

# Data Platforms, Clouds and Spaces: Trends in Contemporary Data Processing

**Prof. Vassil Vassilev & Team**  
*Cyber Security Research Centre*

# Content

## 1 Contemporary Data Processing

## 2 Data Platform of Cyber Security Research Centre

## 3 Alternatives and Choices

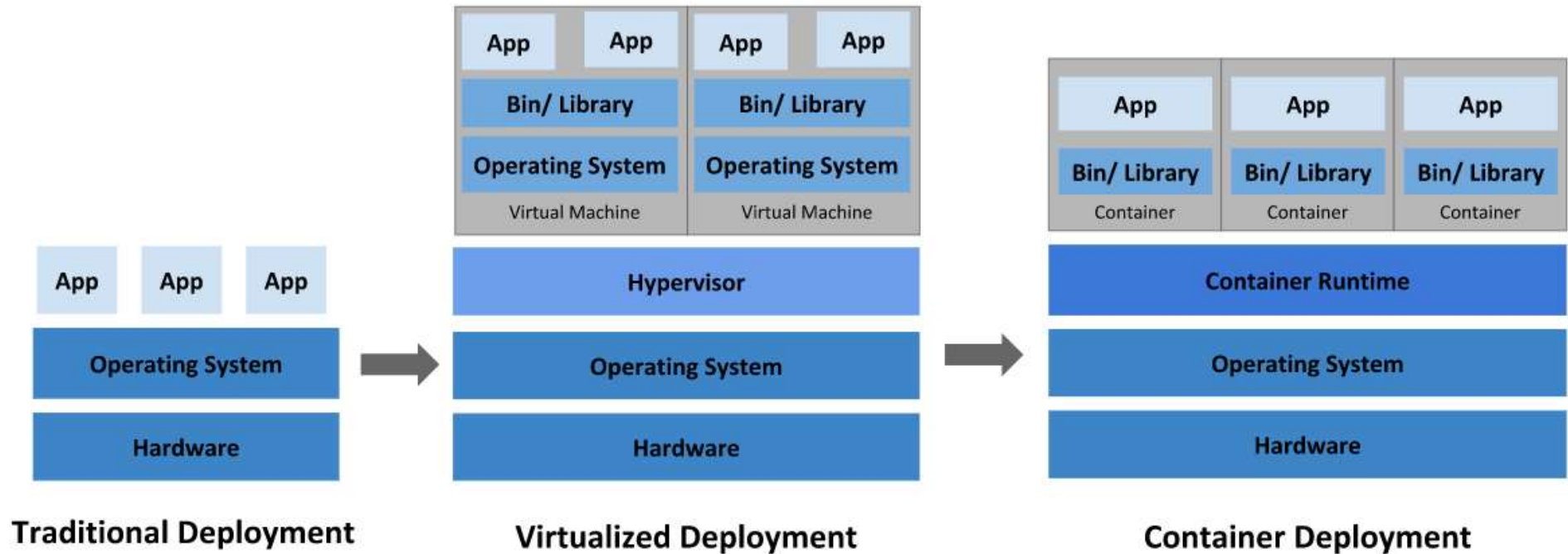
- ❖ **The Data**
- ❖ **The Metadata**
- ❖ **The Technologies**
- ❖ **The Tools**

## 4 The Lifecycle

- **Design**
- **Development**
- **Deployment**
- **Orchestration**
- **Monitoring**
- **Auditing**

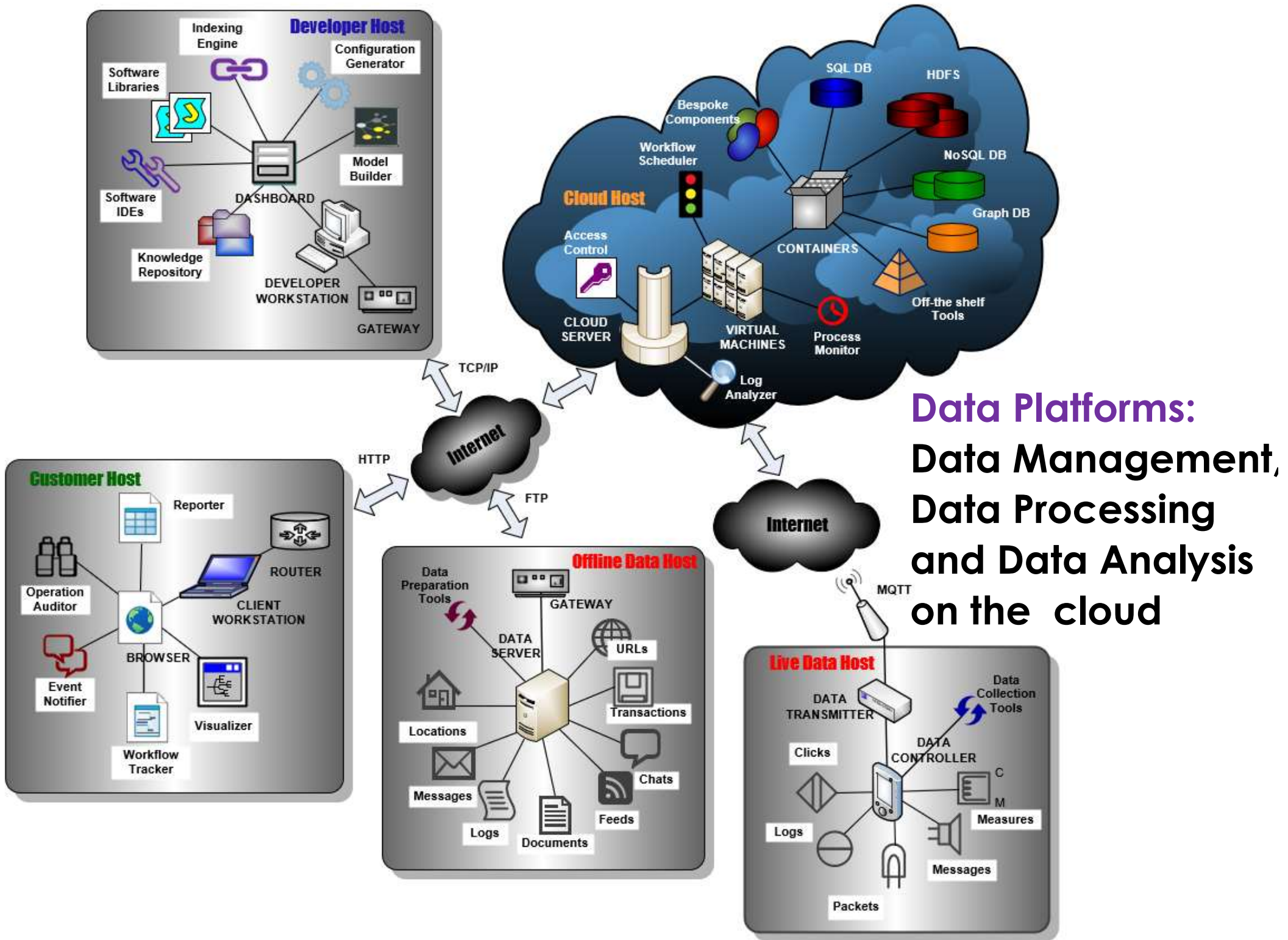
## 5 Where to go from here?

