# Bridging Cloud and HPC for Scalable Event-driven Processing of AI Workloads

Estíbaliz Parcero, Sergio Langarita, Germán Moltó

Instituto de Instrumentación para Imagen Molecular (I3M), Universitat Politècnica de València, Valencia, Spain.
esparig@i3m.upv.es

# Bridging Cloud and HPC for Scalable Event-driven Processing of AI Workloads

- Bridging Cloud and HPC ⇦ **What** we want

- Scalable Event-driven Processing ⇦ **How** we run workloads

- AI Workloads ⇦ **Why** we need that bridge

# Index

OSCAR

# AI workloads

The number of cloud-native machine learning tools is rapidly increasing, offering more options and capabilities to AI engineers:

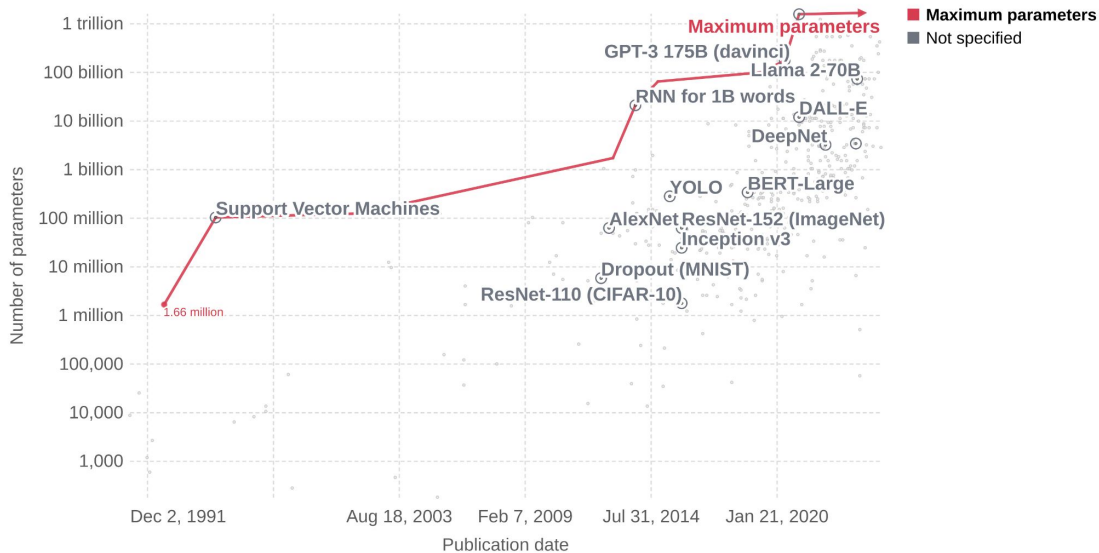# AI workloads

AI engineers use Cloud-native ML tools:

- **Accessibility and Cost-Effectiveness**: Provides necessary compute, storage, and services
- **Scalability**: Handles large datasets and complex models
- **Flexibility**: From programming languages, or pre-trained models to custom algorithms
- **Speed and Efficiency**: Speeds up building, training, and deploying models
- **Integration and Interoperability**: Integrates with other tools and services

