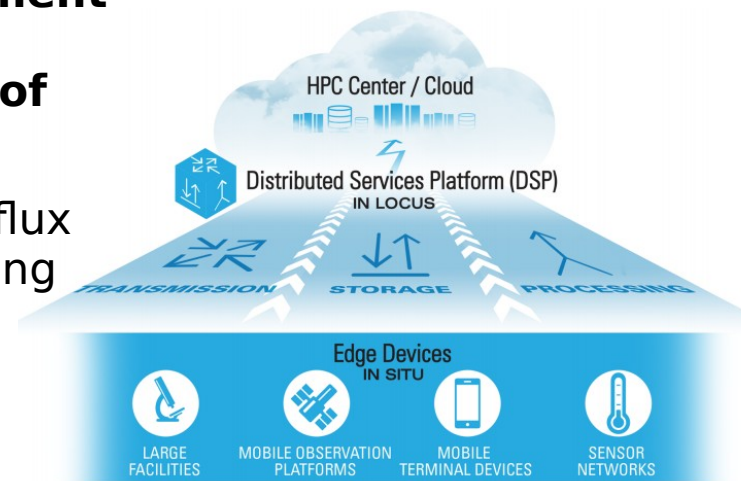
HPC and Big Data collocation Study



Michael Mercier, David Glesser, Yiannis Georgiou, Olivier Richard: Big data and HPC collocation: Using HPC idle resources for Big Data analytics. BigData 2017: 347-352

Towards Hybrid Edge-Cloud

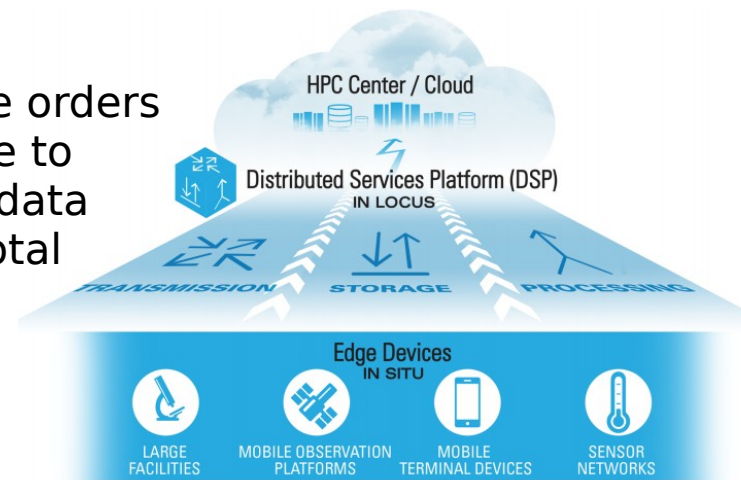
- The explosive growth and dispersion of digital data producers in **edge environments** creates various **challenges**.
- Big instruments such as LHC and the Argonne Photon Source (APS) illustrate how **managing the movement and staging of data** from where it is collected to where it needs to be analyzed **can take up most of the time to solution**.
- Neuroimaging and genetics data shows that the influx is doubling every year, with some estimates reaching **over 20 petabytes per year** by 2019



[1]Tech Report: BIG DATA AND EXTREME-SCALE COMPUTING: PATHWAYS TO CONVERGENCE
Toward a Shaping Strategy for a Future
Software and Data Ecosystem for Scientific Inquiry

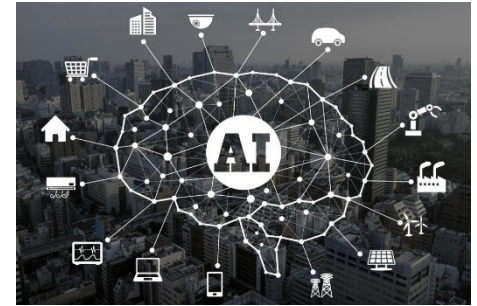
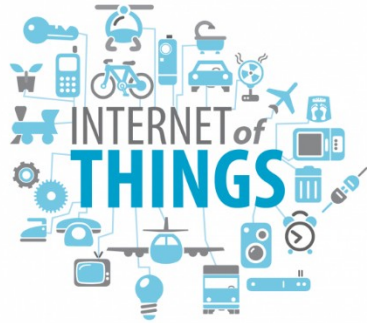
Towards Hybrid Edge-Cloud

- **Light Detection And Ranging (LIDAR)** survey technology routinely produces **terabyte-level datasets**, with huge cumulative volumes.
- **Autonomous vehicles** will generate and consume roughly **40 terabytes of LIDAR** data for every **8 hours of driving**, and LIDAR prices have come down 3 orders of magnitude in 10 years
- **Mobile data traffic** has increased more than three orders of magnitude over the last decade and will continue to **grow at more than 50% annually**; by 2020 this data traffic is projected to surpass **500 zettabytes** in total



