# AI Inference Pipeline Composition with AI4Compose and OSCAR

Amanda Calatrava, Diego Aguirre, Vicente Rodríguez, Caterina Alarcón, Miguel Caballer and Germán Moltó

Universitat Politècnica de València (UPV)

Co-funded by
the European Union

# AI4EOSC Overview

Artificial Intelligence for the #EOSC

- Evolution of the DEEP Hybrid DataCloud platform.

- HORIZON-INFRA-2021-EOSC-01-04 call.

- Runs September 1st 2022 – August 2025 (36 months).

- 7 academic + 2 SME + 1 non-profit organization.

Advanced features for distributed, federated, composite learning, metadata provenance, MLOps, event-driven data processing, and provision of AI/ML/DL services.
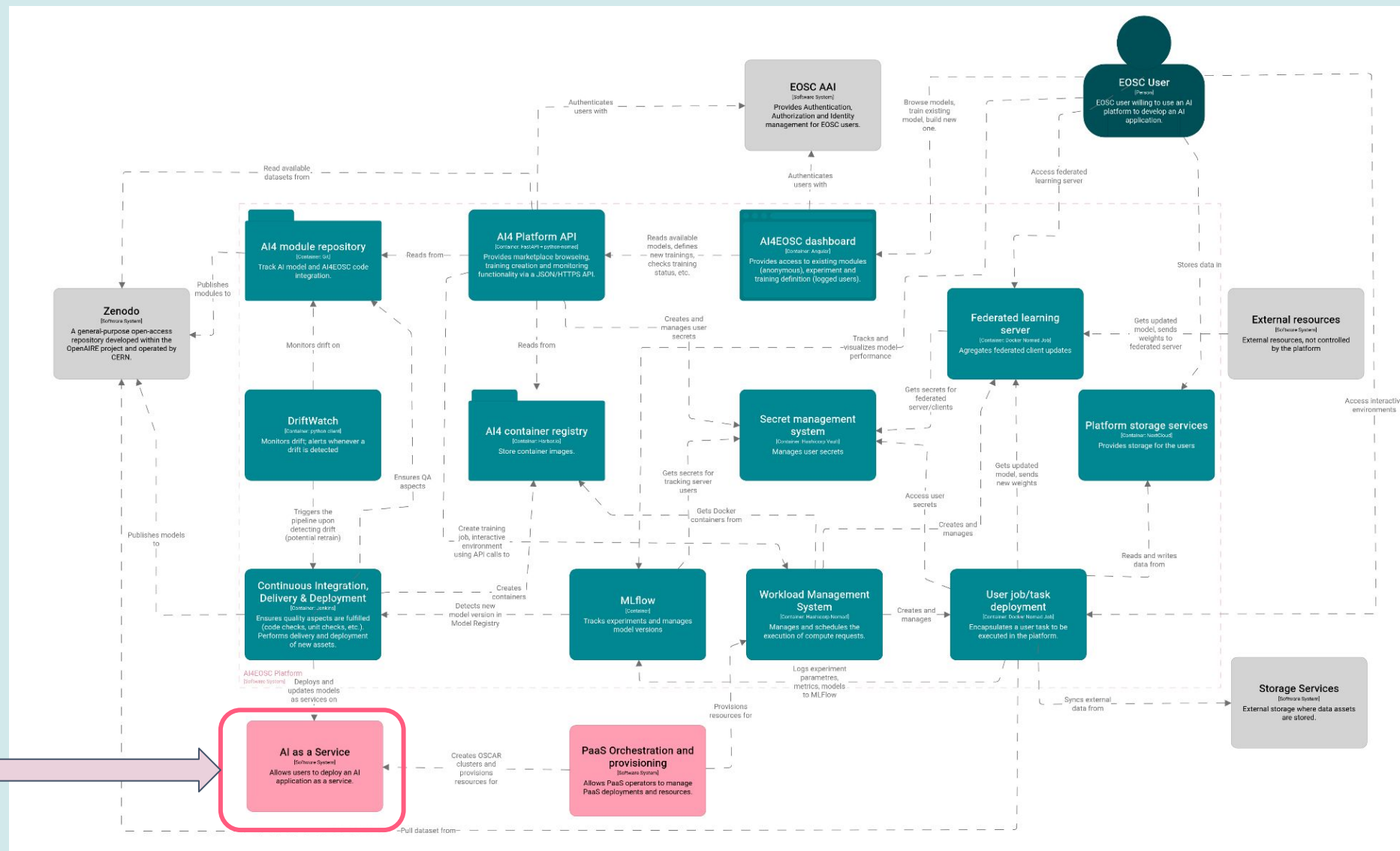
**Objective 3**
Services to compose AI tools, enabling the development of complex data-driven applications