# 1 Contemporary Data Processing: From the Desktop to the Cloud



Source: <https://kubernetes.io/docs/concepts/overview/what-is-kubernetes/>

- ✓ Hardware infrastructure
- ✓ Operating system
- ✓ Middleware software
- ✓ Execution environment



**Data Platforms:**  
 Data Management,  
 Data Processing  
 and Data Analysis  
 on the cloud

# Data Spaces: Decentralized Supply and Consumption of Data Services

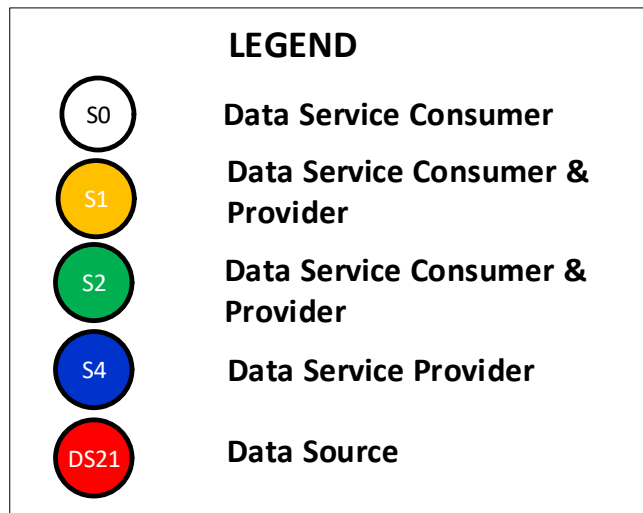
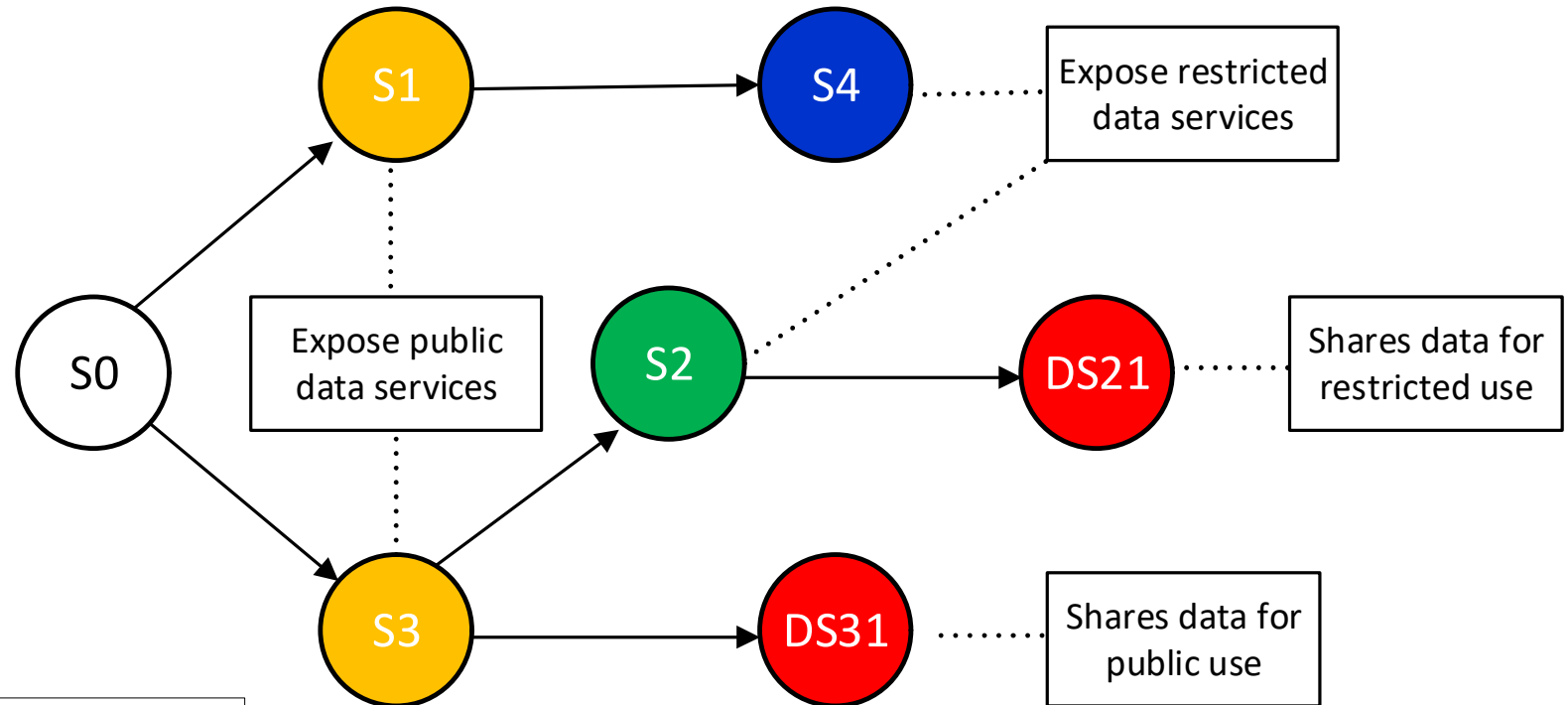
- Initiative of Fraunhofer Institute in Germany to serve the needs of Germany's industry, heavily based on B2B
  - Addresses the need for information from external sources in digital format over the Internet (**accessibility**)
  - Does not require the data to be available at the place of its use (**distribution**)
  - Retain the ownership while sharing the data and services (**control**)
  - Requires only two participants: provider (**sharing** the data it owns and/or **exposing** some data services) and consumer (**consuming** the shared data and/or **utilizing** the exposed data service for its own use); they can be combined in the same service-oriented architecture (client/server+server/server+peer-to-peer)

***Adopted across EU (International Data Space Association, IDSA)***

## Example: Urban Mobility Data Space

	City Mobility Centre	Environment Control Agency	Urban Planning Department
<b>Data Shared</b>	routes, places, vehicles, times	pollutions, standards, polluters	
<b>Data Access Permissions</b>	public (routes, places, times), restricted (vehicles, locations)	public (pollutions, standards), restricted (polluters)	
<b>Operations Supported</b>	locating, placing, timing	pollution, polluter determination	
<b>Operations Rights</b>	public (placing, timing), restricted (locating)	public (pollution), restricted (polluter)	
<b>Data Consumed</b>		routes, places, vehicles	vehicles, places, routes, pollutions, polluters
<b>Operations Executed</b>		placing, locating	place pollutions, vehicle pollutions, route pollutions

# Data Space as Service-oriented Architecture



# Data Platform support needed for Data Spaces

## ✓ Data service consumption

- » Identification of the services (addressing)
- » Identification of the consumers (authenticating)
- » Requesting the services (consuming)
- » Requesting the consumption (reporting)

## ✓ Data service provision

- « Registration of providers and consumers (identity management)
- « Assignment of responsibilities for data sharing and service provisioning to providers (provider profiling)
- « Assignment of rights for data access and privileges for operation execution to consumers (consumer profiling)
- « Identification of the communications (session tracking)
- « Logging of the operations (event logging)
- « Estimation of the consumptions (service reporting)
- « Calculation of the cost of consumption (billing)