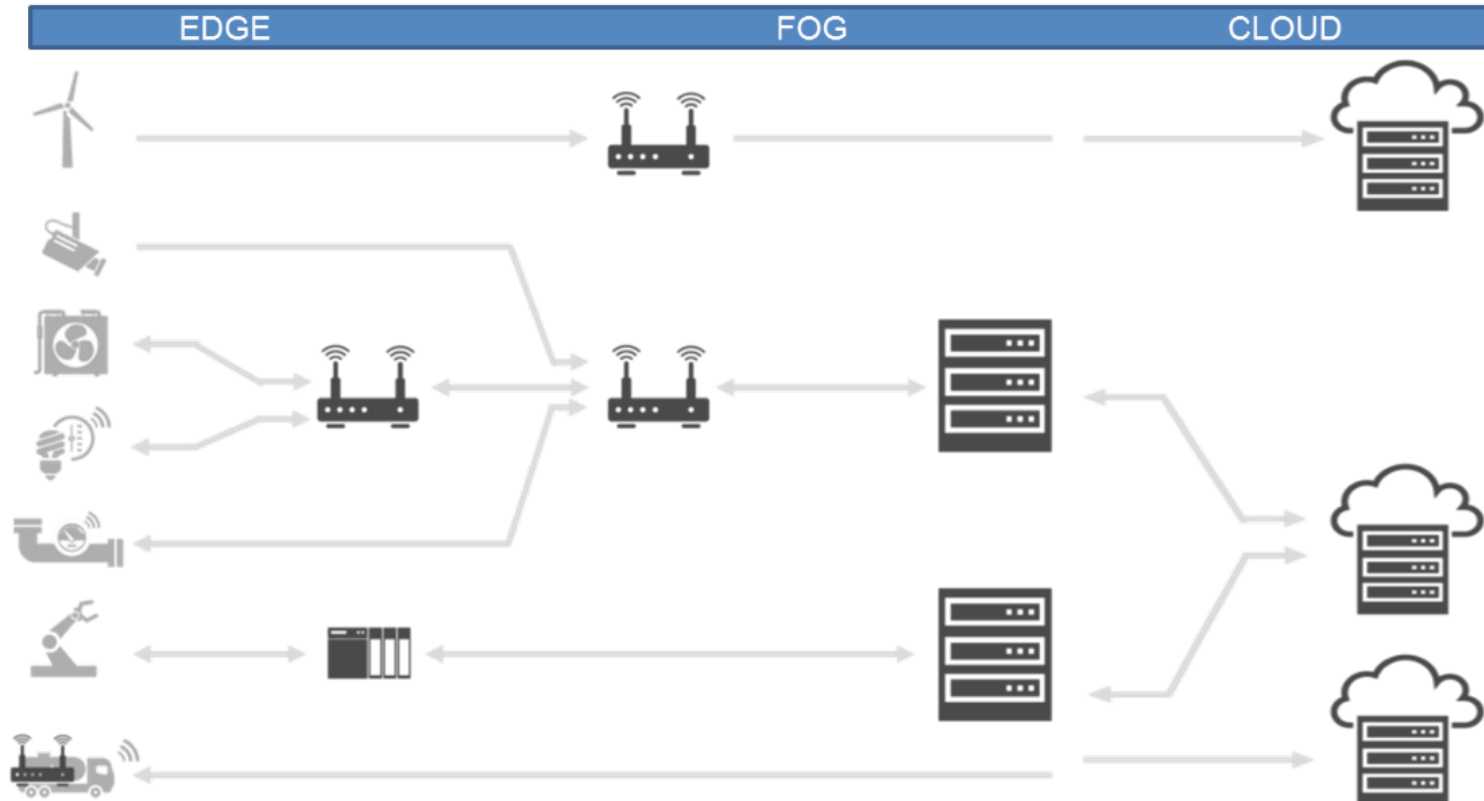
[1]Tech Report: BIG DATA AND EXTREME-SCALE COMPUTING: PATHWAYS TO CONVERGENCE
Toward a Shaping Strategy for a Future
Software and Data Ecosystem for Scientific Inquiry

Towards Hybrid Edge-Cloud



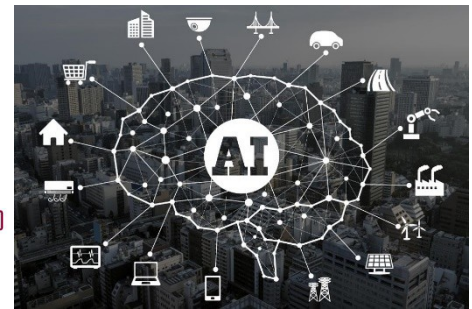
- For reasons such as **cost** of data transfers, needs of **low latency**, data **locality**, data **privacy**, Cloud, HPC and centralized facilities are not enough.
- Compute/Analyze data closer to sources leveraging on **Edge and Fog Computing**.

Hybrid Infrastructures

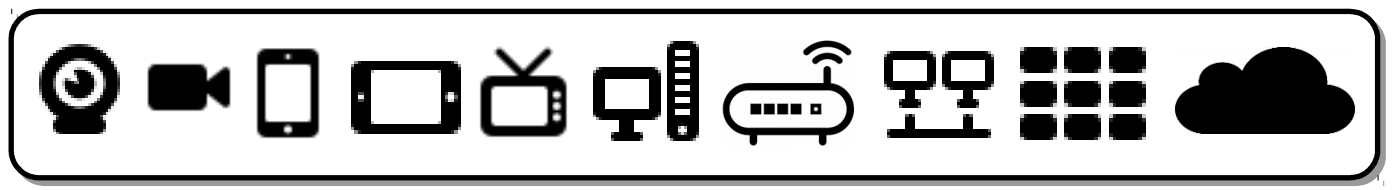


Applications Orchestration on Hybrid Infrastructures

Increasing needs for Compute and Data Intensive Applications



Move towards Hybrid Infrastructures



Edge

Fog

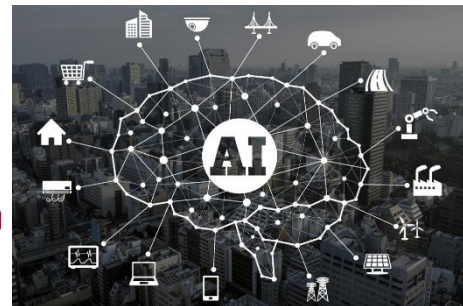
Cloud

Challenges for Orchestration on Hybrid Infrastructures

- **Network complexity / Hardware Heterogeneity** : How to manage them within hybrid edge/ fog/ cloud environments?
- **Seamless code execution** : How to execute applications without needing to develop each part differently depending on where it is going to be executed?
- **Deployment and Monitoring** : How to deploy environments and monitor resources even on resources with low capabilities of compute/memory?
- **Multi-level, multi-user data privacy** : How to manage it on a distributed system potentially owned by different owners and used by big number of users?
- **Isolation, security, billing**: How to guarantee data transfers security, isolation of computations and precise accounting/billing based on resources consumption?
- **Content-aware, offline support, scalability**, etc

Applications Orchestration on Hybrid Infrastructures

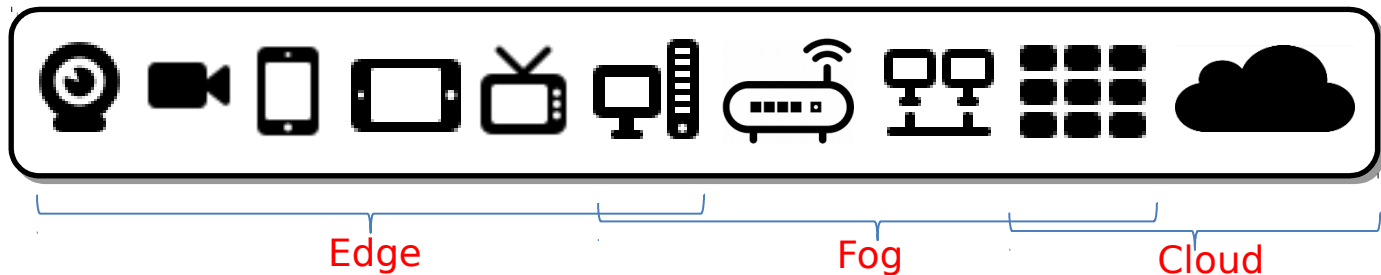
Increasing needs for
Compute and Data
Intensive Applications