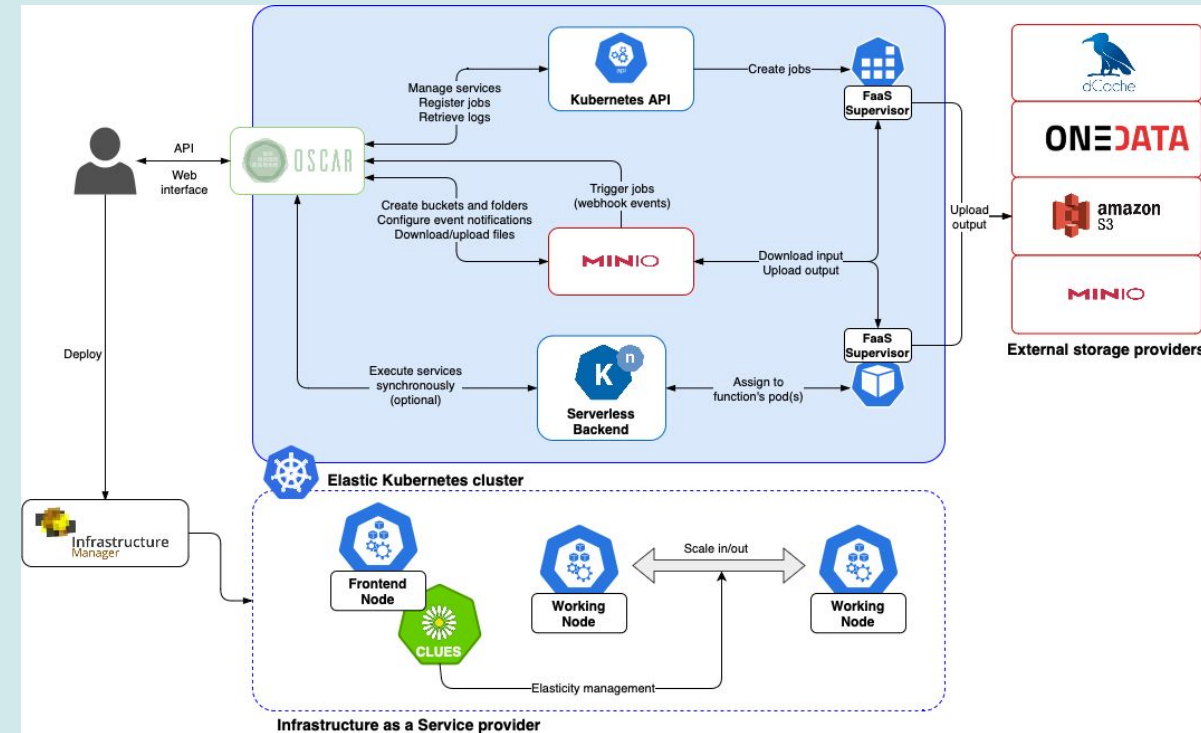→ Graphical Composition of **AI Inference Pipelines**

# AI4EOSC Platform architecture



Interactive C4 model: https://structurizr.com/share/73873/2f769b91-f208-41b0-b79f-5e196435bdb1/diagrams#ai4eosc_container_view

