**Vicente Rodriguez Benitez**

vrodben1@i3m.upv.es

Director: Germán Moltó Martínez

Programa de Doctorado en Informática

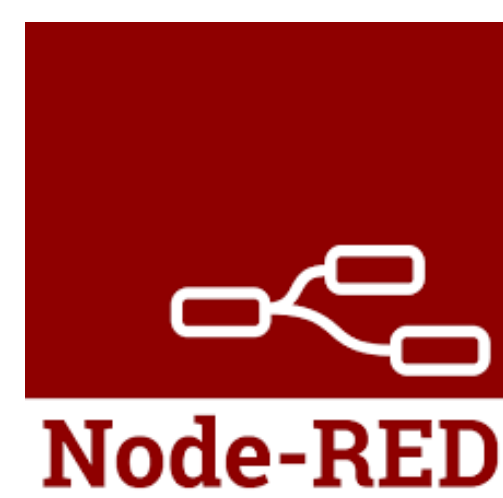Instituto de Instrumentación para Imagen Molecular (I3M)

# Composite AI through Serverless Orchestration

## Introduction

Recent years have witnessed how technologies for cloud services are advancing at a very fast pace. Serverless services allows you to create and run applications quickly and with a lower total cost of ownership, since it is not necessary to provision and manage infrastructure. For the creation of this type of services, two tools will be used: OSCAR developed at the UPV and Node-RED developed by IBM, both open source.

- **Node-RED** (www.nodered.org**)** is a flow-based programming tool, originally developed by IBM's Emerging Technology Services team and now a part of the OpenJS Foundation. It is a powerful tool that serves to communicate hardware and services in a fast and easy way.
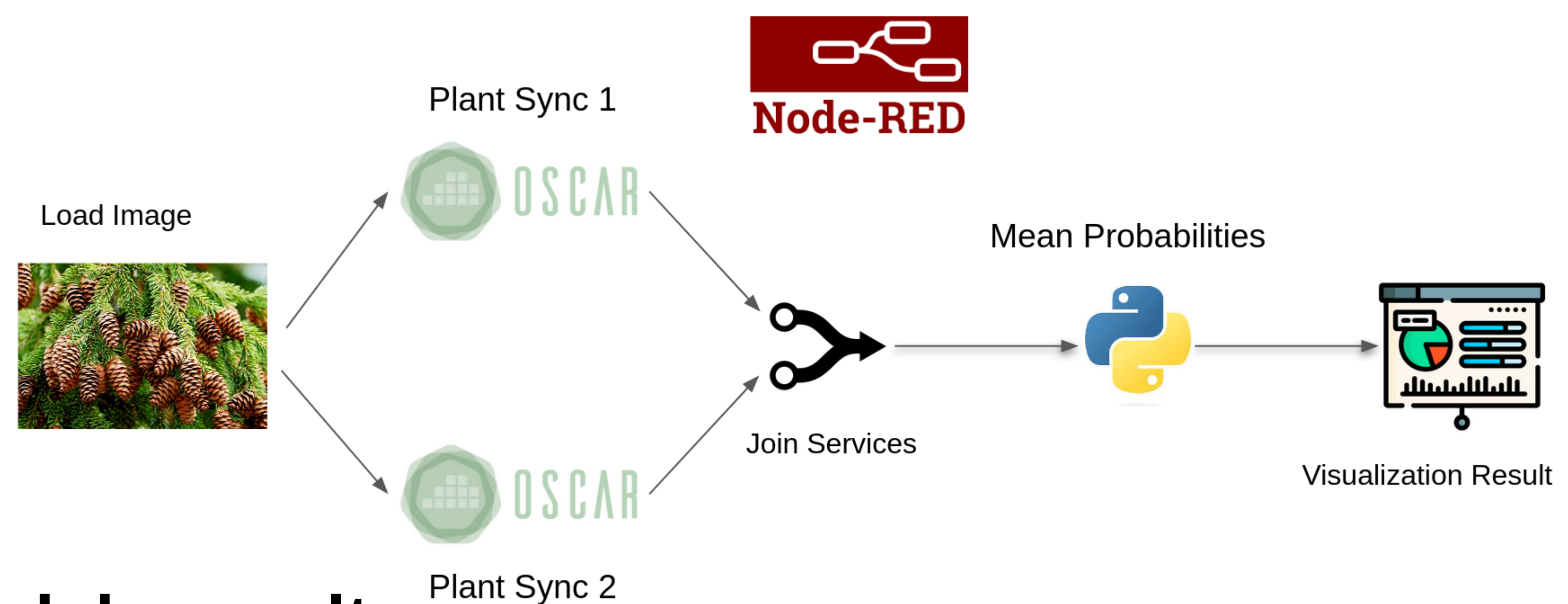
**OSCAR** (www.oscar.grycap.net) is a framework to efficiently support on-premises serverless applications for general-purpose data-processing computing applications. It supports a High Throughput Computing Programming Model to create highly-parallel event-driven file-processing serverless applications that execute on customized runtime environments provided by Docker containers run on AWS Lambda. [2]
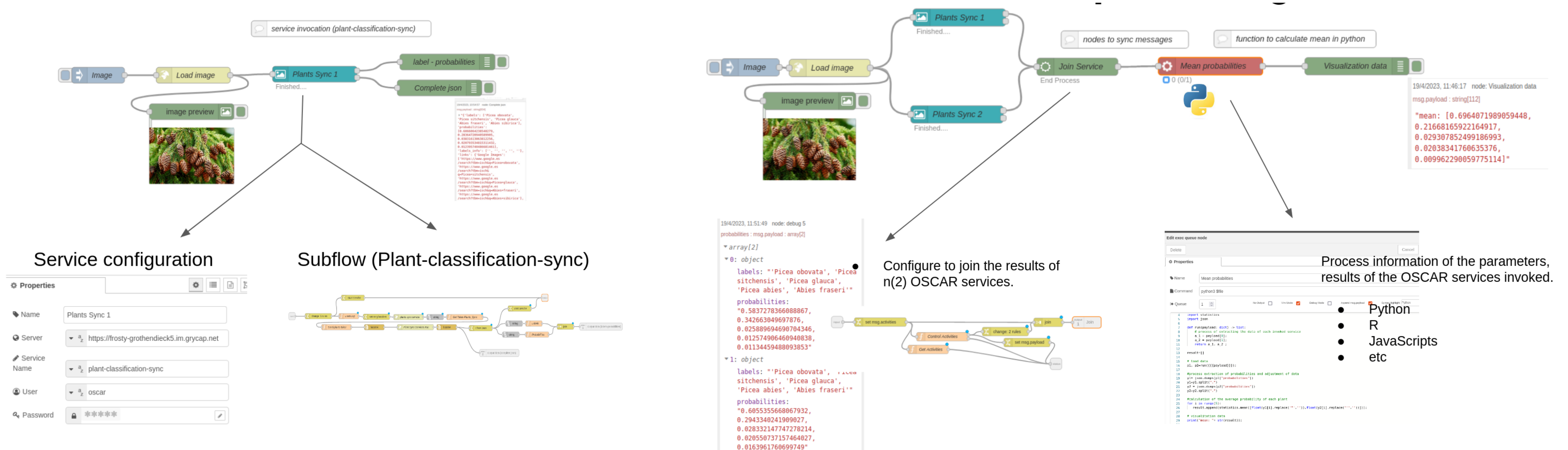
## General Objective