Challenges:

- Network complexity and Hardware Heterogeneity
- Programming and Executing
- Multi-level, multi-user data privacy
- Isolation, security, billing
- Content-aware, offline support, scalability

Move towards
Hybrid
Infrastructures

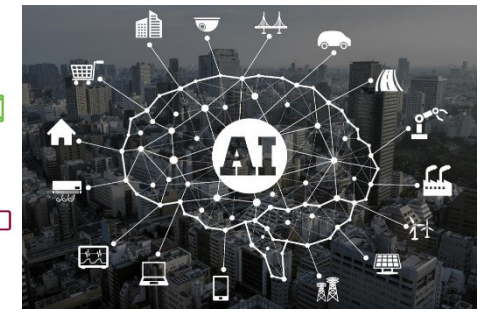


Our product - Ryax

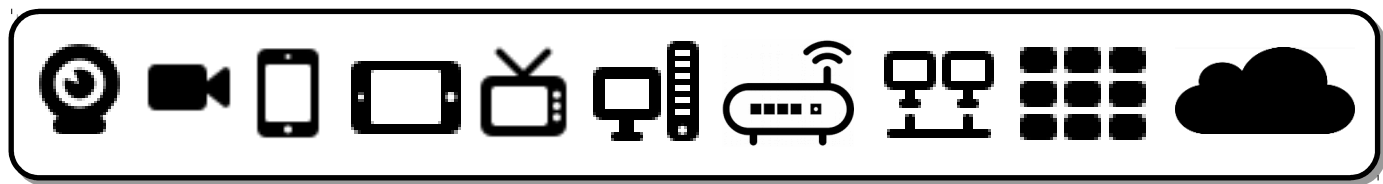
We propose a **middleware** to deal with the above challenges and provide the **Applications Orchestration** and **Compute management** for **hybrid Edge-Cloud infrastructures** to facilitate the deployment of **Compute and Data Intensive Workloads**

Applications Orchestration on Hybrid Infrastructures

Increasing needs for Compute and Data Intensive Applications



Move towards Hybrid Infrastructures

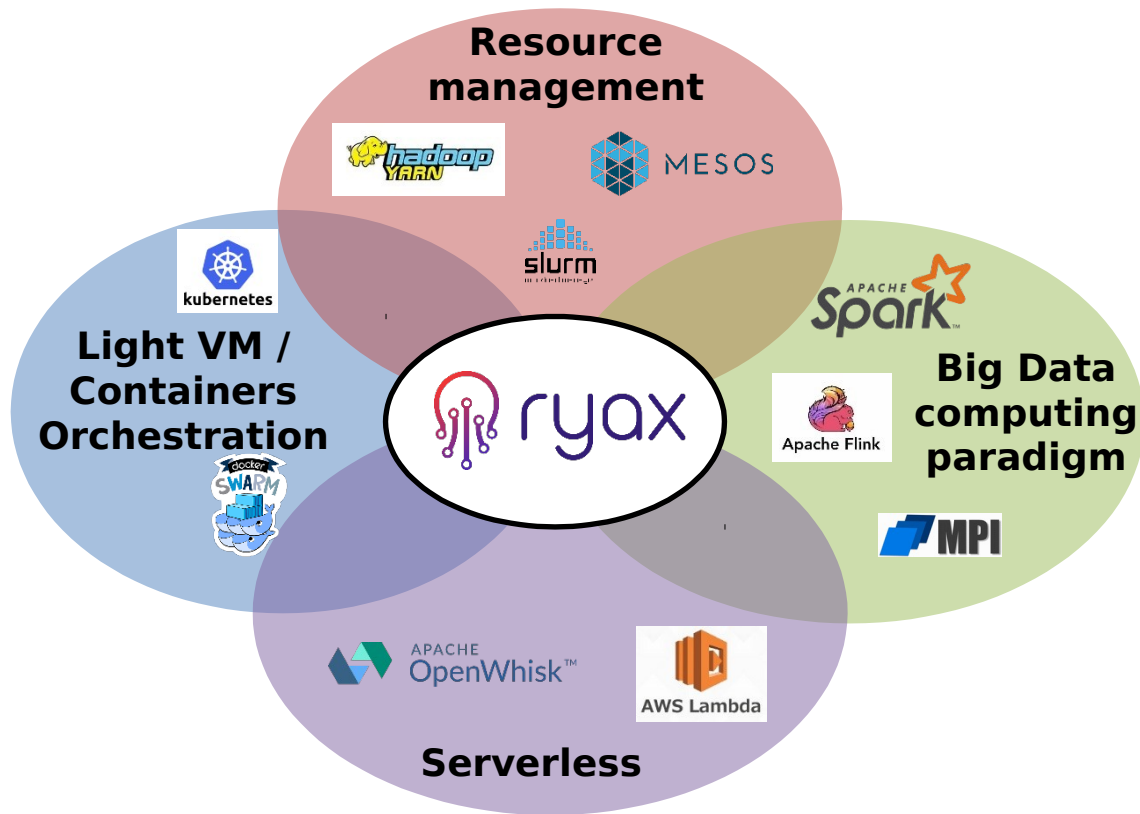


Edge

Fog

Cloud

Technological landscape



To efficiently manage and compute on Hybrid Edge-Cloud Infrastructures we need to **tightly integrate functionalities** emerging from the 4 pillars of modern stacks.

Our solution - overview

Resource and workload management for hybrid Edge-Cloud environments



Hybrid Environments Edge-Fog-Cloud

Resource Management on hybrid networks of private nodes, mini data-centers, public gateways, clouds.



Performance

Low overhead, smart allocation decisions and data flow dynamic optimizations.



Ease of use

Seamless deployment and management of containers and bare metal software.
Support multiple programming languages.

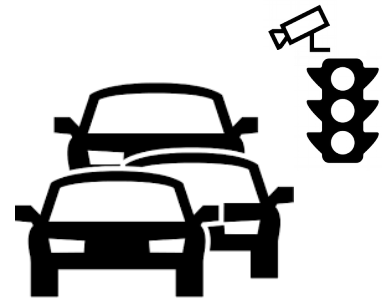


Security and Privacy Management

Integrated security and data privacy management with data tracking and usage limitation.