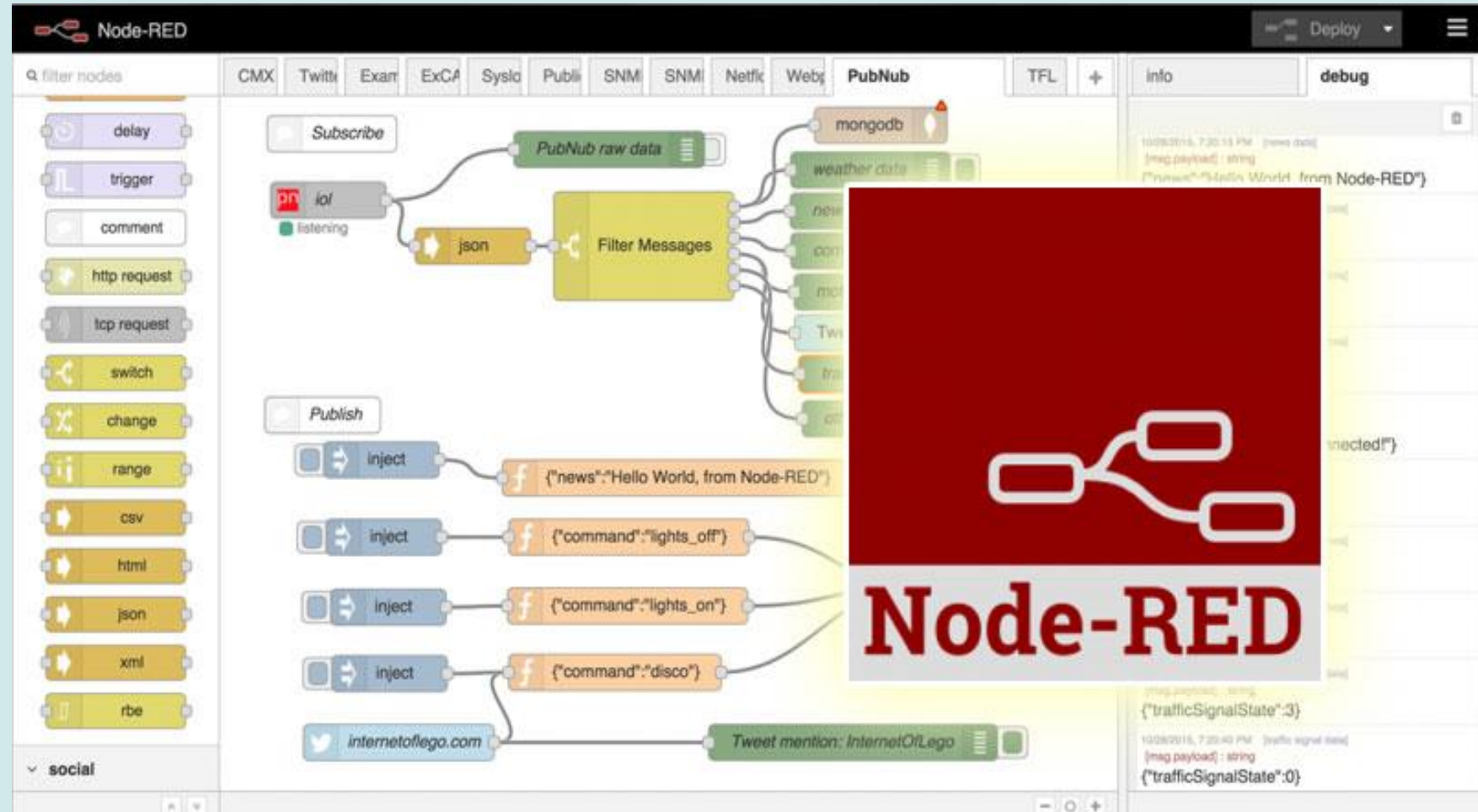# What is OSCAR?

OSCAR is…

- Open-source: https://github.com/grycap/oscar

- Serverless (event-driven + automated elasticity) computing model for data-processing Docker-based applications on Kubernetes.

  - Users do not need to manage the capacity provisioning of the cluster.

- Easily deployed on multi-Clouds via the Infrastructure Manager (IM).

- OSCAR allows to execute AI models packaged as Docker images.

- Users can upload files to an OSCAR cluster to automatically trigger the inference.

- OSCAR also allows programmatic access to trigger the inference of AI models.