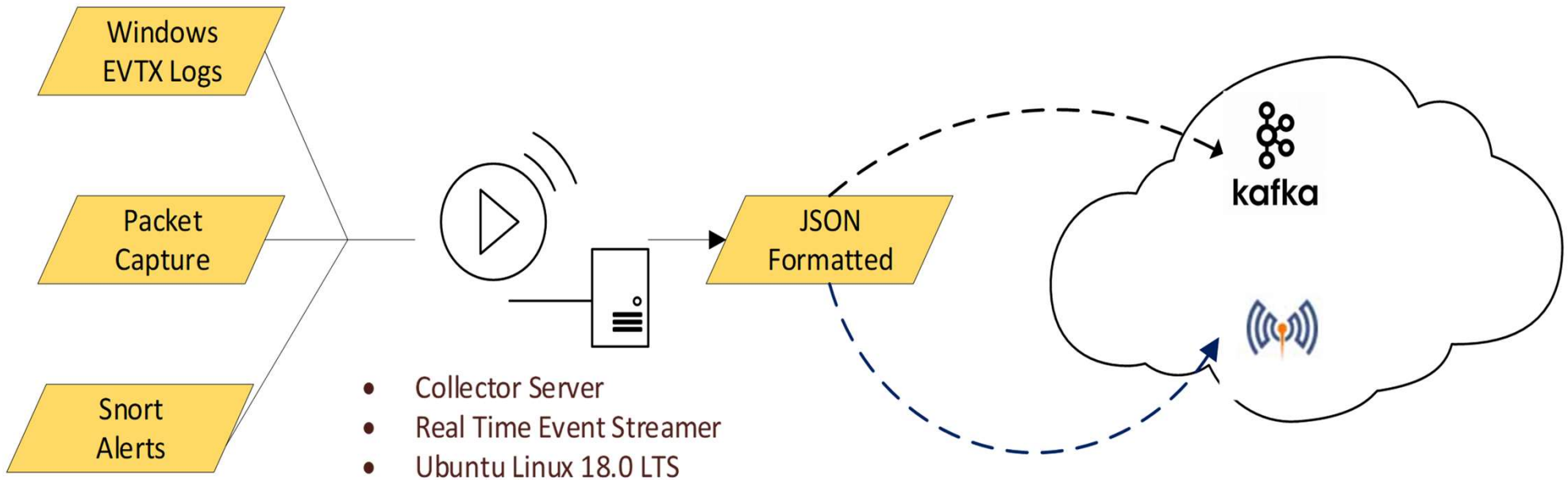
*Data platforms certification under GAIA-X initiative of IDSA*



## 2 Data Platform of Cyber Security Research Centre

- ❑ Piloted at the **Cyber Security Research Centre** of London Metropolitan University with funding from DCMS and HEIF
  - Private **cloud**, based on commodity server architecture
  - Public domain software for **virtualization, containerization** and **orchestration** of the server-side applications
  - Communal editions of enterprise software products for **data management** on the server
  - Free software for **application development**
  - Web-based interfaces for **service deployment** and **use**
- ❑ Replicated by the partner organization, **GATE Institute** of Sofia University
- ❑ Tested in three different scenarios for data processing in real-time: security analytics, outdoor air pollution in London and indoor air pollution factor analysis

# Project 1: Real-time Security Data Analytics



# Classification of Network Packets using Regression, NN and SVM

Model	Predicted regular packets:	Regular packets in test set:	Predicted ACK packets:	ACK packets in test set:	Predicted SYN packets:	SYN packets in test set:	Accuracy:
Neural Network	129	303	2023	1851	8	6	79%
Support Vector Machine	107	276	2050	1877	3	7	90%
Logistic Regression	417	258	1743	1893	0	9	71%
Linear Regression	791	255	1369	1897	0	8	60%

## Legend

- ACK** - acknowledgement flag confirming normal exchange of packets between two sites
- SYN** - synchronization flag signaling initiation of normal communication between two sides

# Detection of Potential Unauthorized Intrusions using CNN

## Actual flags within the dataset

## Predicted Suspicious Flags

ACK packets in test set: 39913

Predicted ACK packets: 43647

SYN packets in test set: 13

Predicted SYN packets: 12

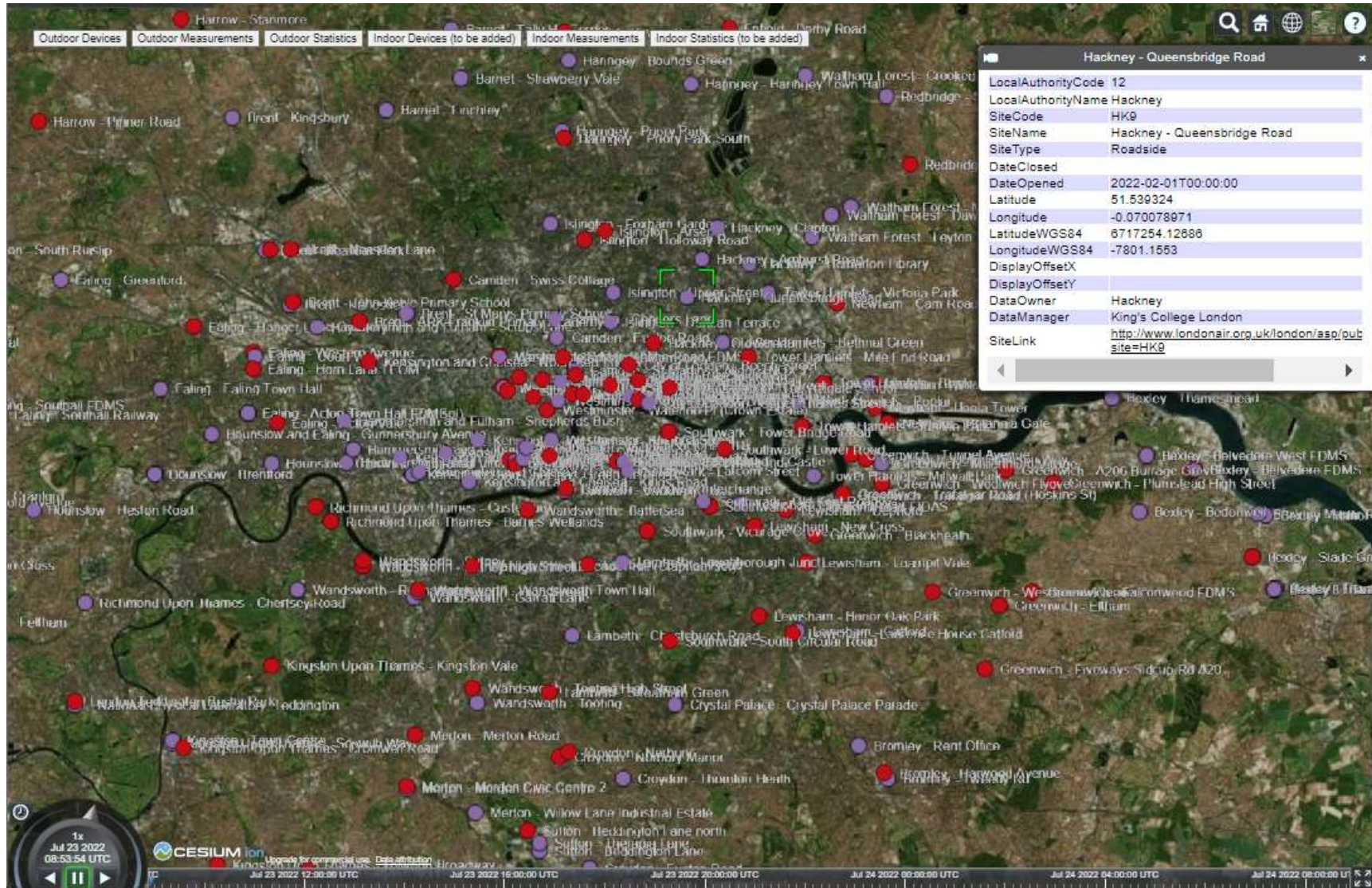
RST packets in test set: 165

Predicted RST packets: 164

## Legend

- ACK** - acknowledgement flag confirming normal exchange of packets between two sites
- SYN** - synchronization flag signaling initiation of normal communication between two sides
- RST** - warning flag sent after anomaly has been detected in the previous communication

# Project 2: Outdoor Air Pollution in London



# Logical Analysis of Air Pollution by combining Ontological and Sensor Data

[http://www.semanticweb.org/oem/ontologies/2021/10/untitled-ontology-2#Railway\\_Station\\_Sofia'>>](http://www.semanticweb.org/oem/ontologies/2021/10/untitled-ontology-2#Railway_Station_Sofia'>>)

The location longitude is 23.31639720.  
The location latitude is 42.72653500.