# AI workloads

## Parameters in notable artificial intelligence systems

Parameters are variables in an AI system whose values are adjusted during training to establish how input data gets transformed into the desired output; for example, the connection weights in an artificial neural network.
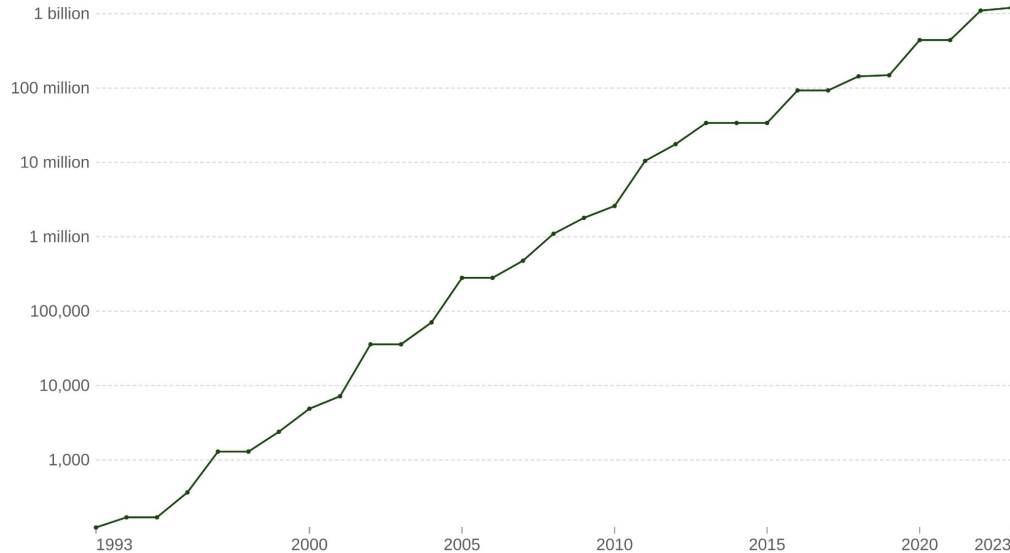


- Exponential growth in the number of parameters

- Training tasks take a lot of resources

- Also Inference, eg:
  - LLMs
  - Diffusion models
  - GANs

**Data source:** Epoch (2024)   OurWorldinData.org/artificial-intelligence | CC BY
**Note:** Parameters are estimated based on published results in the AI literature and come with some uncertainty. The authors expect the estimates to be correct within a factor of 10.

# AI workloads

## Computational capacity of the fastest supercomputers

The number of floating-point operations[1] carried out per second by the fastest supercomputer in any given year. This is expressed in gigaFLOPS, equivalent to $10^9$ floating-point operations per second.



**Data source:** Dongarra et al. (2023)

OurWorldinData.org/technological-change | CC BY

**1. Floating-point operation**: A floating-point operation (FLOP) is a type of computer operation. One FLOP represents a single arithmetic operation involving floating-point numbers, such as addition, subtraction, multiplication, or division.

- Exponential growth in computational capacity

- Available on HPC Clusters

# AI workloads

How to keep benefiting from Cloud ML tools and leverage the HPC resources?

Bridging Cloud and HPC seems a good idea

# Bridging Cloud and HPC



- Hoefler et al., 2024, highlight the importance of bridging Cloud and HPC for resource-intensive workloads like **climate simulations** or **machine learning** processing

- Some key ideas:
  - Leverage containers
  - Improve communication
  - Enable access to data (I/O)