Use Case Scenario – Smart Traffic Lights

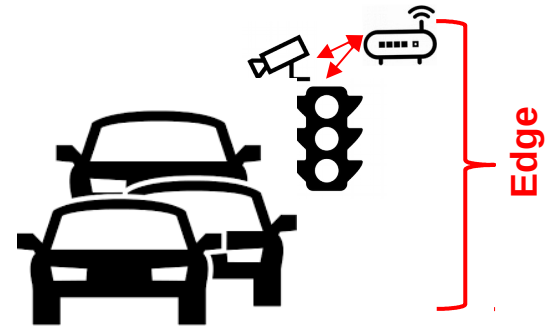
- 1) Single **crossroad** managed by a traffic light.
- 2) Traffic light equipped with **IoT light state** sensors and **cameras**.
- 3) Data **streams** of light states and visual evidence of congestion.
- 4) First level of data **analytics** to merge data streams of light states and cameras.



Use Case Scenario – Smart Traffic Lights

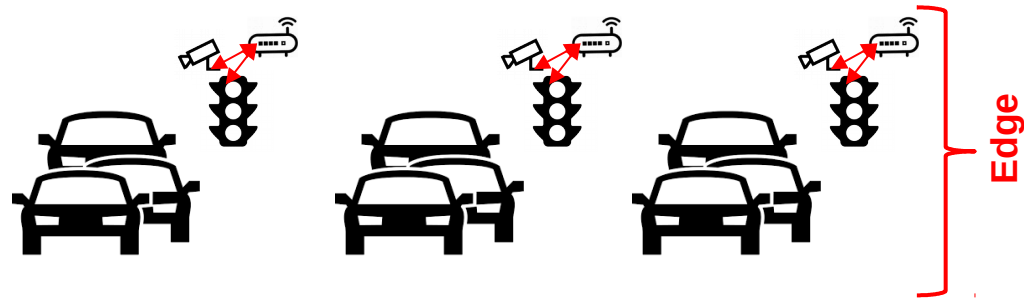
- 1) Single **crossroad** managed by a traffic light.
- 2) Traffic light equipped with **IoT light state** sensors and **cameras**.
- 3) Data **streams** of light states and visual evidence of congestion.
- 4) First level of data **analytics** to merge data streams of light states and cameras.

Use of **edge** infrastructure (**on-board gateway**) for fast response time, simple data analytics no need for important compute power.



Use Case Scenario – Smart Traffic Lights

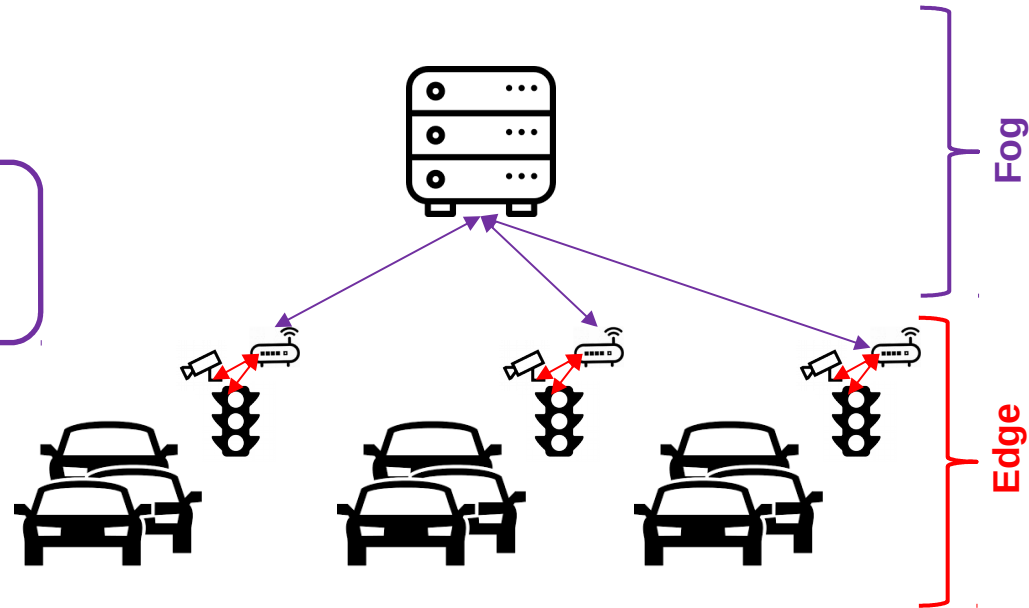
- 5) Multiple **crossroads** managed by smart traffic lights.
- 6) Need to **combine** output from first level of data analytics with **traffic flow model** to develop aggregated knowledge for a whole area of roads and **adapt** traffic lights accordingly.



Use Case Scenario – Smart Traffic Lights

- 5) Multiple **crossroads** managed by smart traffic lights.
- 6) Need to **combine** output from first level of data analytics with **traffic flow model** to develop aggregated knowledge for a whole area of roads and **adapt** traffic lights accordingly.

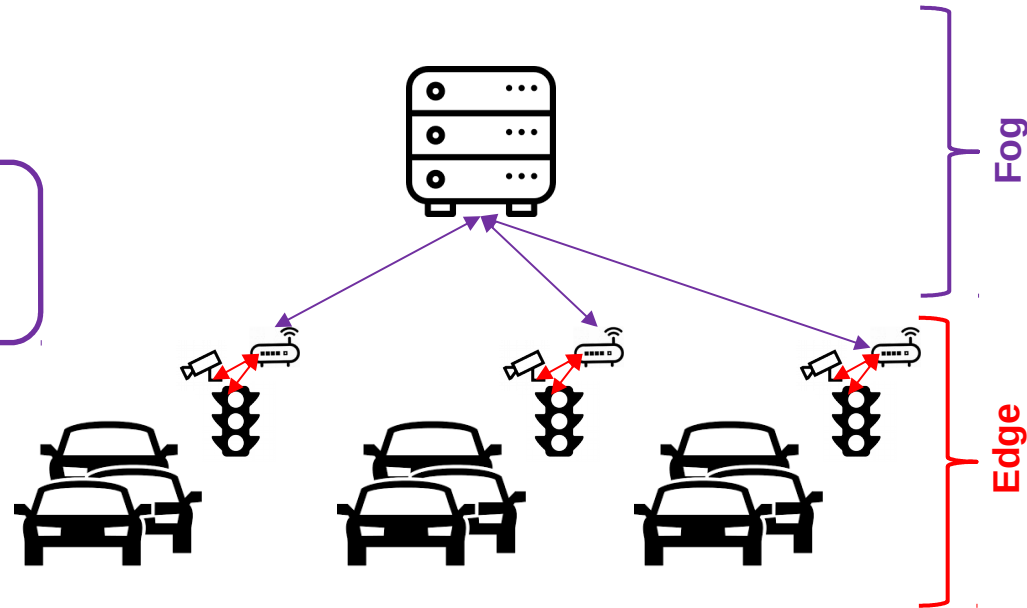
Use of local **fog** infrastructure (**nearby mini-datacenter**) for complex data analytics that need compute power.