The nearest air quality station is AT22532689.  
The distance to air quality station is 1.8 kilometers.

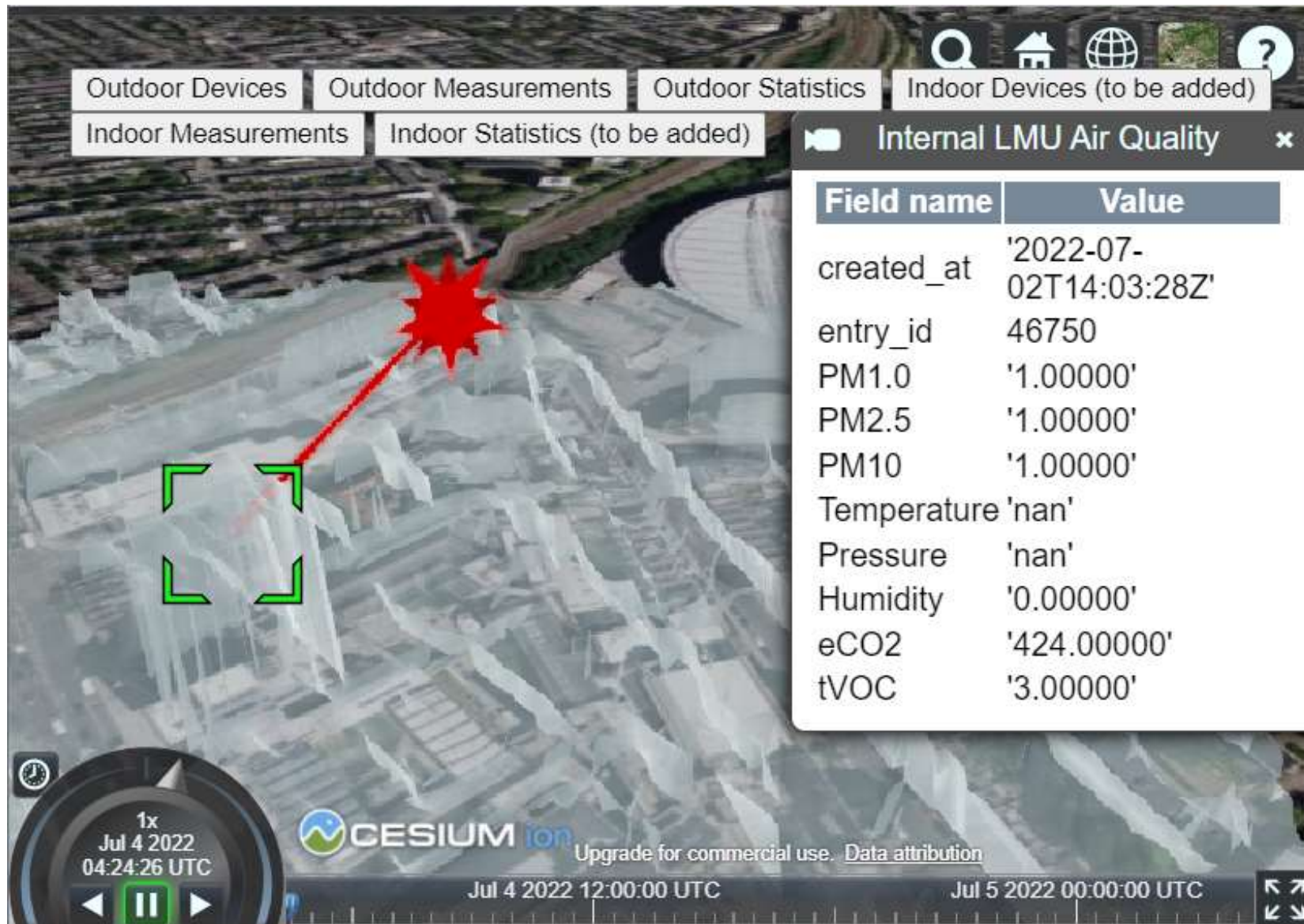
The CO level is 1.0268903624234662.  
The NO2 level is 27.976839433671802.  
The O3 level is 26.340781837908203.  
The PM10 level is 4.171716448194848.  
The PM2 level is 2.282572653193867.  
The SO2 level is 6.925875484004399.

<Record n=<Node id=14 labels=frozenset({'Resource', 'owl:NamedIndividual', 'ns0:Railstation'}) properties={'ns0:Has-coordinates': '23.31639720,42.72653500', 'uri': 'http://www.semanticweb.org/oem/ontologies/2021/10/untitled-ontology-2#Railway\_Station\_Sofia'>>

submit  
Go back

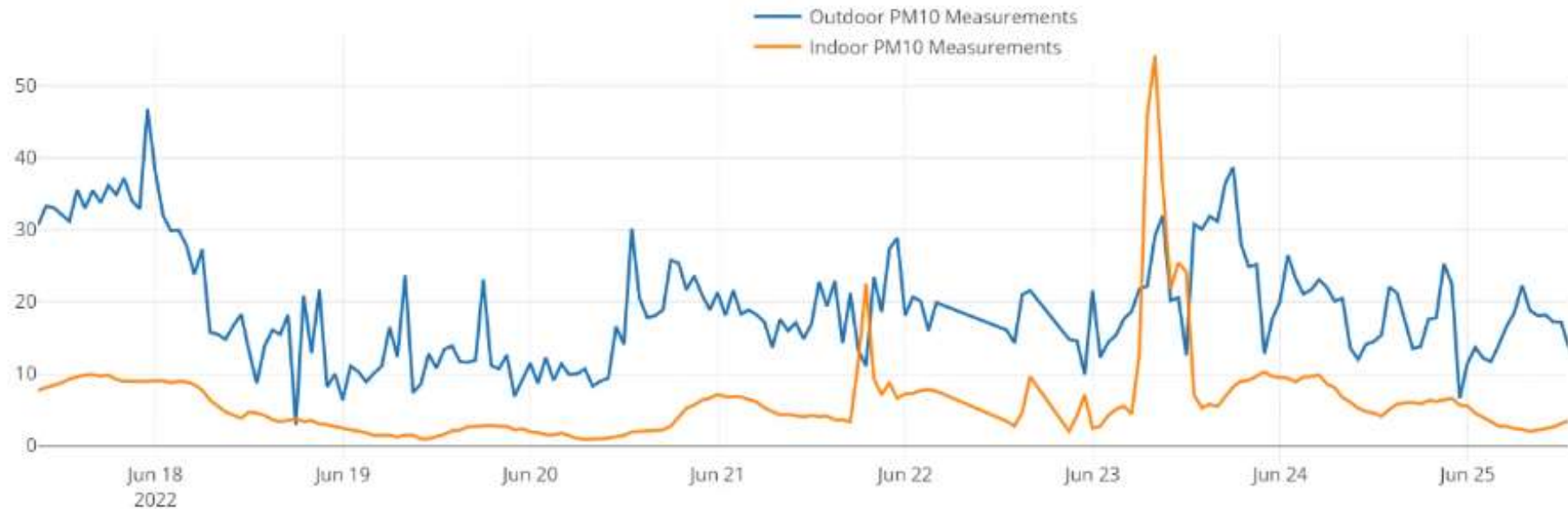
- Identifying location of the sensor station meta-data stored in **MongoDB** Database
- Loading individual descriptions of the objects from the city ontology stored in **Neo4J** Database
- Calculating the distance between air quality station and the objects in **Cypher** query language
- Analysing the air quality measurements at this location in **Python**