# Scalable Event-driven Processing

- **OSCAR**, an open source platform for serverless event-driven computing

- **OSCAR** is a cloud-native tool, so all the previous benefits apply

https://oscar.grycap.net/

# Scalable Event-driven Processing

**Multi-Cloud Support:**
Provision OSCAR clusters on on-premises, public, and federated Clouds

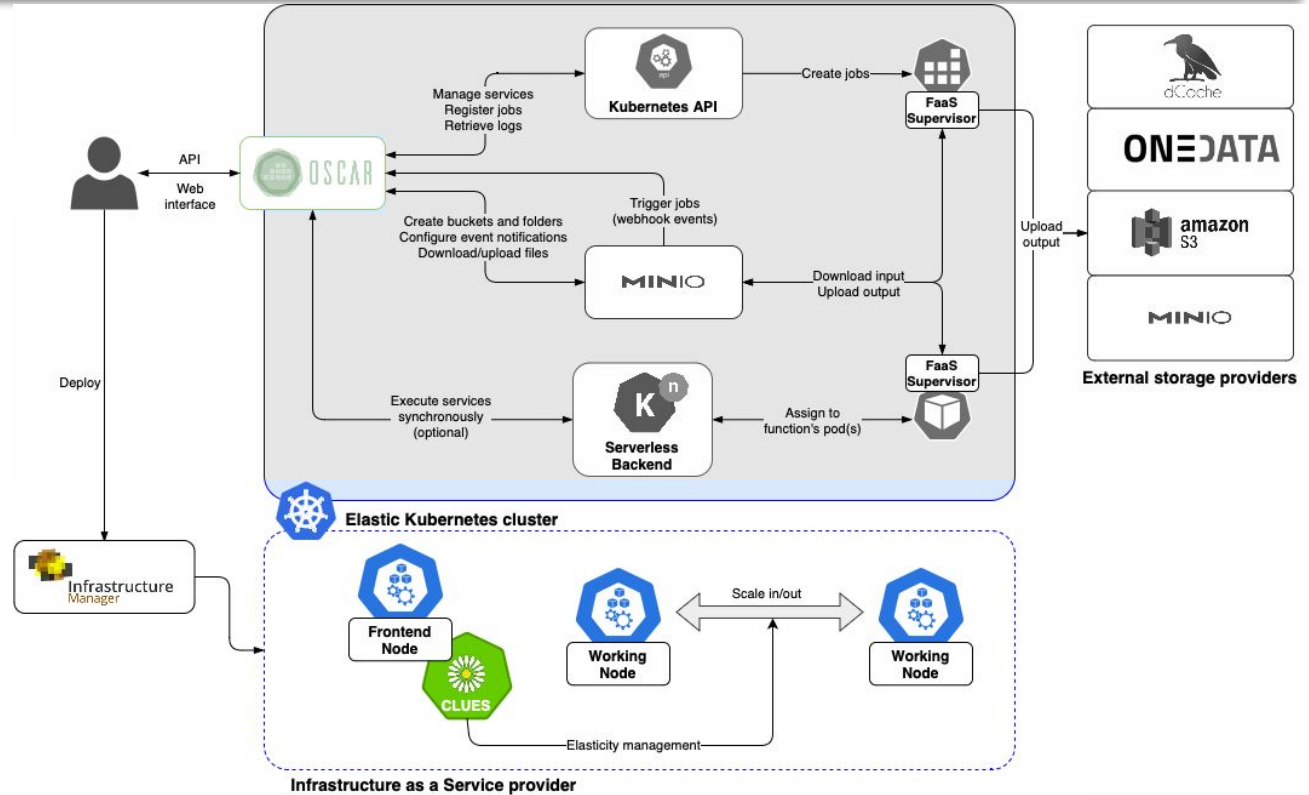# Scalable Event-driven Processing

**Multi-Cloud Support:**

Provision OSCAR clusters on on-premises, public and federated Clouds

Easily deployed via **IM** web interface:

- Add infra credentials
- Configure (step by step) your elastic OSCAR cluster



https://im.egi.eu

# Scalable Event-driven Processing

**Flexible Interfaces:**

- API REST
- Web Interface
- Command Line Interface

# Scalable Event-driven Processing

**Flexible Interfaces:**

- API REST
- Web Interface
- Command Line Interface





Apply a FDL file to create or edit services in clusters.

```
Usage:
  oscar-cli apply FDL_FILE [flags]

Aliases:
  apply, a

Flags:
      --config string   set the location of the config file (YAML or JSON)
  -h, --help            help for apply
```

# Scalable Event-driven Processing

**Support for Multiple Storage Back-ends:**

- MinIO, Amazon S3, OneData…

# Scalable Event-driven Processing

**Support for Multiple Storage Back-ends:**

- MinIO, Amazon S3, Onedata…
- Any Storage Service using WebDAV protocol via DCNiOS



https://github.com/interTwin-eu/dcnios

# Scalable Event-driven Processing

**Support for Multiple Storage Back-ends:**

- MinIO, Amazon S3, Onedata…

- Any Storage Service using WebDAV protocol via DCNiOS

- Decoupled FaaS Supervisor is the I/O manager

# Scalable Event-driven Processing

**Container-based services:**

- Customized runtime environments

# Scalable Event-driven Processing

**Event-driven processing:**

- When the user uploads data to the selected storage, that will generate an event that will trigger an OSCAR Job

# General Architecture

**interLink:**

- Developed by INFN, provides an abstraction for the execution of a Kubernetes pod on any remote host supporting containers