Use Case Scenario – Smart Traffic Lights

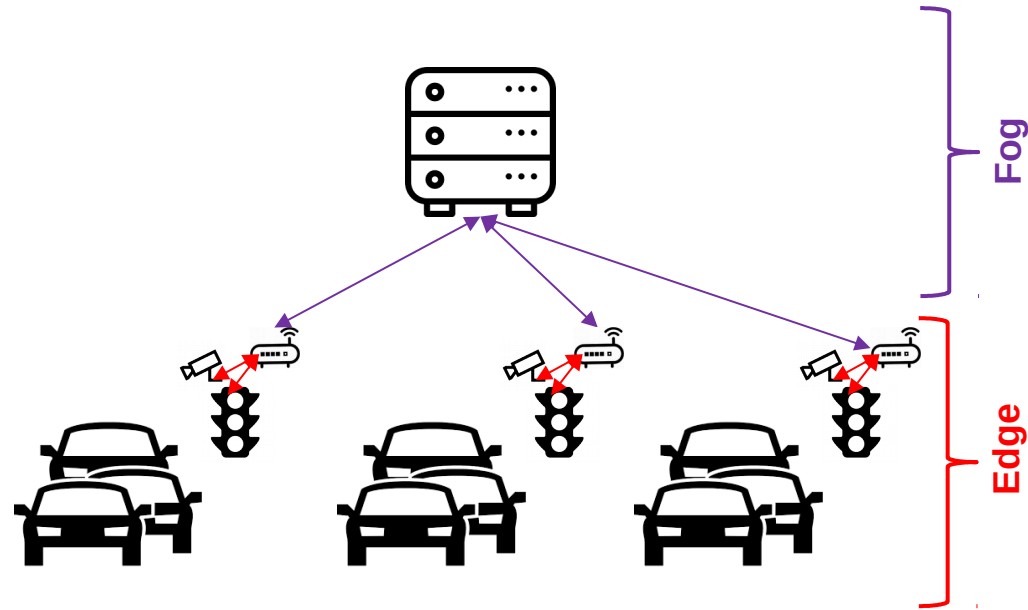
- 5) Multiple **crossroads** managed by smart traffic lights.
- 6) Need to **combine** output from first level of data analytics with **traffic flow model** to develop aggregated knowledge for a whole area of roads and **adapt** traffic lights accordingly.

Use of local **fog** infrastructure (**nearby mini-datacenter**) for complex data analytics that need compute power.



Use Case Scenario – Smart Traffic Lights

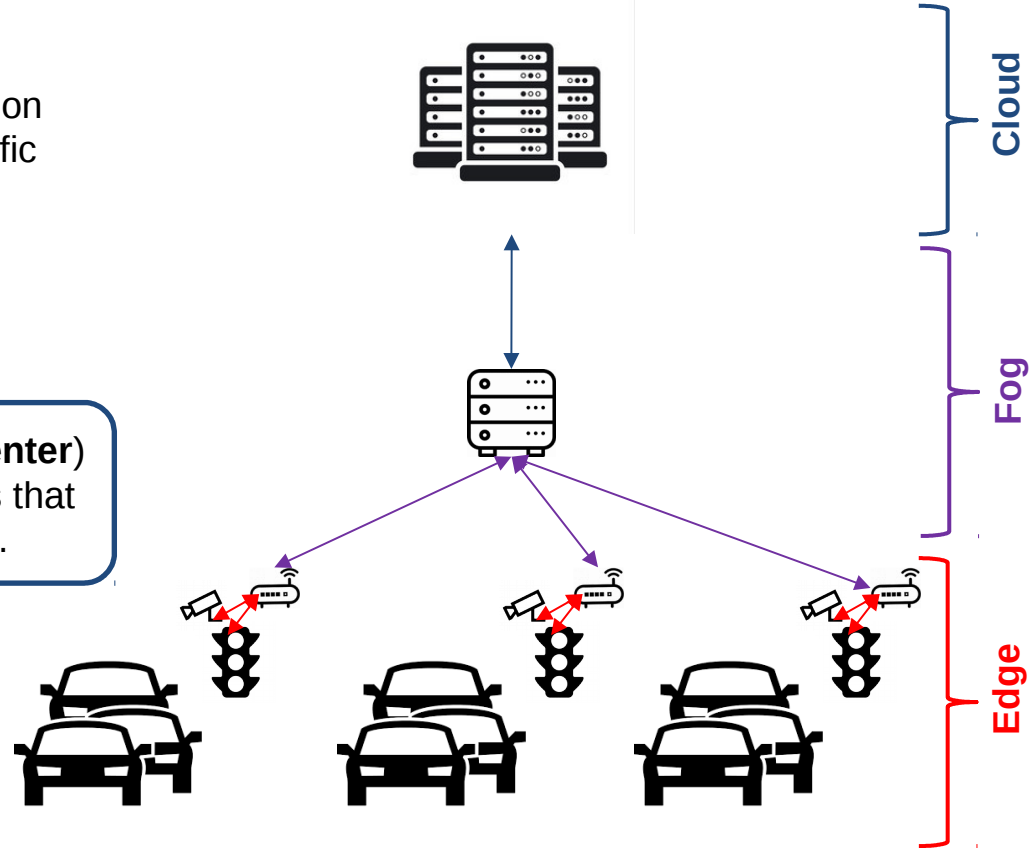
- 7) Need to store particular aggregated information or perform further analytics demanding specific type of computational resources.



Use Case Scenario – Smart Traffic Lights

- 7) Need to store particular aggregated information or perform further analytics demanding specific type of computational resources.

Use of remote **cloud** infrastructure (**datacenter**) for storing or highly complex data analytics that need specific computational resources.