- OSCAR can help you to expose your trained AI model to the community.
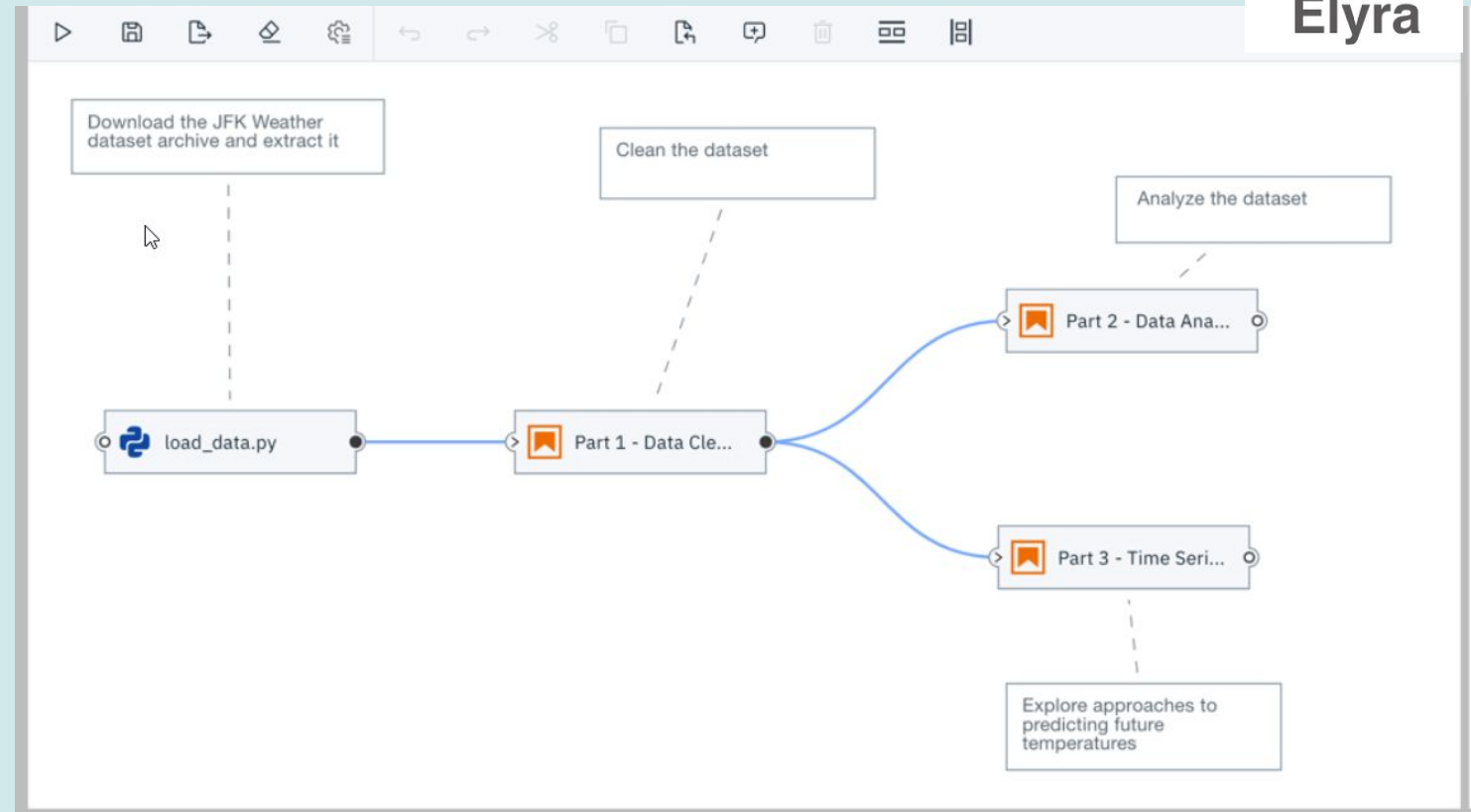


https://oscar.grycap.net

# What is Node-RED?

- Node-RED is a flow-based programming tool for connecting hardware devices, APIs, and services.

- Built on NodeJS and the D3.js library.

- The minimal structure are the nodes.