# Project 3: Indoor Air Pollution Analysis



# Correlation between Outdoor and Indoor levels of Particles in the Air

PM10 Measurements - Holloway Road



## Correlation Matrix

Indoor Value	Outdoor Value
1	0.3634504520543431
0.3634504520543431	1



### 3 Alternatives and Choices: Data, Metadata, Technologies and Tools

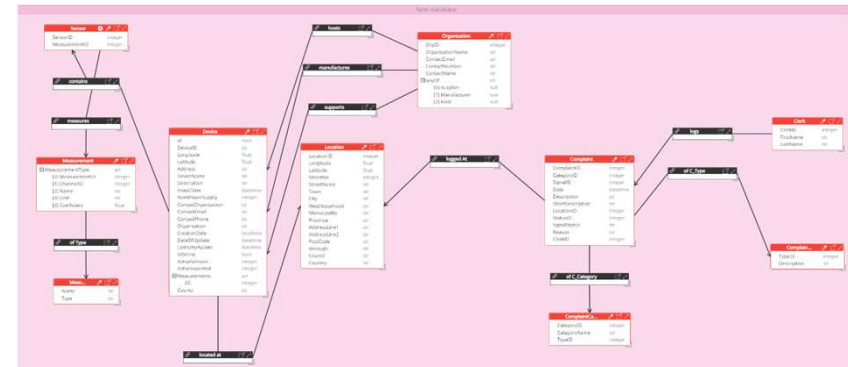
- The data is on **multiple scales** – different formats, granularity, volume, noise, location...
- The tasks for data analysis have **large diversity** – detection, recognition, classification, correlation, factorisation, prediction...
- There is a **variety of methods** with different applicability – temporal, structural, logical, model-driven, behavioural, hybrid
- Data analysis is performed as part of **complex workflows** – sampling, aggregation, buffering, feature selection, training, validation, analysis, merging, interpretation, explanation...
- The applications may require **significant resources** (both in terms of memory and computing power).
- AI technologies for data processing need to be comprehensive to reach **wide community** of users.

# All about Data

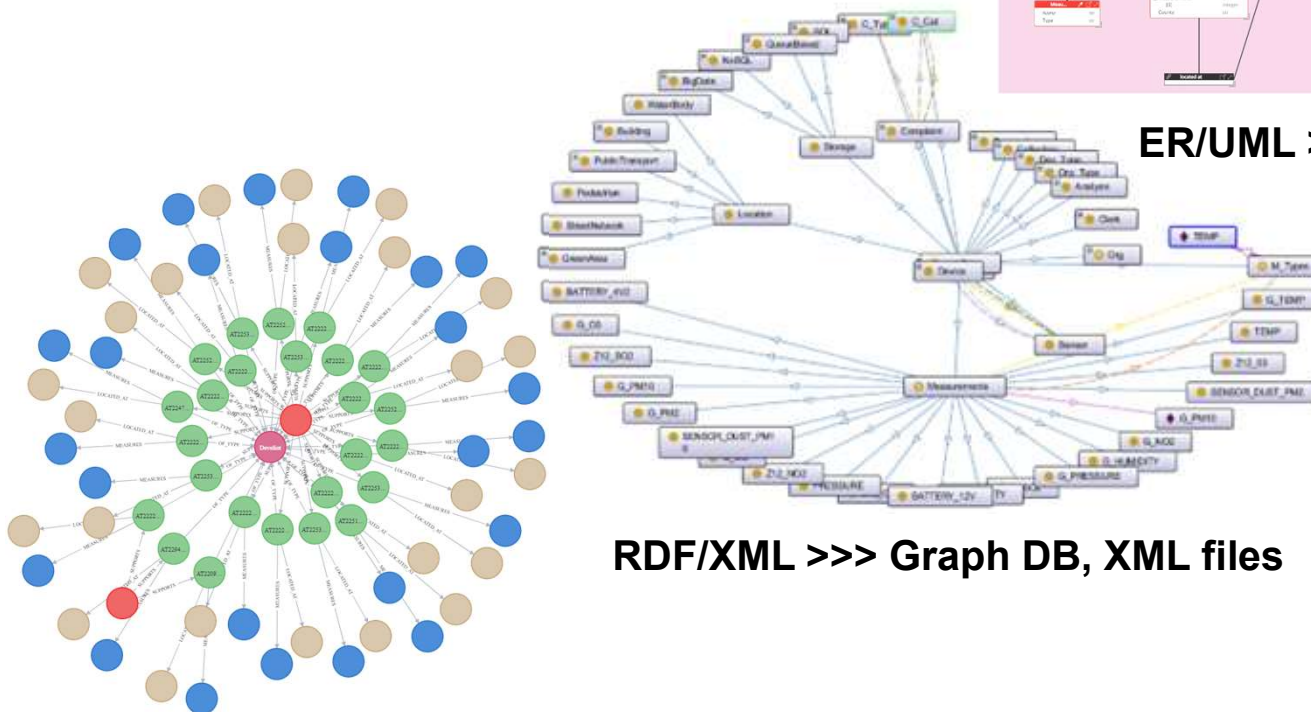
Data Types	Data Sources	Ingestion Methods	Transport Protocols
<b>Samples</b>	<b>Hardware</b> (external devices, infrastructure)	One-off	memory sharing, parameter passing
<b>Files</b>	<b>OS</b> (clients, suppliers)	One-off, Batch	FTP, HTTP, SCP, WebDAV, etc.
<b>Messages</b>	<b>Events</b> (Messengers, Listeners, Loggers)	One-off, Batch, Continuous	MQTT, AMQP, SMS/MMS, RCS, SOAP, etc.
<b>Repository Collections</b>	<b>Drivers</b> (databases, data warehouses, data lakes)	One-off, Batch	native to the repository
<b>Streams</b>	<b>APIs</b> (sensors, service providers)	Continuous	native to the streaming

# The Metadata and its Utilization

- ✓ Understand the data for **better design** of applications
- ✓ Enrich semantically the data for **more informative** and **more convenient handling**
- ✓ Prepare the data for **further** storage and processing



ER/UML >>> SQL DB, ER files



RDF/XML >>> Graph DB, XML files

JSON >>> NoSQL DB, JSON files

# Technologies for Data Processing

**Different Stages along the Data Processing Pipeline:** At the source, Before transmission, During transmission, On arrival, Within repository, After retrieval, etc.

**Different Structure and Formats of the Data:** structured (CSV, SQL), semi-structured (JSON, XML, RDF, SVG, etc.) and unstructured (binary, text, graphics, video)

**Different Preparation of the Rough Data:** Filtering, Formatting, Anonymisation, Normalization, Enrichment, Aggregation, Reconciliation, Buffering, Accumulation, etc.

**Different Methods for Data Analysis:** Statistics, Regression, Correlation, Clustering, Graph based, Rule based, Neural, Genetic, Swarm, Deep Learning, Reinforcement Learning, etc.

