# Easy Development and Execution of Workflows with eFlows4HPC

Rosa M Badia – BSC

21ST INTERNATIONAL SYMPOSIUM ON PARALLEL & DISTRIBUTED COMPUTING,
JULY 11-13, 2022 BASEL, SWITZERLAND

# Complex workflows and complex infrastructures

- EuroHPC aims at developing a World Class Supercomputing Ecosystem in Europe
    - Procuring and deploying pre-exascale and petascale systems in Europe
- These systems will be capable of running large and complex applications
- Applications demand the composition of HPC, artificial intelligence and data analytics

# EuroHPC systems

**eFlows4HPC**

| | Status | Country | Peak performance | Architecture |
|---|---|---|---|---|
| **Pre-Exascale** | | | | |
| LUMI | Operational | Finland | 552 petaflops | 64-core AMD EPYC™ CPUs + AMD Instinct™ GPU |
| Leonardo | Under construction | Italy | 322.6 petaflops | Intel Ice-Lake, Intel Sapphire Rapids + NVIDIA Ampere |
| MareNostrum 5 | Contract signed | Spain | 314 petaflops | Bull Sequana XH3000 and Lenovo ThinkSystem |
| **Petascale** | | | | |
| Meluxina | Operational | Luxembourg | 15 petaflops | AMD ... NVIDIA Ampere A100 |
| Vega | Operational | Sl... | | ...12 + Nvidia A100 |
| Karolina | | | | ...vidia A100 |
| Discoverer | ... | Bulgaria | 6 petaflops | AMD EPYC |
| Deucalion | Under construction | Portugal | 10 petaflops | A64FX, AMD EPYC +Nvidia Ampere |

*First European Exascale Supercomputer announced to be installed in Jülich*

# Main objectives

- Software tools stack that make it easier the development of workflows
  - HPC, AI + data analytics
  - Reactive and dynamic workflows
  - Efficient resource management
- HPC Workflows as a Service:
  - Mechanisms to make it easier
    the use and reuse of HPC
    by wider communities

# Outline

- Project architecture
  - Software stack
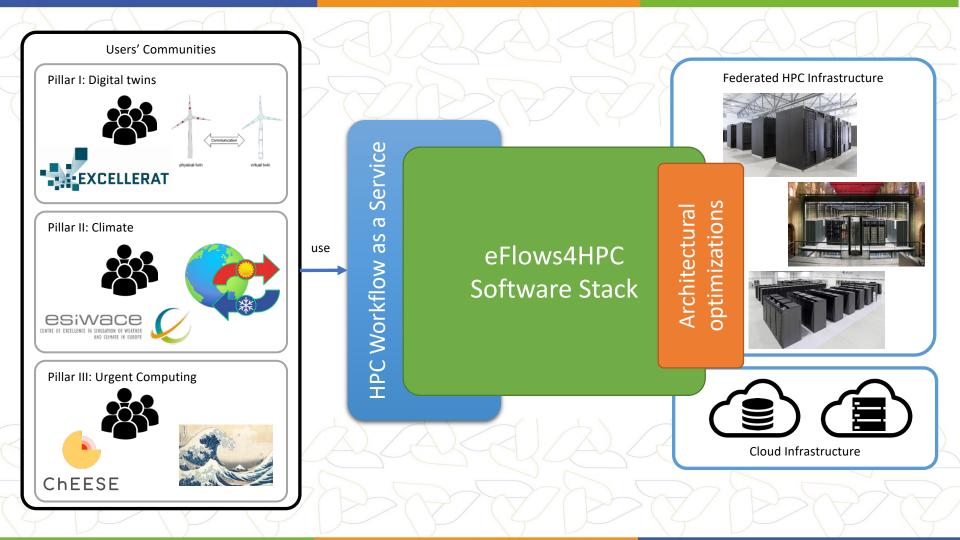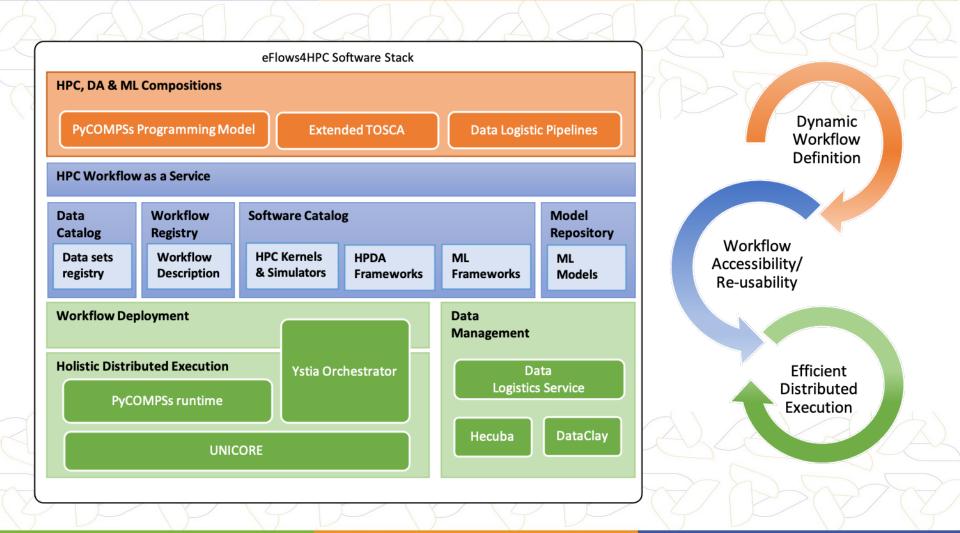  - HPC Workflows as a Service
- Pillar applications