- Provides a gateway from OSCAR to delegate pod executions into HPC supercomputers

https://github.com/interTwin-eu/interlink

# Technical aspects of the integration

Changes in OSCAR for the integration with interLink:

- The pod needs to be assigned to the Virtual Kubelet node so interLink can offload it (via node selector)

- The pod needs to include some annotations for SLURM (HPC system)

- When running in the Cloud, the FaaS Supervisor is mounted as a volume in the pod. To be able to be exported by interLink to HPC, the FaaS Supervisor is passed as a sidecar

- Environment variables in the pod have to be codify (base64) for interLink to be able to export them

- Once in HPC, decodify the environment variables to run the job

All these changes are transparent to the user

# The use case

- Execution of an AI pipeline using itwinai (ML tool developed by CERN)

- The AI pipeline is configured to run the inference of a 3DGAN pretrained model to simulate particles in the HL-LHC (CERN)

- Data input from dCache, a distributed storage system for storing and retrieving huge amounts of data

https://github.com/interTwin-eu/itwinai/

https://www.dcache.org/

OSCAR

# The workflow

# The workflow

# The workflow

# The workflow

# The workflow

# The workflow

# The workflow

# The workflow

# The workflow

# The workflow

# The workflow
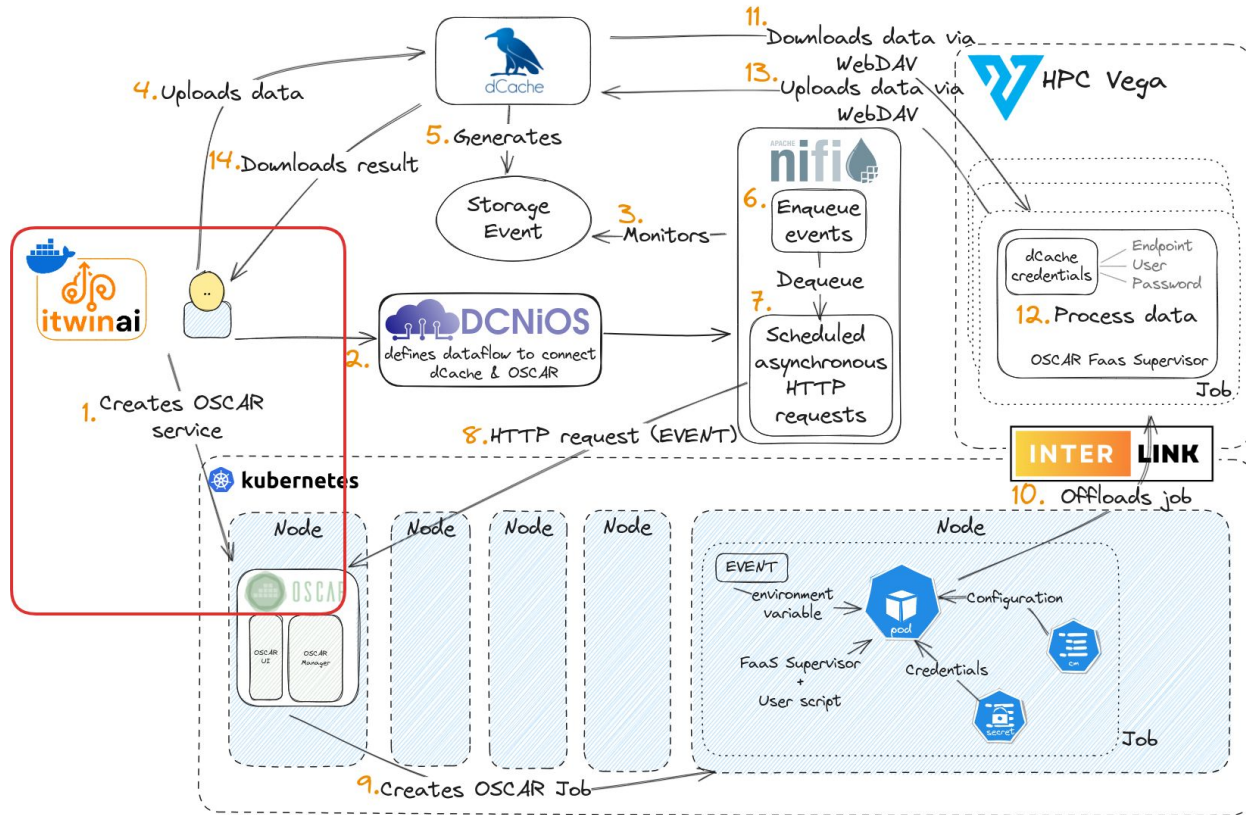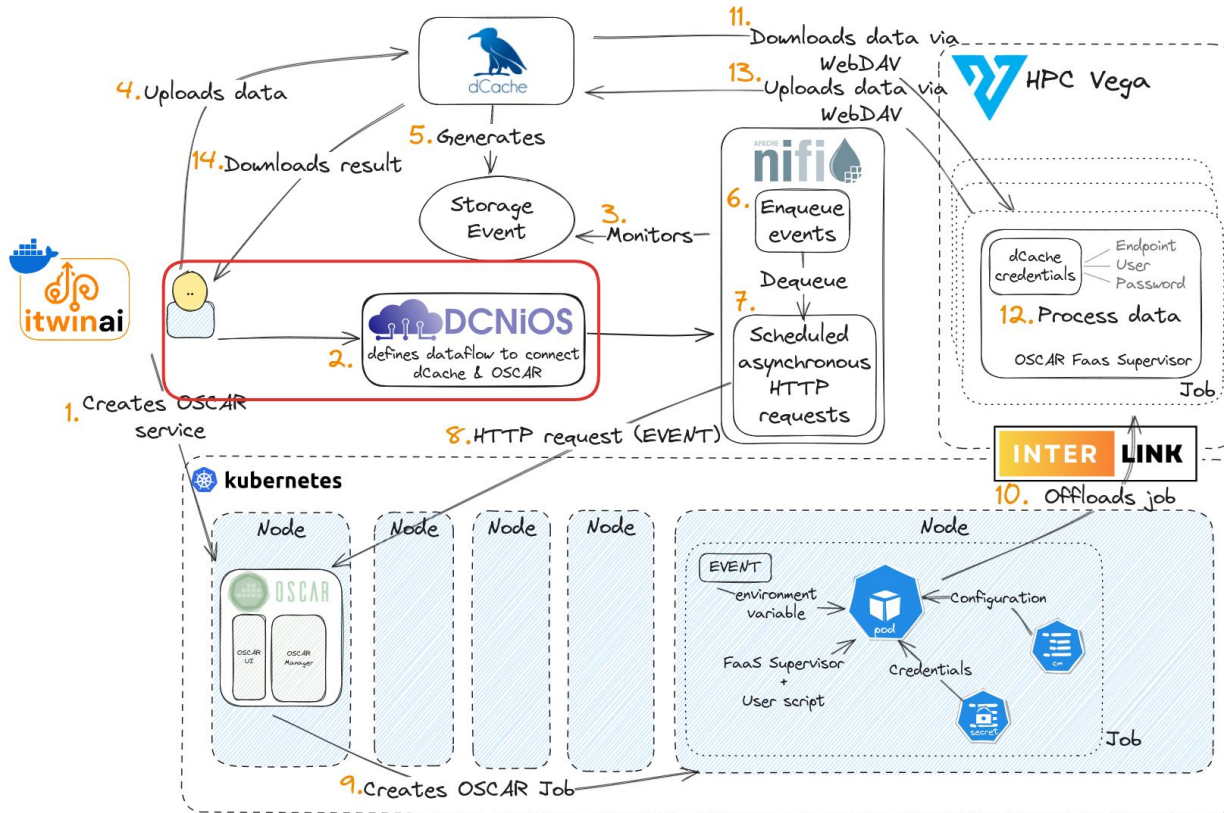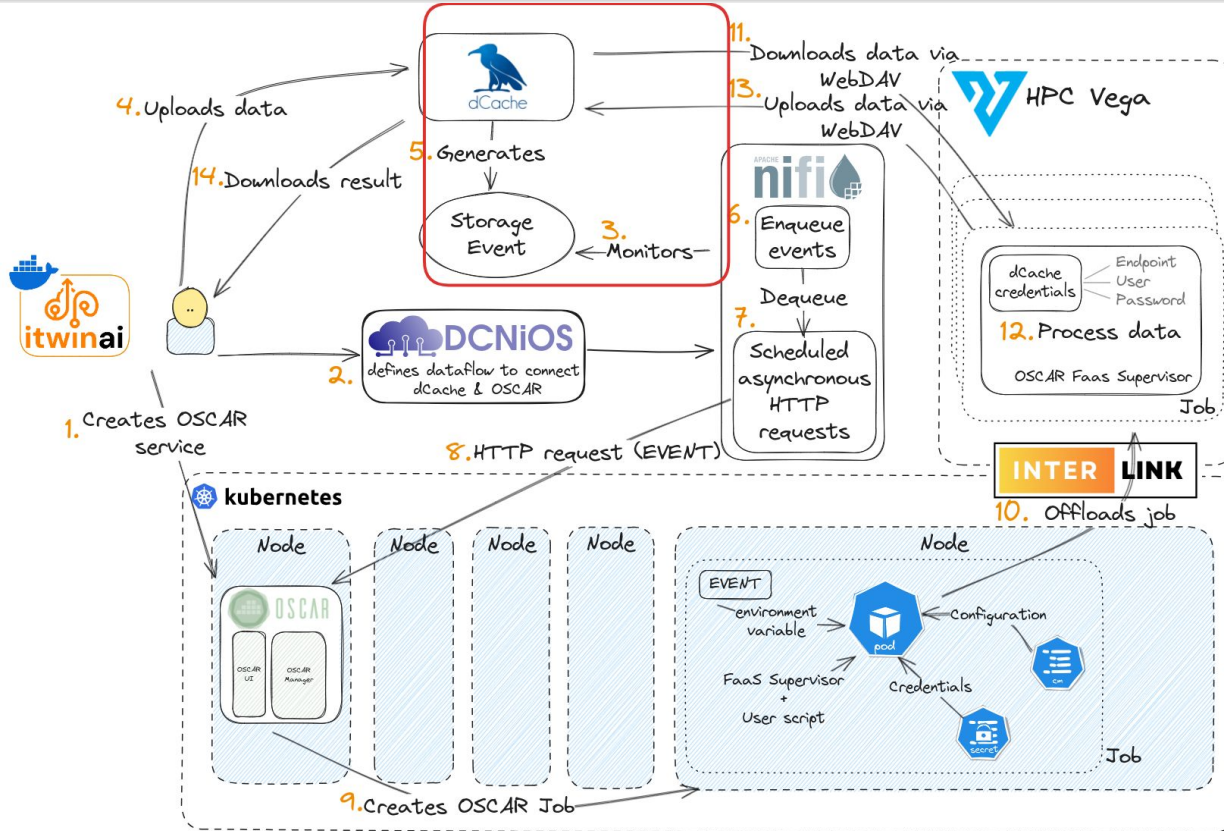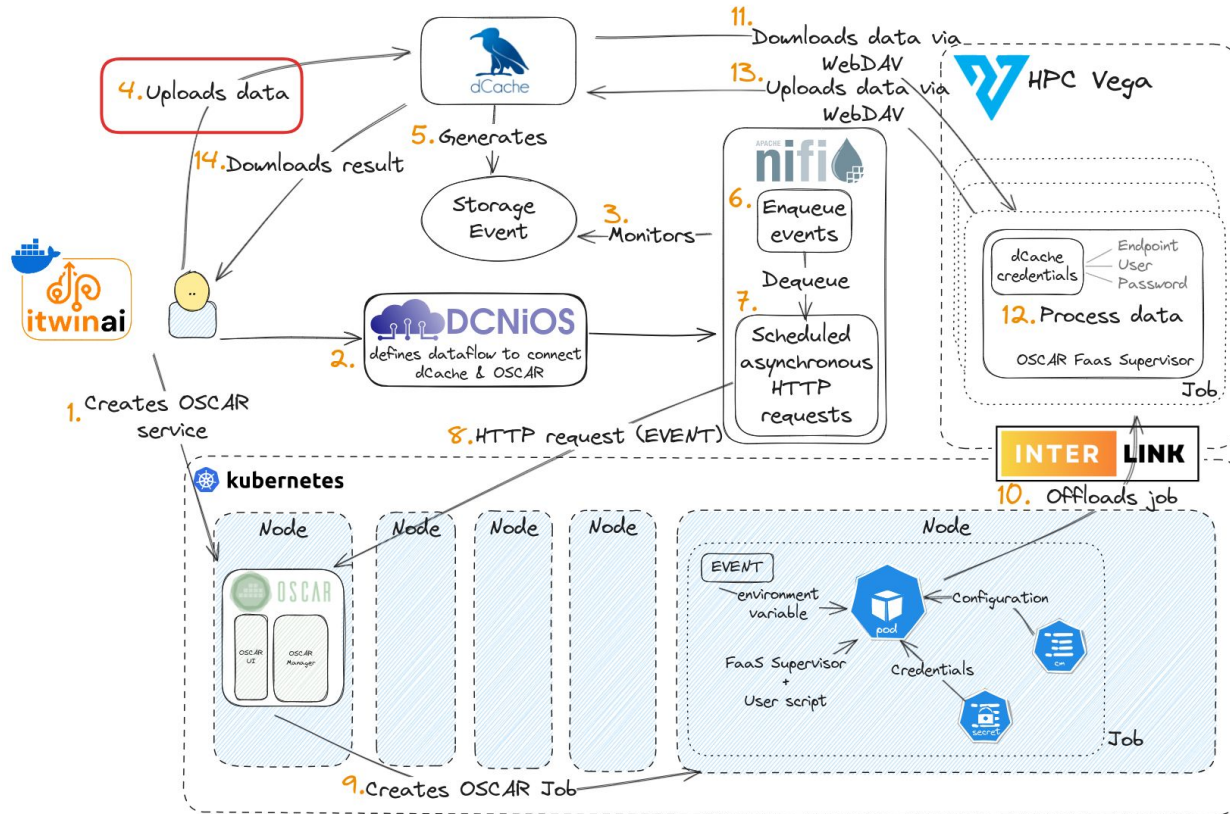
# The workflow