**Different Interpretations of the Results:** Simple reporting, Blackbox explanation, Whitebox explanation, Causal explanation, Impact factor analysis, etc.

# Software for Data Pipelines on the Cloud

Type	Software	Context
<b>Virtual Machine</b>	VMWare Workstation, Oracle VM, KVM Windows 10, MS Linux, etc.	OS or hypervisor
<b>Hypervisor</b>	VMWare vSphere, Oracle Virtual Box, MS Hyper-V, Linux KVM	OS
<b>Container</b>	Docker, LXC, Windows Containers, Portainer, Podman	OS, VM or container manager
<b>Container Manager</b>	Google Kubernetes, Apache Mesos, Docker Swarm, HashiCorp Nomad	OS
<b>Engine</b>	code interpreter (i.e., Python)	OS, VM or container
<b>Server</b>	off-the-shelf software (i.e., MongoDB)	OS, VM or container
<b>Application</b>	general server-side component (i.e., service registry)	Engine or server, deployed to OS
<b>Service</b>	domain-specific server-side component (i.e., sensor data filter)	Engine or server, deployed to VM
<b>Microservice</b>	application-specific server-side component (i.e., 2D city map)	Engine or server, deployed to container

# Advantages of Cloud-based Data Pipelines

## Containerization

- ✓ **Modularization** with no dependencies to set
- ✓ **Efficiency** in memory, CPU, and storage usage
- ✓ Application containers are **portable** across platforms without code changes
- ✓ Support for **configuration generation** through the use of parametrization and templates
- ✓ Full **traceability** of the operations for testing and debugging purposes

## Orchestration

- ✓ **Model-driven** application development
- ✓ Support for **reusability** of existing solutions in the form of process workflows
- ✓ Support for **auditing** of data processing pipelines for monitoring, analysing and billing purposes
- ✓ Support for **reproducibility** by preserving data dependencies
- ✓ Possibility for process **automation** based on workflow models and planning heuristics

# Take away: Always Informed Decisions with Optimal Choices

- ✓ **covering the full extent:** generation, pre-processing, transportation, post-processing, analysis, interpretation, reporting, etc.
- ✓ **constructing the richest models:** features, types, structures, constraints, annotations, indexes, etc.
- ✓ **selecting the most appropriate methods:** statistic, clustering, rule-based, graph-based, ML, RL, etc.
- ✓ **using the most convenient tools:** programming languages, engines, libraries, APIs, software tools, development methodologies
- ✓ **placing in the right context:** operating systems, virtual machines, containers, runtime engines
- ✓ **managing with an adequate organisation:** resources, policies, tasks, users

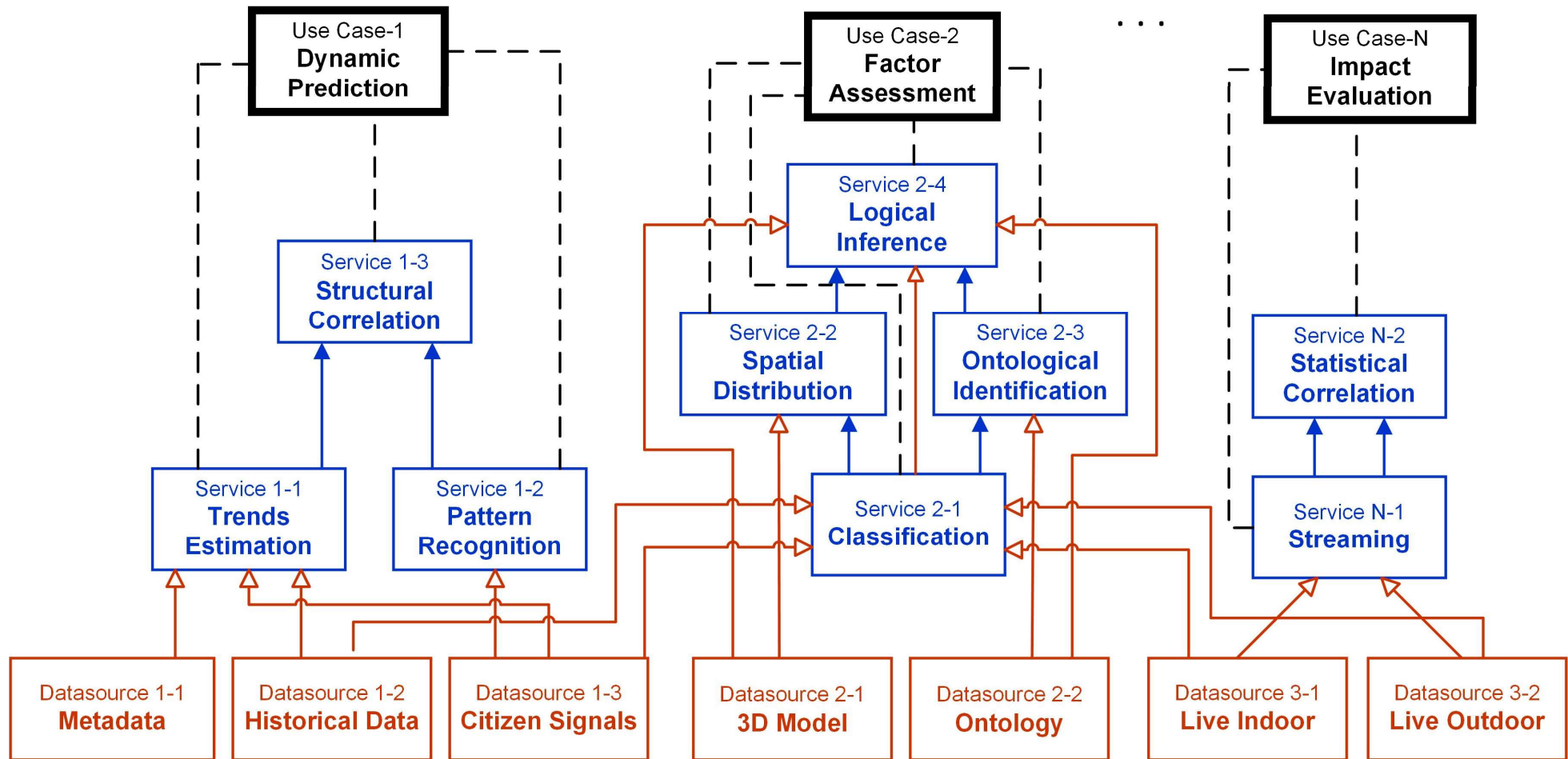
# 4 Data Services Lifecycle: DataOps in Action

**DataOps:** *“A set of practices, processes and technologies that combines an integrated and process-oriented perspective on data with automation and methods of agile software engineering to improve quality, speed, and collaboration and promote a culture of continuous improvement in the area of data analytics”*

- Logical **specification**
- Physical **design**
- Software **development**
- Application **deployment**
- Service **orchestration**
- Process **monitoring**
- Operations **auditing**