# PROJECT ARCHITECTURE

# eFlows4HPC software stack and HPCWaaS

- Gateway services

  - Components deployed outside the computing infrastructure.

  - Managing external interactions and workflow lifecycle

- Runtime Components

  - Deployed inside the computing infrastructure to manage the  workflow execution

**Gateway Services**

**HPC WaaS Interface**

Alien4Cloud | Execution API

**Repositories for re-usability**

Data Catalog | Software Catalog | Workflow Registry

Ystia Orchestrator | Data Logistics Service | Image Creation

**Runtime Components**

PyCOMPSs | Hecuba | DataClay

**HPDA/ML Frameworks**

**Workflow**

Computing Infrastructure

# Workflow development overview



eFlows4HPC

eFlows4HPC Gateway Services

Data Catalog

Software Catalog

Workflow Registry

Alien4Cloud

Data Logistics Pipelines

PyCOMPSs Code

TOSCA Description

Dynamic Workflow Description

2. Store

3. Deploy

1. Create Workflow

4. share

endpoint to invoke the Workflow

Description of data movements as python functions. Input/output datasets described at Data Catalog

Computational Workflow as a simple python script. Invocation of software described in the Software Catalog

Topology of the components involved in the workflow lifecycle and their relationship.
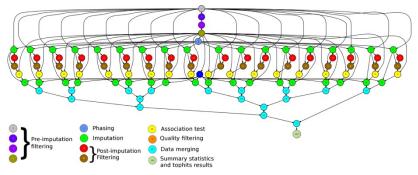
# Main element: Workflows in PyCOMPSs

- Sequential programming, parallel execution

- General purpose programming language + annotations/hints
  - To identify tasks and directionality of data

- Builds a task graph at runtime that express potential concurrency

- Tasks can be sequential, parallel (threaded or MPI)

- Offers to applications a shared memory illusion in a distributed system
  - The application can address larger data than storage space: support for Big Data apps
  - Support for persistent storage

- Agnostic of computing platform
  - Enabled by the runtime for clusters, clouds and container managed clusters

```
@task(c=INOUT)
def multiply(a, b, c):
    c += a*b
```
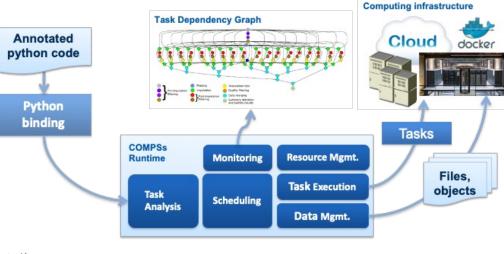


Pre-imputation filtering
Post-imputation Filtering

Phasing
Imputation
Data merging
Association test
Quality filtering
Summary statistics and tophits results

**Barcelona Supercomputing Center**
Centro Nacional de Supercomputación
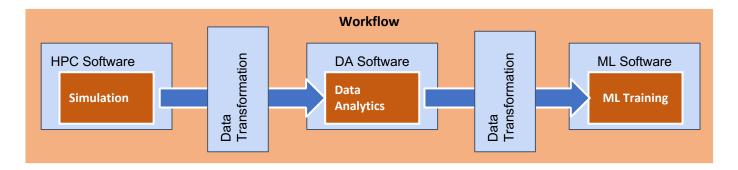
# PyCOMPSs features and runtime