- Nodes are organized in flows that connect nodes.
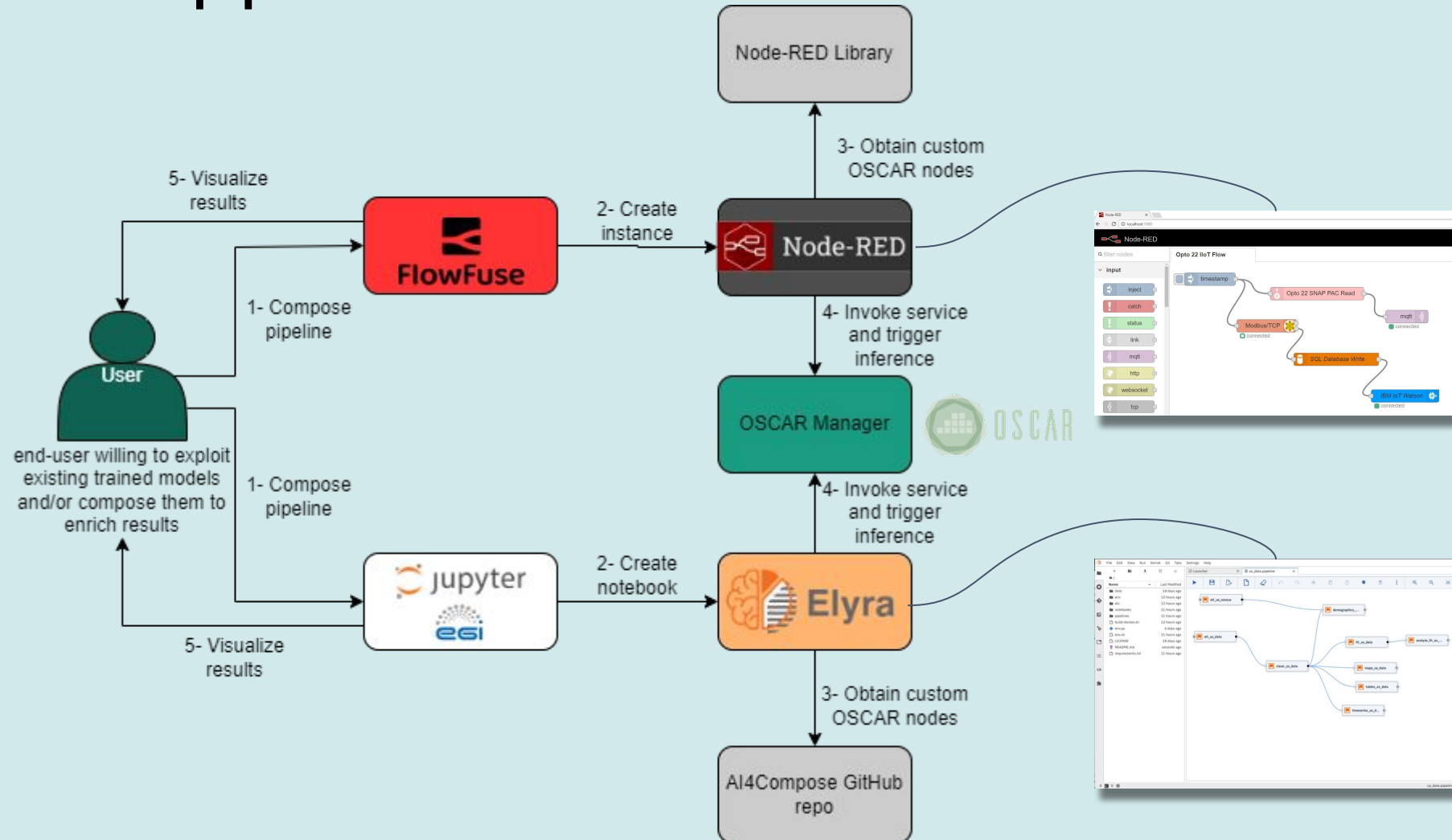


https://nodered.org

# What is Elyra?

- Elyra is an open-source project for developing and running machine learning workflows in JupyterLab.

- Provides tools for users to create visual pipelines for machine learning workflows.

- Elyra is programmed in Python and allows the use of many libraries focused on data analysis.

- The minimum structure for a workflow are the nodes (Python script, R script, and Jupyter notebook).