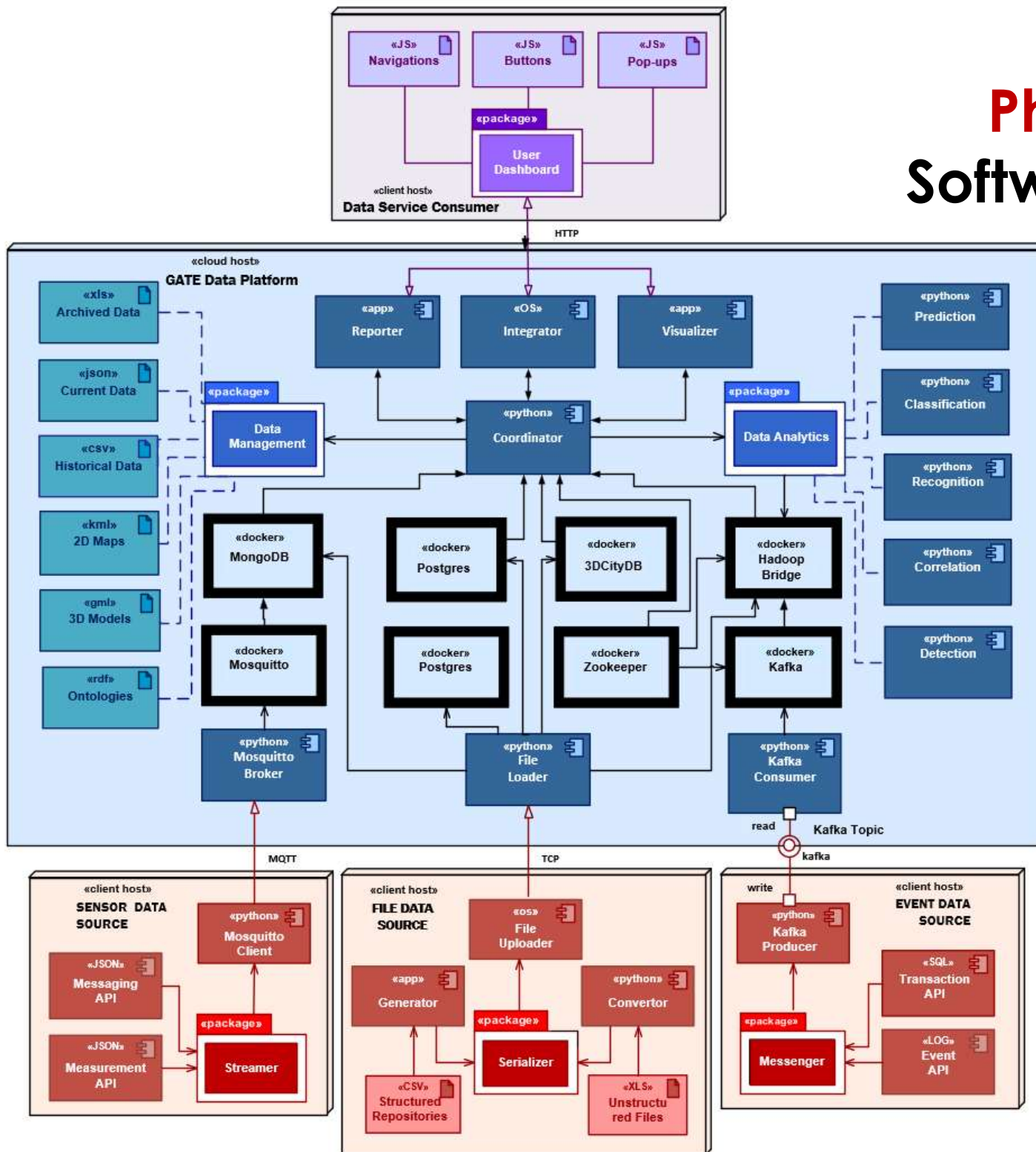


# Logical Specification: Data sources, Information Flows & Services



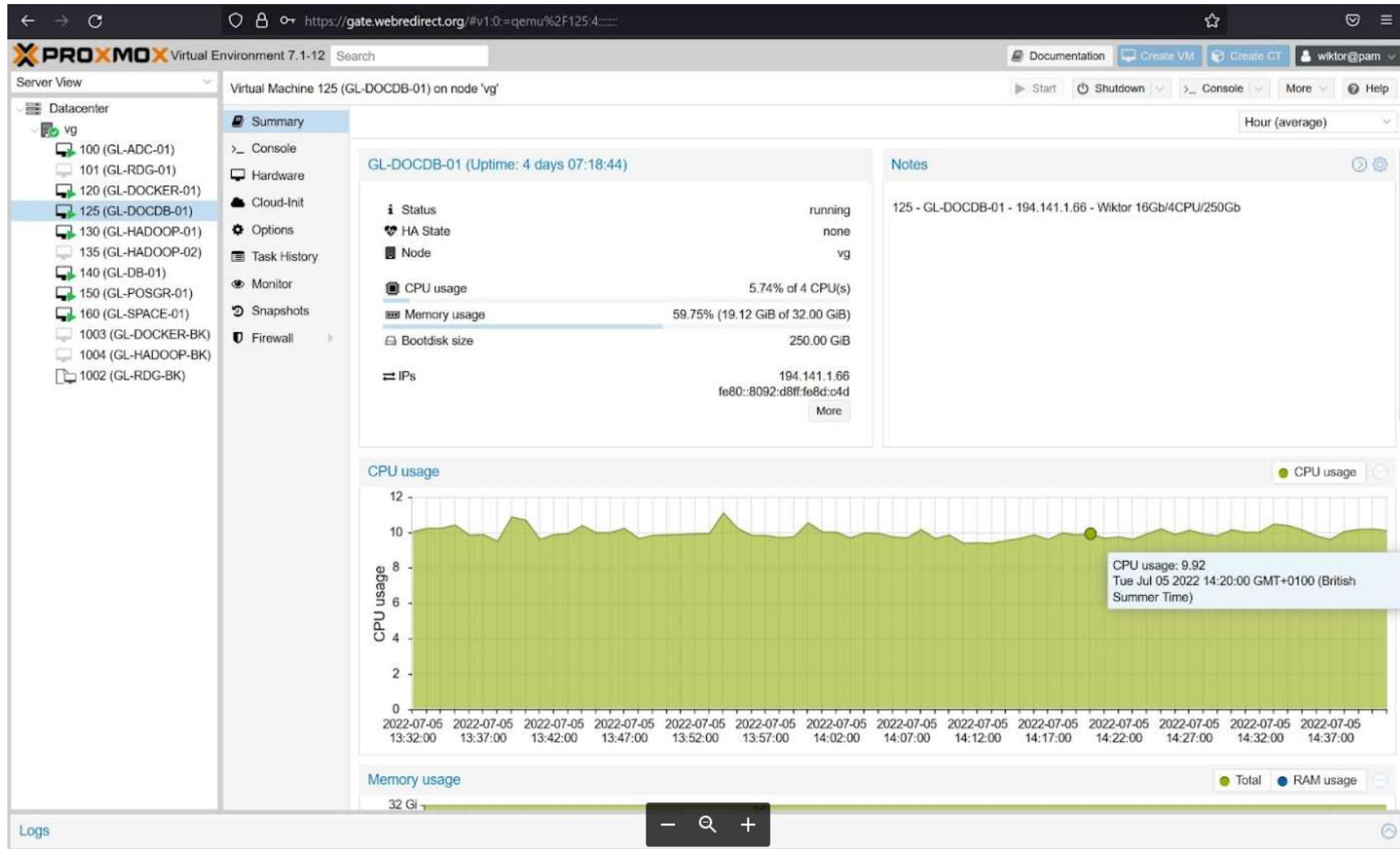
- ✓ Data sources
- ✓ Data processing services
- ✓ Parametric dependencies
- ✓ Information flows between services