- Support for tasks' constraints – support for heterogeneous infrastructure
- Support for tasks' faults and tasks' exceptions
  - Enlarges the dynamicity of the type of workflows that we support
- Streamed data
  - … and many others
- Runtime deployed as a distributed master-worker
- All data scheduling decisions and data transfers are performed by the runtime
- Support for elasticity



**Barcelona Supercomputing Center**
Centro Nacional de Supercomputación

# Interfaces to integrate HPC/DA/ML



- Goal:
  - Reduce glue code
  - Developer focuses in the functionality, not in the integration
  - Reusability

- Two paradigms:
  - Software integration
  - Data transformations

```
@data_tranformation(input_data, transformation description)
@software(invocation description)
def data_analytics (input_data, result):
    pass

#Worflow

simulation(input_cfg, sim_out)
data_analytics(sim_out, analysis_result)
ml_training(analysis_result, ml_model)
```
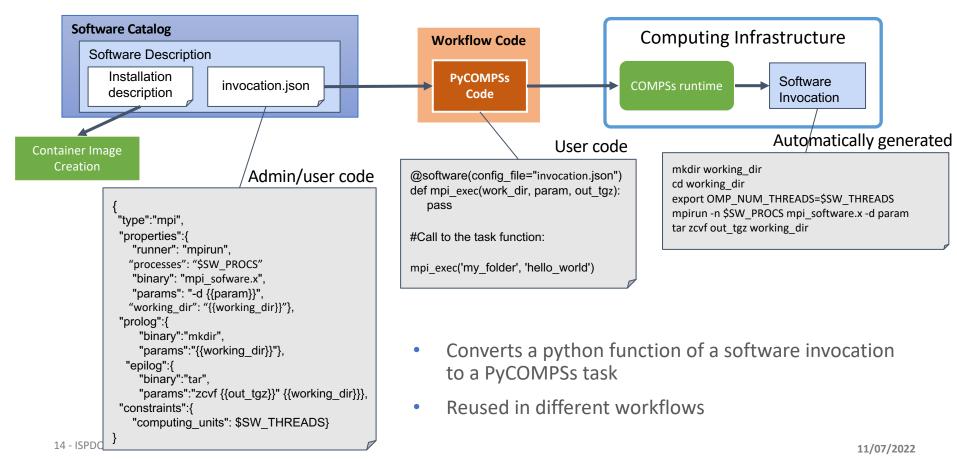
# Software Invocation description

eFlows4HPC

**Software Catalog**

Software Description

| Installation description | invocation.json |

**Workflow Code**

PyCOMPSs Code

## Computing Infrastructure

COMPSs runtime → Software Invocation

Container Image Creation

Admin/user code

```
{
  "type":"mpi",
  "properties":{
    "runner": "mpirun",
    "processes": "$SW_PROCS"
    "binary": "mpi_sofware.x",
    "params": "-d {{param}}",
    "working_dir": "{{working_dir}}"},
  "prolog":{
    "binary":"mkdir",
    "params":"{{working_dir}}"},
  "epilog":{
    "binary":"tar",
    "params":"zcvf {{out_tgz}}" {{working_dir}}},
  "constraints":{
    "computing_units": $SW_THREADS}
}
```

User code

```
@software(config_file="invocation.json")
def mpi_exec(work_dir, param, out_tgz):
    pass

#Call to the task function:

mpi_exec('my_folder', 'hello_world')
```

Automatically generated

```
mkdir working_dir
cd working_dir
export OMP_NUM_THREADS=$SW_THREADS
mpirun -n $SW_PROCS mpi_software.x -d param
tar zcvf out_tgz working_dir
```

- Converts a python function of a software invocation to a PyCOMPSs task

- Reused in different workflows

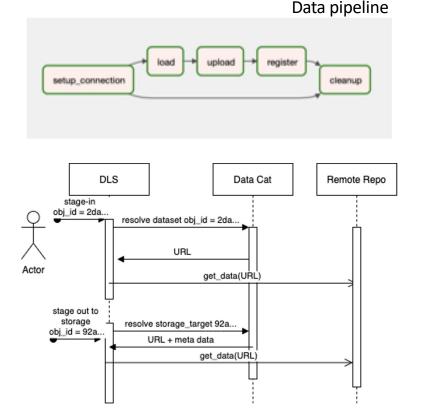# Data Catalogue and Data Logistics Service