Develop workflows capable of orchestrating the distributed inference of AI models on OSCAR clusters with easy interaction by users through the usage of Node-RED.

## Stages of development

The work is based on the implementation of a workflow on Node-RED [1][2] in which two services are called on OSCAR in parallel. This service is Plant Classification with Lasagne/Theano [3]. Once the results are obtained, the results are aggregated for enhanced accuracy.



## Composite AI models result



Service configuration

Subflow (Plant-classification-sync)

Configure to join the results of n(2) OSCAR services.

Process information of the parameters, results of the OSCAR services invoked.

- Python
- R
- JavaScripts
- etc

## Expected Results

- Specific nodes (or subflows) in Node-Red can be created for the different AI Models for easier definition of the workflows.
- Each node can be configured to invoke an OSCAR service within specific OSCAR clusters.
- Pre-defined workflows can be created to facilitate interaction among the AI models in from the AI4EOSC project.
- Event-driven serverless workflows can be used to combine the outputs of different AI Models.
- Dashboards can be created to facilitate output data processing within the framework.

## References

[1] Kousiouris, G., Ambroziak, S., Costantino, D., Tsarsitalidis, S., Boutas, E., Mamelli, A., & Stamati, T. (2022). Combining node-red and openwhisk for pattern-based development and execution of complex faas workflows. *arXiv preprint arXiv:2202.09683*.

[2] Kousiouris, G., Ambroziak, S., Zarzycki, B., Costantino, D., Tsarsitalidis, S., Katevas, V., ... & Stamati, T. (2023, April). A Pattern-based Function and Workflow Visual Environment for FaaS Development across the Continuum. In *Companion of the 2023 ACM/SPEC International Conference on Performance Engineering*.

[3] Heredia, I. (2017, May). Large-scale plant classification with deep neural networks. In *Proceedings of the Computing Frontiers Conference* .