# Physical Design: Software Components

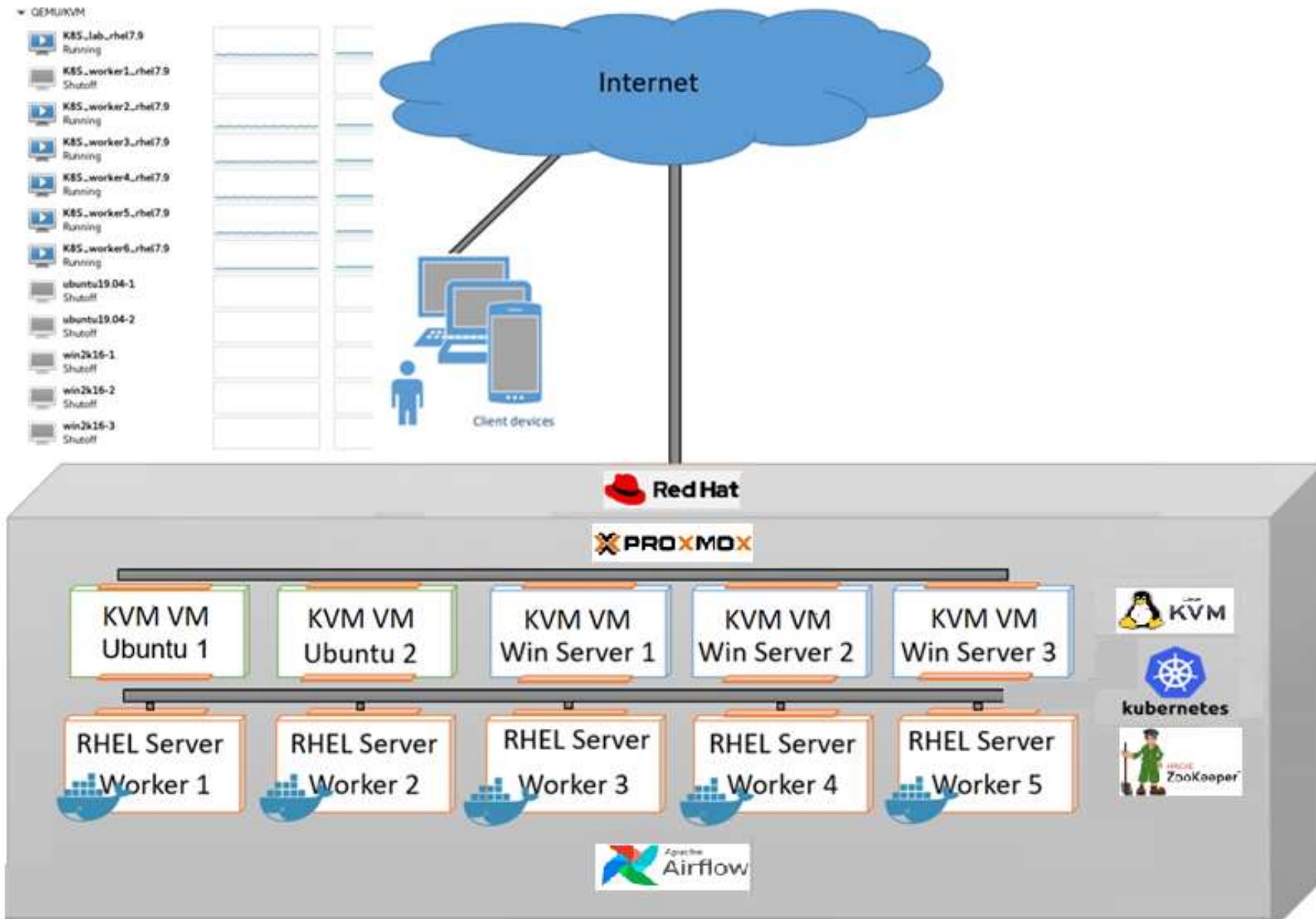


- Data formats, files and protocols
- Client-side and third-party provider APIs
- Server-side services and containerized microservices
- Language interpreters and runtime engines
- Configuration files, templates & scripts

# Software Development: Resource access using Proxmox VE



# Application Deployment: Component Management using OS, VM and Docker





# Process Monitoring: Execution Statistics gathering using MLFlow

← → ↻ ⚠ Not secure | 194.141.1.6:5000/#/experiments/0/runs/352ad13144e4476ab3c888ccf2

Number of new ACK packets streamed <a href="#">🔗</a>	974
Number of new RST packets streamed <a href="#">🔗</a>	2
Number of new SYN packets streamed <a href="#">🔗</a>	23
Number of new regular packets streamed <a href="#">🔗</a>	0
Number of predicted ACK <a href="#">🔗</a>	977
Number of predicted RST <a href="#">🔗</a>	2
Number of predicted SYN <a href="#">🔗</a>	20
Number of predicted regular packets <a href="#">🔗</a>	0
Number of ACK packets in initial test set <a href="#">🔗</a>	2008
Number of RST packets in initial test set <a href="#">🔗</a>	0
Number of SYN packets in initial test set <a href="#">🔗</a>	0
Number of new samples used for training <a href="#">🔗</a>	999
Number of regular packets in initial test set <a href="#">🔗</a>	4260
test accuracy - current model <a href="#">🔗</a>	0.983
test accuracy - updated model <a href="#">🔗</a>	0.32

## ▼ Tags

Name	Value	Actions
------	-------	---------

# Operations Auditing: Event Log Analysis and Reporting using Elasticsearch and Kibana



## 5 Where to go from here?

- **Cross-domain integration** of both data and analytics (i.e., environment and transport, environment and healthcare, healthcare and social services, environment and social services, etc.)
- Combining of real-time data with historical data for trends analysis and **investigation of process dynamics** (retrospective and predictive analytics)
- Localization of the data sources using rich data formats which allow **combining geolocation data with sensor data** (using richer representations such as semantic in GeoJSON and layering in KML)
- **Spatial navigation** to data sources by linking spatial ontologies and sensor data (using VR and games engines)
- Combining sensor data, symbolic meta-data and ontological information for **logical analysis** of the data beyond the pure data patterns (hybridization)



# Publications

- [1] **V. Vassilev, B. Virdee, K. Ouazzane, D. Maryanayagam, V. Sowinski-Mydlarz, et al.**, “Data Platform and Urban Data Services on Private Cloud”, in: *Proc. Int. Conf. Smart Trends in Computing and Communications* (SmartCom2023), 24-25 Jan 2023, Jaipur, India, IEEE, 2023 (in print).
- [2] **V. Sowinski-Mydlarz, V. Vassilev, K. Ouazzane, and A. Phipps**, “Security analytics framework validation based on threat intelligence”, in: *Proc. 9th Annual Conf. Computational Science and Computational Intelligence* (CSCI2022), Dec 14-16, 2022, Las Vegas, USA, IEEE, 2023 (in print).
- [3] **V. Vassilev, V. Sowinski-Mydlarz, D. Mariyanayagam, et al.**, “Towards first urban data space in Bulgaria”, in: *Proc. Int. Smart Cities Conference* (ISC2), 26-29 Sep 2022, Paphos, Cyprus, IEEE, 2022.
- [4] **V. Vassilev, S. Ilieva, D. Antonova, V. Sowinski-Mydlarz, et al.**, “AI-based Hybrid Data Platforms”, in: *Curry, E., Scerri, S. and Tuikka, T. (eds.), Data Spaces: Design, Deployments and Future Directions*, pp. 147-172, Springer, 2021.
- [5] **V. Vassilev, V. Sowinski-Mydlarz, P. Gasirowski, K. Ouazzane, et al.**, “Intelligence Graphs for Threat Intelligence and Security Policy Validation of Cyber Systems”, in: *P. Bansal et al. (eds.), Advances in Intelligent Systems and Computing*, Vol. 1164, pp. 125-140, Springer, 2020.

**Questions?**