Data pipeline

## Data Catalogue:

- Lists datasets used and created by the workflow according to FAIR principles

- Provides metadata to make data movement pipelines more generic

## Data Pipelines:

- Formalization of data movements for transparency and reusability

- Stage-in/out, image transfer

## Data Logistics Services (DLS):

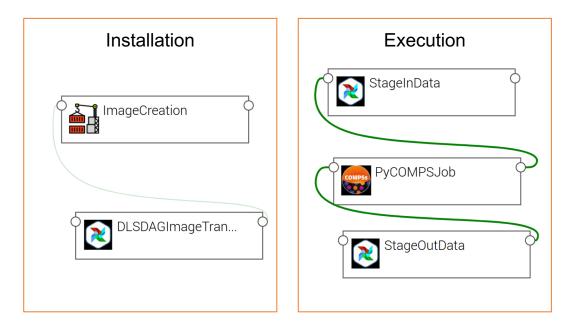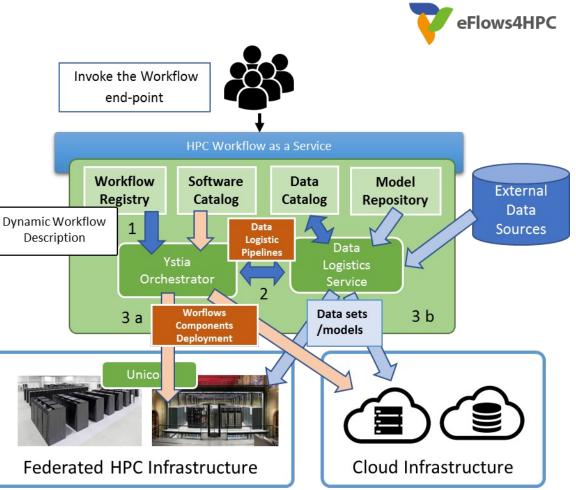- Performs the execution of data pipelines to fuel Pillars' computations

# TOSCA Modelization

Topology of the different components involved in the Workflow lifecycle
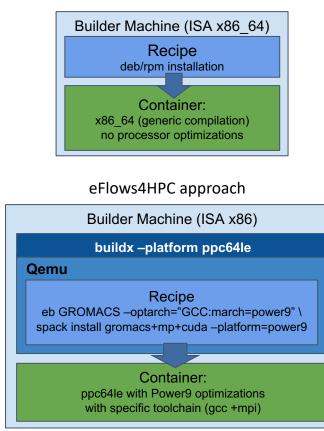
# Deployment

- Invocation of a workflow

- Deployment orchestrated by Ystia Orchestrator (Yorc)

- Workflow retrieved from registry

- Data Logistic Service – data stage-in and stage-out

  - Periodical transfers of data outside HPC systems

- Deployment of workflow components in the computing infrastructures

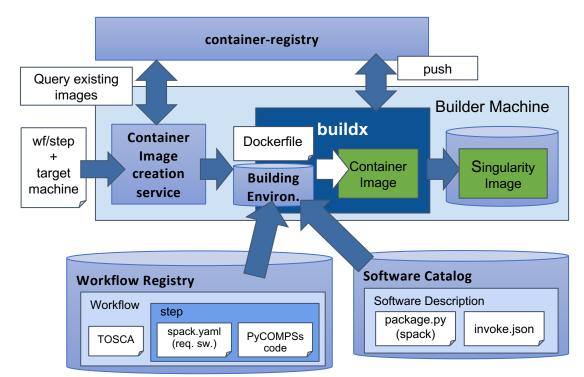  - HPC containers built with easybuild/Spack

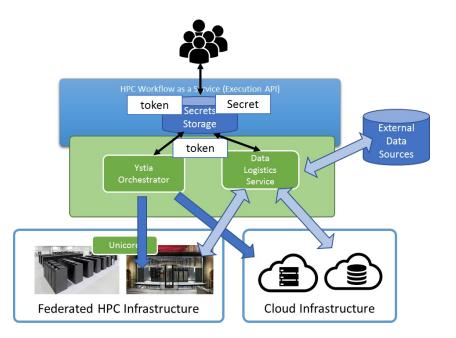# HPC Ready Containers



Standard container image creation

**Builder Machine (ISA x86_64)**

Recipe
deb/rpm installation

Container:
x86_64 (generic compilation)
no processor optimizations

eFlows4HPC approach

**Builder Machine (ISA x86)**

**buildx –platform ppc64le**

Qemu

Recipe
eb GROMACS –optarch="GCC:march=power9" \
spack install gromacs+mp+cuda –platform=power9