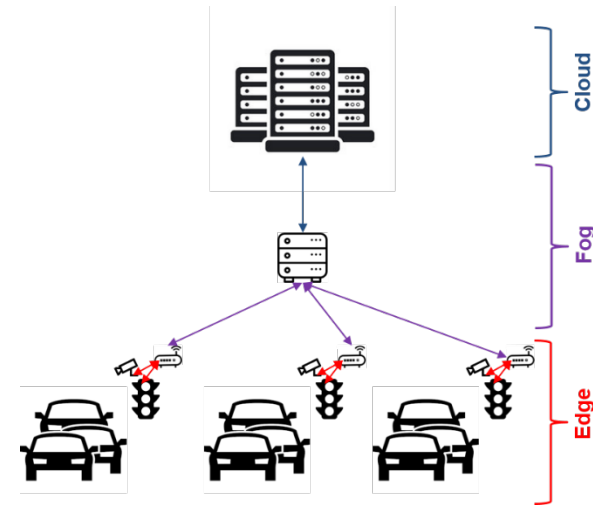
Use Case Scenario – Smart Traffic Lights

Solutions with Ryax

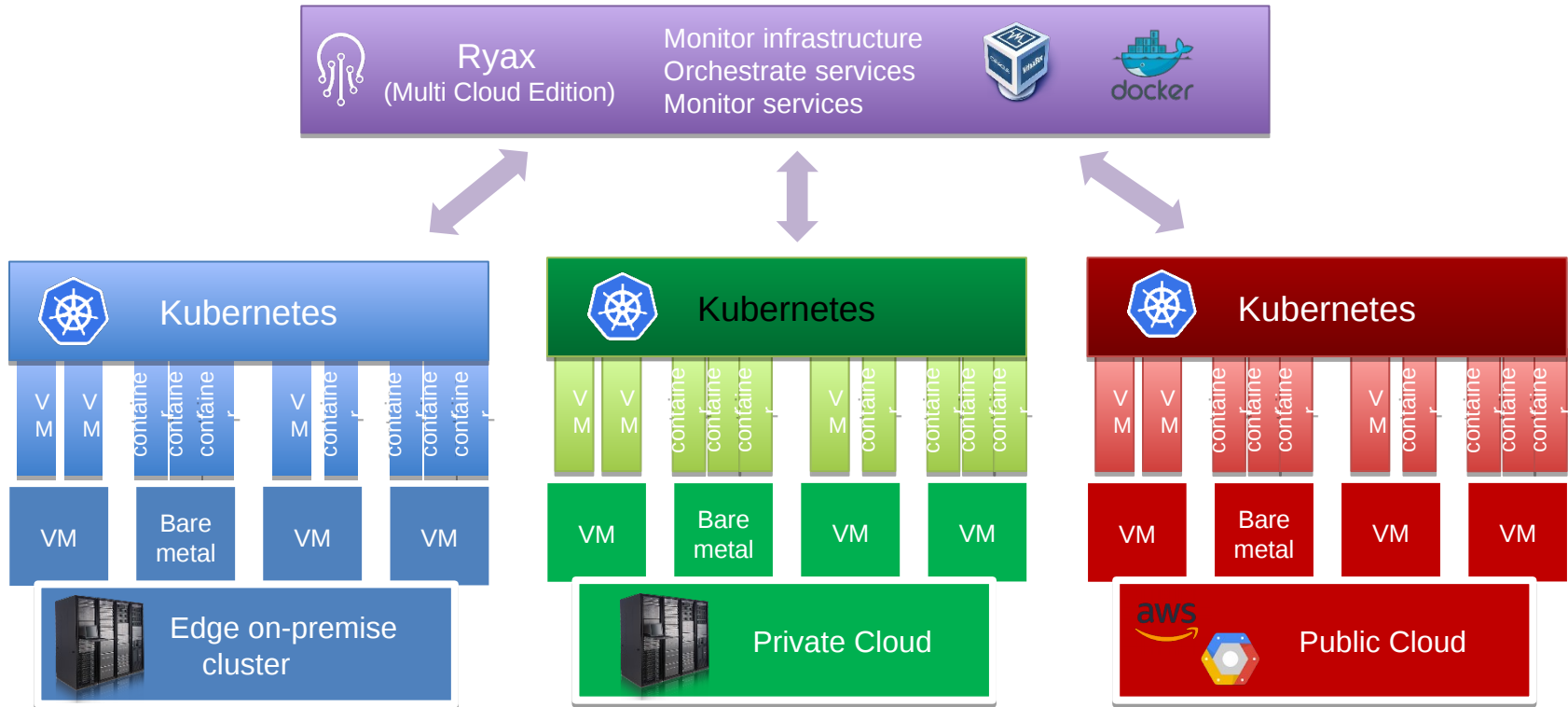
- Low overhead and low latency resource management.
- Seamless data stream processing across edge, fog, cloud.
- Optimized resource allocation based on different objectives.
- Adapted Privacy Management.
- Functions offline.
- Accurate billing per application based on resource consumption.

Benefits with Ryax

- Data streams can be used for multiple applications (i.e. traffic flow data used for environmental impact monitoring).
- Take advantage of the advanced isolation, SDN, efficiency and billing functionalities offered by Ryax.
- Edge/Fog gateways and mini-datacenters can be operated by Ryax for efficient usage of computational resources by smart-city applications.

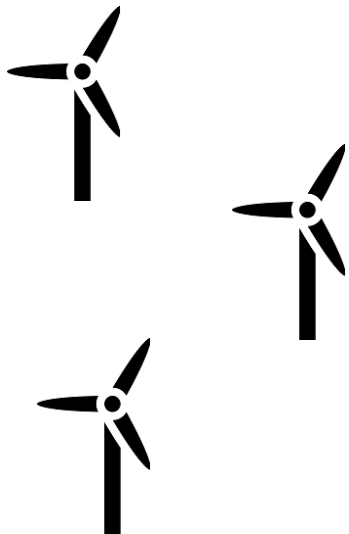


Ryax Multi-Cloud Usage



Demo Windmills Use Case

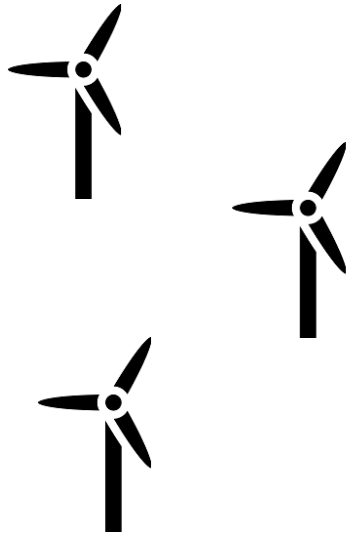
This demo presents a windmills use case on global edge-fog-cloud hybrid infrastructure



Main fonctionnalites demonstrated:

- Deployment upon a set of distributed clusters
- Centralized CLI and graphical user interface
- Advanced scheduling technique

Product demo - Windmills



- Cloud
- Monthly reports



Global efficiency?

- On premise
- Detailed view
- Predictive maintenance



Everything works?