https://elyra.readthedocs.io/en/stable/

# AI4Compose: Low-code composition of AI inference pipelines.

# Composing AI pipelines: an example

# Composing AI pipelines: Node-RED demo

# Composing AI pipelines: Elyra demo

# **AI4Compose:** Innovation in AI4EOSC

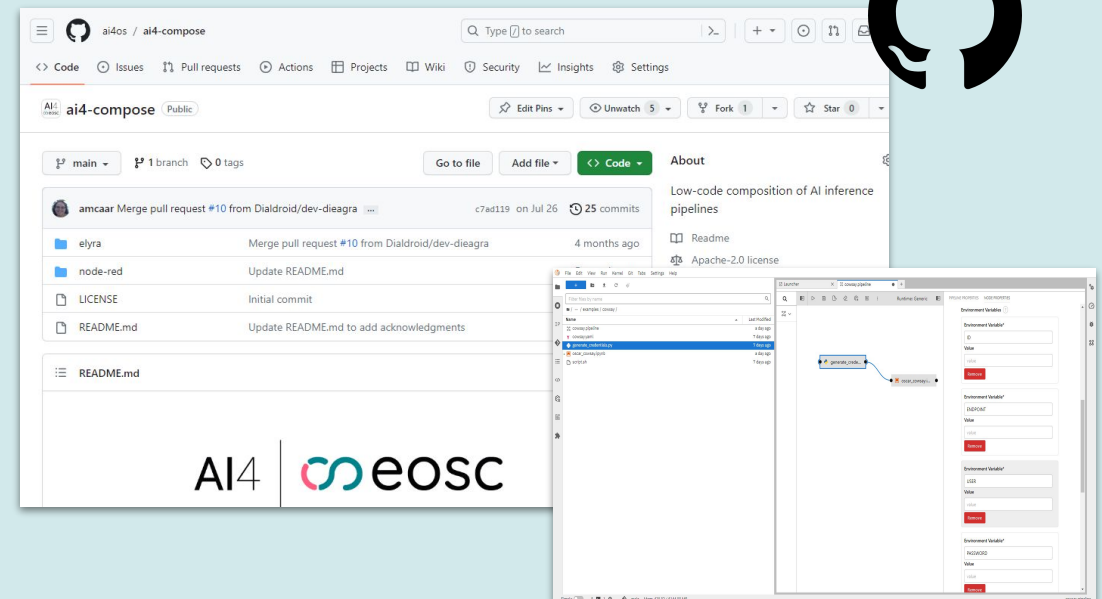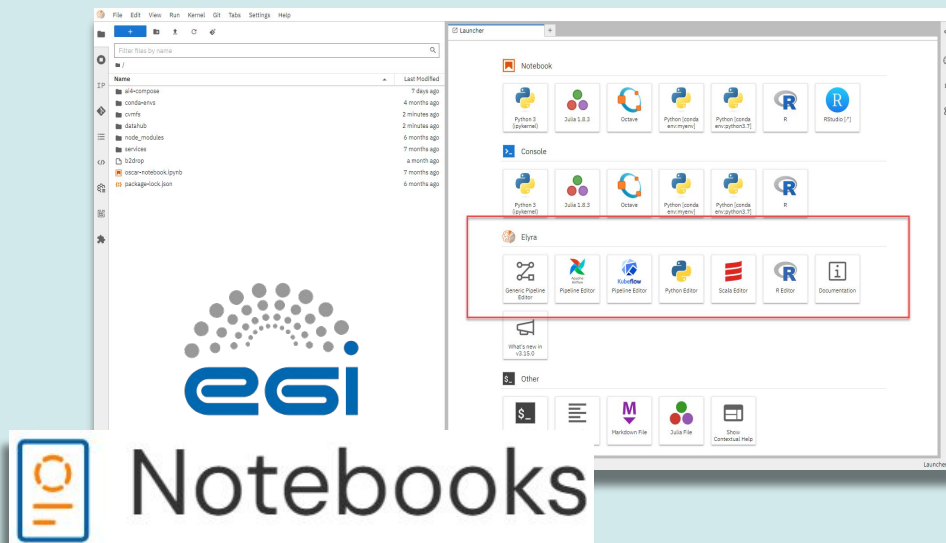- New modules to interact with OSCAR, published in the official Node-RED library (*more coming soon*):
  - https://flows.nodered.org/search?term=AI4EOSC

- Flowfuse: to manage users and multiple instances of Node-RED.
  - Dedicated instance for the project.
  - OSCAR templates ready to offer a pre-configured Node-RED instance with OSCAR support.
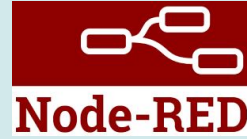
# AI4Compose: Innovation in AI4EOSC

- Elyra in [EGI Notebooks](#): we have contributed to have official support of Elyra inside the service, so users can easily access the tool.

- Notebooks to be used in Elyra for AI models of the DEEP Open Catalog:

  ○ https://github.com/ai4os/ai4-compose/tree/main/elyra/examples

- Documentation in the official AI4Docs repo: https://docs.ai4os.eu/en/latest/user/index.html#deploy-a-model-in-production

# Conclusions

- AI4Compose offers a visual support (drag & drop + customization) to compose AI inference pipelines.

- It minimizes the orchestration effort:
  - Multiple AI models can be triggered for inference and later aggregate the results for enhance accuracy.

- Reusable functions:
  - Pre-defined workflows can be created to facilitate interaction among the AI models in the DEEP Open Catalog.
  - Specific nodes can be created for the different AI models for a simpler definition of workflows.

- Workflows along the computing continuum can be supported (e.g. OSCAR clusters in disparate computing infrastructures).

- Each node can be configured to invoke an OSCAR service within a specific OSCAR clusters.

- Deployment of pre-defined Node-Red instances to facilitate AI composition of workflows (multi-tenant support thanks to FlowFuse).
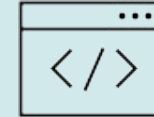
AI4 | eosc

Co-funded by
the European Union

AI4EOSC    amcaar@i3m.upv.es    ai4eosc.eu

# Thank you! Any questions?

Amanda Calatrava (on behalf of all the authors)