# The workflow

# The workflow

# The workflow

# The workflow



*This image is a creative interpretation based on a descriptive prompt and does not represent an actual scene or object. It is a product of imagination and artistic expression.*

# The workflow

# Results



- Hoefler et al., 2024, highlight the importance of bridging Cloud and HPC for resource-intensive workloads like **climate simulations** or **machine learning** processing

- Some key ideas:
  - Leverage containers
  - Improve communication
  - Enable access to data (I/O)

# Results



XaaS: Acceleration as a Service to Enable Productive High-Performance Cloud Computing

- Hoefler et al., 2024, highlight the importance of bridging Cloud and HPC for resource-intensive workloads like **climate simulations** or **machine learning** processing

- Some key ideas:
    - Leverage containers
    - Improve communication
    - Enable access to data (I/O)

We employ containers to be able to run the workload in the Cloud and in a HPC cluster

# Results



XaaS: Acceleration as a Service to Enable Productive High-Performance Cloud Computing

- Hoefler et al., 2024, highlight the importance of bridging Cloud and HPC for resource-intensive workloads like **climate simulations** or **machine learning** processing

- Some key ideas:
  - Leverage containers
  - Improve communication
  - Enable access to data (I/O)

We employ interLink to establish communication between Cloud and HPC

# Results



XaaS: Acceleration as a Service to Enable Productive High-Performance Cloud Computing

- Hoefler et al., 2024, highlight the importance of bridging Cloud and HPC for resource-intensive workloads like **climate simulations** or **machine learning** processing