Edge cluster

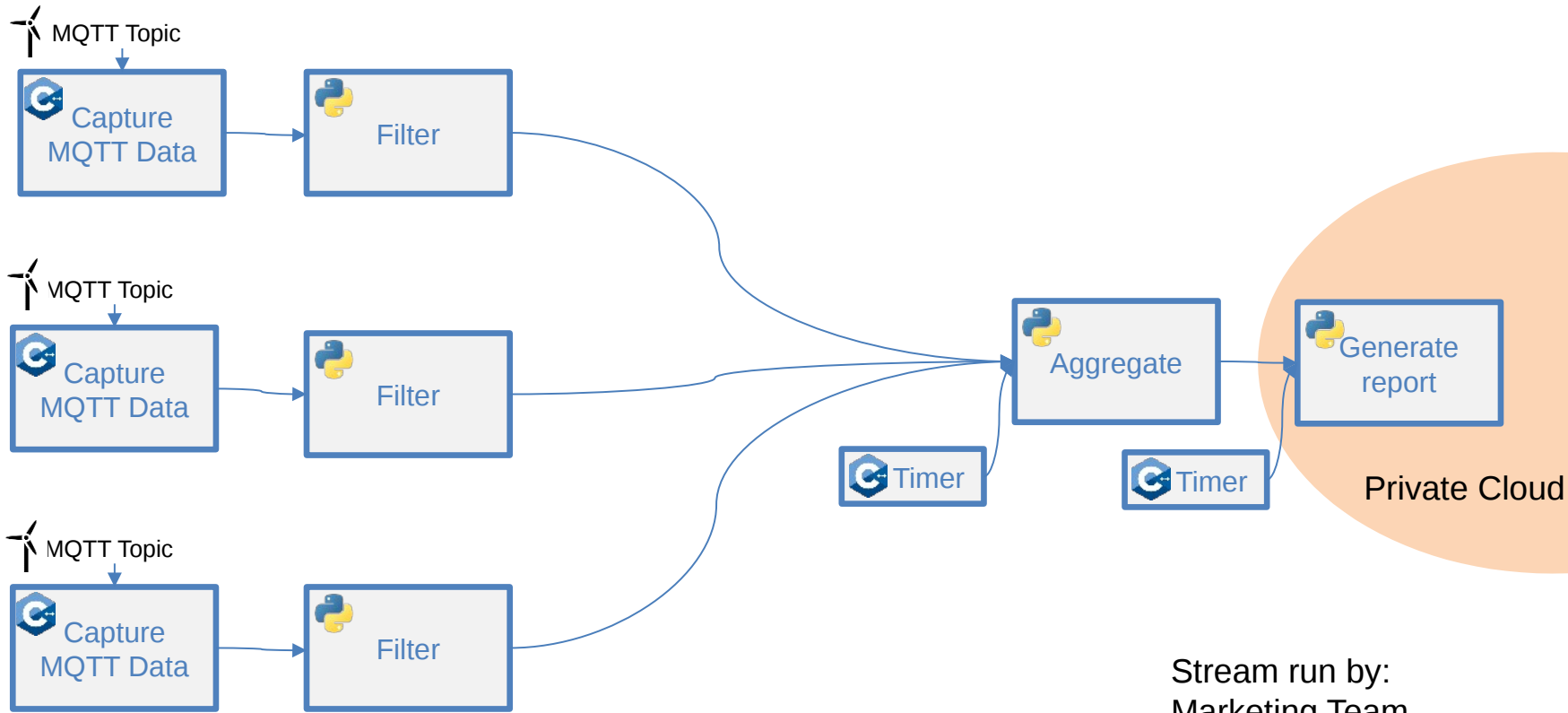


Private cloud



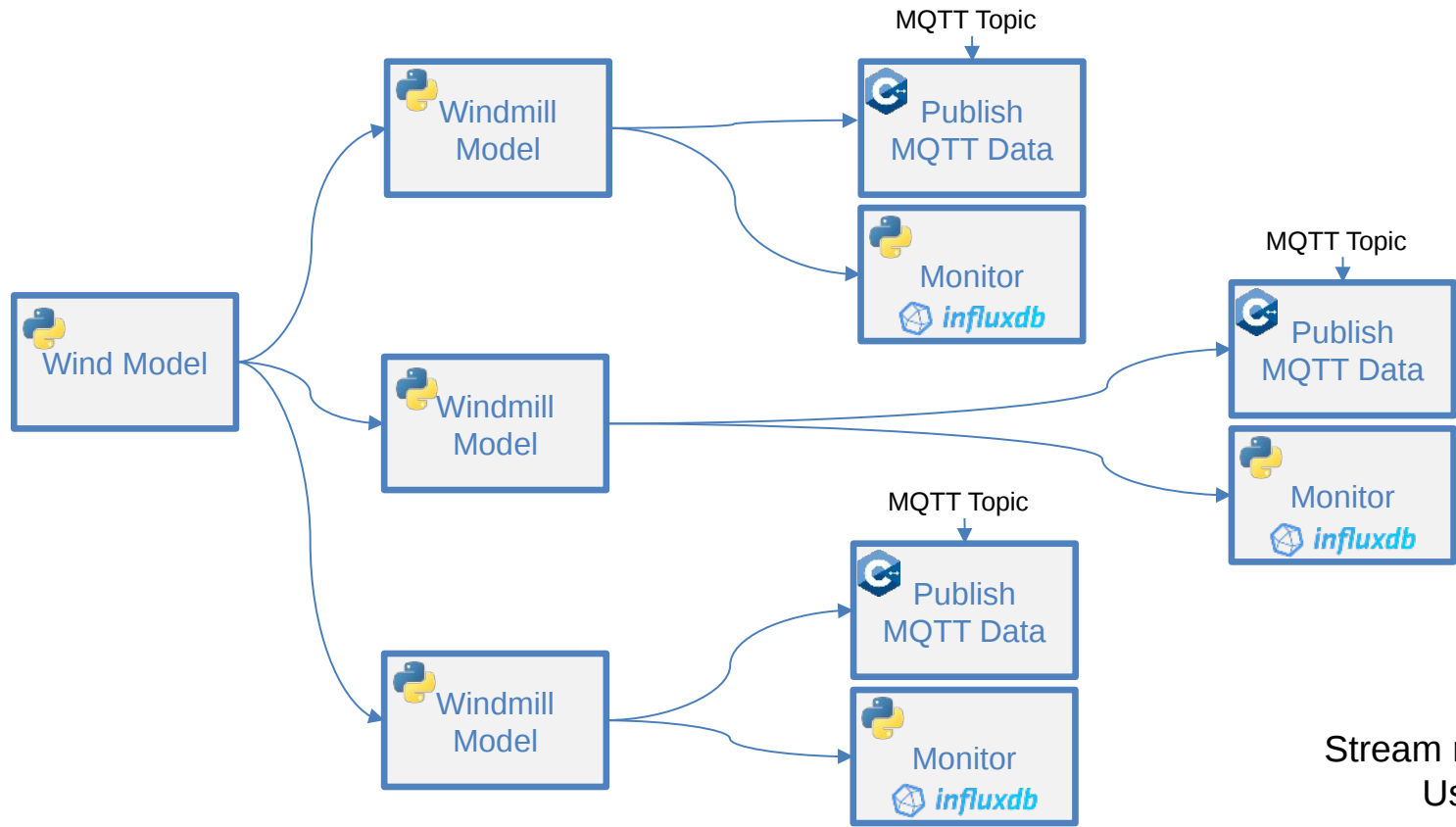
Public cloud

Product demo - Windmills



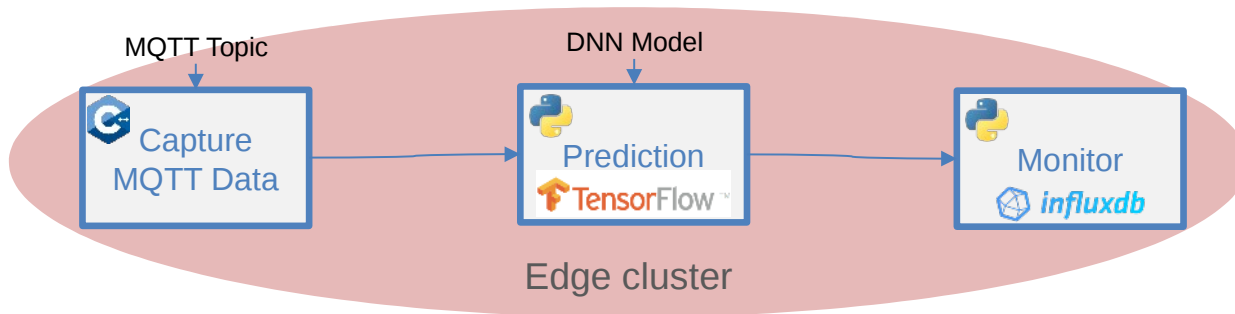
Stream run by:
Marketing Team

Product demo - Windmills



Stream run by:
Us

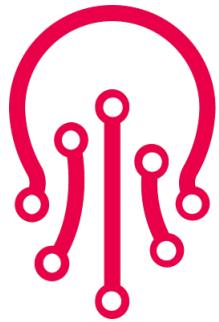
Product demo - Windmills



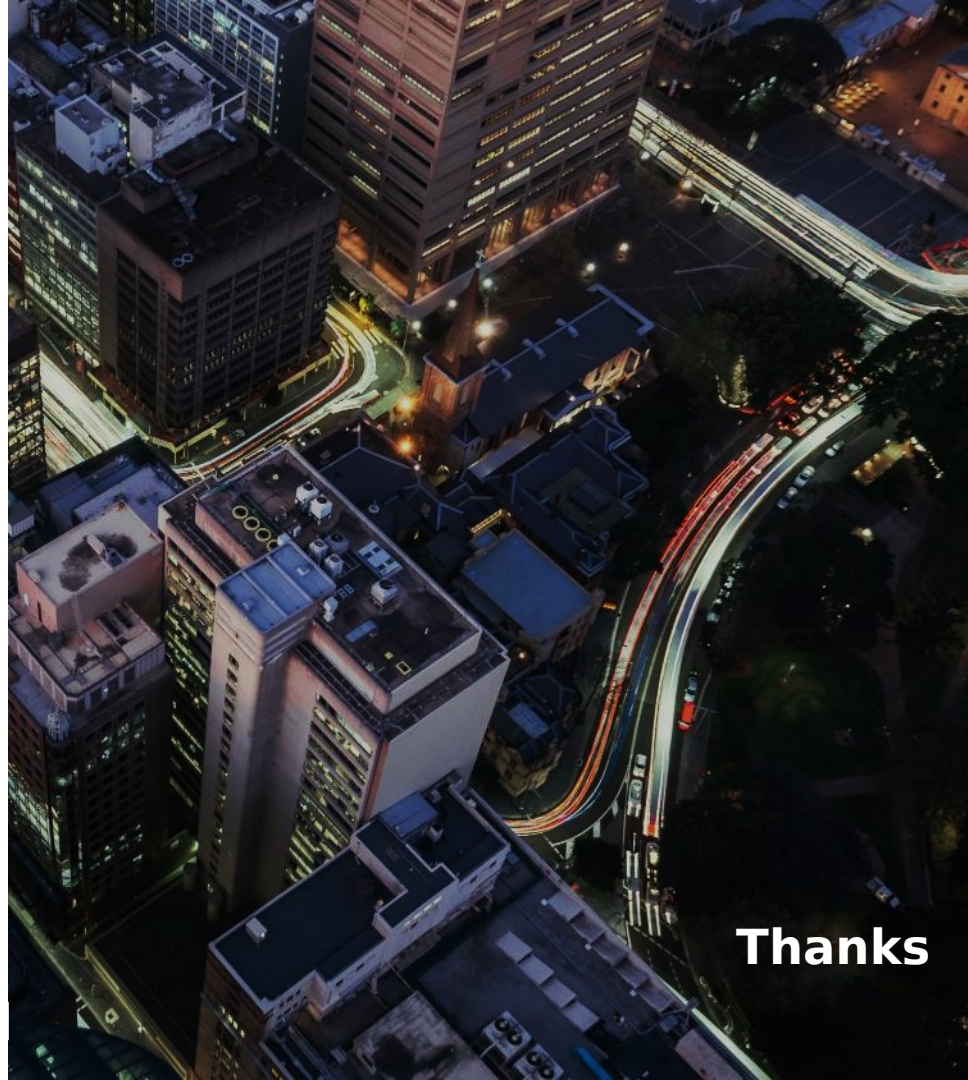
Stream run by:
Maintenance Team

Collaborations and Research

- Funded Project CEF upon “Air Quality and Mobility” , HPC analytics for air quality on transportation vehicles, in collaboration with Irisa, Rennes Metropole, GENCI/Idris, etc
 - Ryax responsible for orchestration and resource management from IoT sensors up to HPC clusters
- Research collaboration between Ryax Technologies and Inria Rhone-Alpes/LIG Datamove team upon simulation of Big Data workloads on hybrid edge/cloud infrastructures (started March 2018).
- Two other H2020 projects submitted and awaiting evaluation for funding:
 - **ICT-16** Software Technologies on Integrated programming models & techniques for exploiting the potential of virtualised and software defined infrastructures with Atos, ICCS, TUK, nviso, etc
 - **ICT-11a** Large-scale HPC-enabled industrial pilot test-beds supporting big data applications with Ubitech, Cineca, ICCS, IBM, Bull/Atos, etc



ryax
technologies



Thanks