- Some key ideas:
  - Leverage containers
  - Improve communication
  - Enable access to data (I/O) ⬅ We employ FaaS Supervisor to enable data access from Cloud and HPC clusters

OSCAR

# Results

- Access to HPC provides better resources and powerful GPUs

- We compared three different scenarios, Cloud-CPU, HPC-CPU, and HPC-GPU

- The processing time for each inference is stable

- HPC-CPU run the inference a ~ 30% faster than Cloud
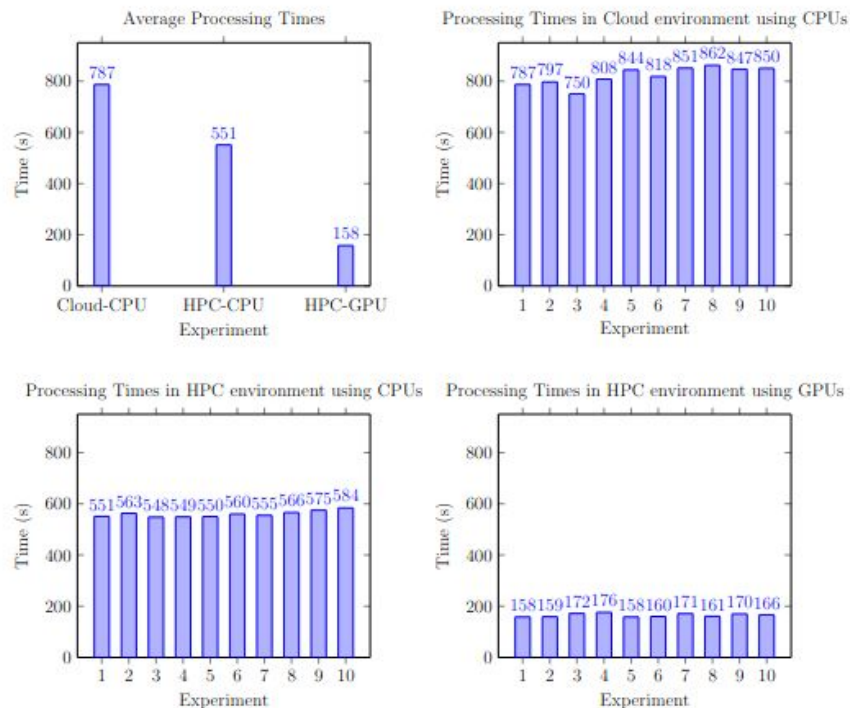
- HPC-GPU run the inference a ~ 80% faster than Cloud



Figure 5: Process time of the itwinai inference use case showing total execution time for each experiment and environment type.
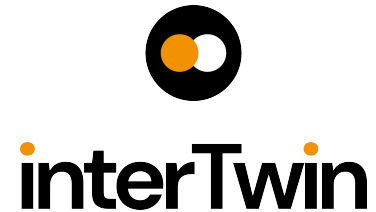
# Conclusions

- In this work we introduced an architecture to facilitate the integration of cloud computing and HPC for the execution of compute-intensive AI model inference, leveraging the event-driven serverless computing paradigm

- This integration facilitates access to HPC resources from OSCAR/Kubernetes clusters, thus making HPC resources seamlessly accessible to a wider range of users and applications via automated offloading through interLink

- A successful use case integrating dCache, Apache NiFi, OSCAR, interLink, and itwinai to support generative AI 3DGAN neural network model has been achieved, demonstrating the benefits of the approach by exploiting remote GPUs from an HPC facility from an OSCAR cluster running on a cloud infrastructure

In summary, the proposed system demonstrates the potential for offloading compute-intensive tasks to an HPC supercomputer, marking a significant step forward in
Bridging Cloud and HPC

# Acknowledgements

**interTwin**

# Thanks for your attention!

## Bridging Cloud and HPC for Scalable Event-driven Processing of AI Workloads

Estíbaliz Parcero, Sergio Langarita, Germán Moltó

Instituto de Instrumentación para Imagen Molecular (I3M), Universitat Politècnica de València, Valencia, Spain.
esparig@i3m.upv.es