Container:
ppc64le with Power9 optimizations
with specific toolchain (gcc +mpi)

Service to automate the Container Image Creation

container-registry

Query existing images

push

wf/step + target machine

**Container Image creation service**

Dockerfile

**buildx**

Builder Machine

**Building Environ.**

Container Image

Singularity Image

**Workflow Registry**

Workflow

step

TOSCA

spack.yaml (req. sw.)

PyCOMPSs code

**Software Catalog**

Software Description

package.py (spack)
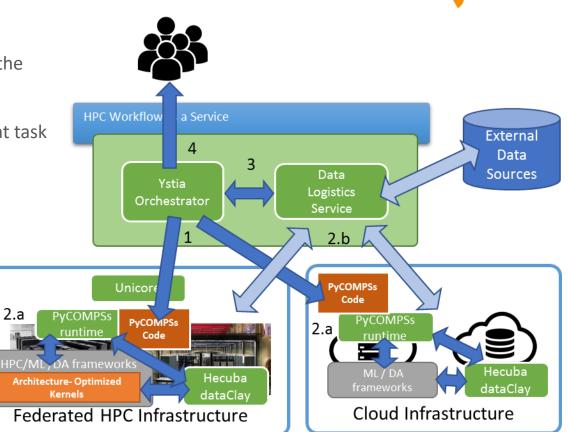
invoke.json

# Credential management

- Prior to executing the registered workflows, the users have to configure the infrastructure access credentials

- Usernames, public-key certificates, passwords

- Users' certificates managed by an Execution API
  - Provides a few methods to register and access credentials or generate a new secret
  - HashiCorp Vault for secret (SSH keys) management

- User authorizes adding credentials in the HPC cluster

- Credentials identified by a token attached to the user's workflow invocation.

eFlows4HPC



HPC Workflow as a Service (Execution API)

token    Secrets Storage    Secret

token

Ystia Orchestrator

Data Logistics Service

External Data Sources

Unicore

Federated HPC Infrastructure

Cloud Infrastructure

# Operation- Workflow Execution

- Submission of the execution of the workflow processes to the HPC infrastructure

- PyCOMPSs orchestrates different task types
  - HPC (MPI), ML, DA

- Dynamic execution
  - Runtime task-graph
  - Task-level FT
  - Exceptions

- Data management
  - Persistent storage
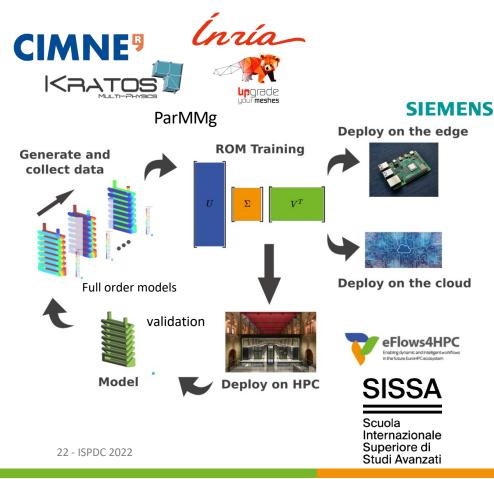
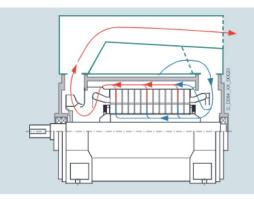- Optimized kernels
  - EPI, GPU, FPGA

# PILLAR APPLICATIONS

# Pillar I: Manufacturing



**ParMMg**

Pillar I focuses on the construction of DigitalTwins for the prototyping of complex manufactured objects:
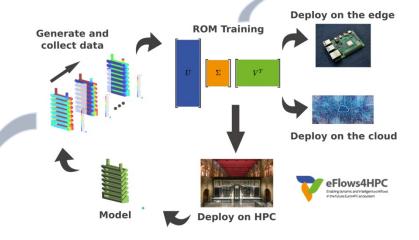
- Integrating state-of-the-art adaptive solvers with machine learning and data-mining

- Contributing to the Industry 4.0 vision

11/07/2022

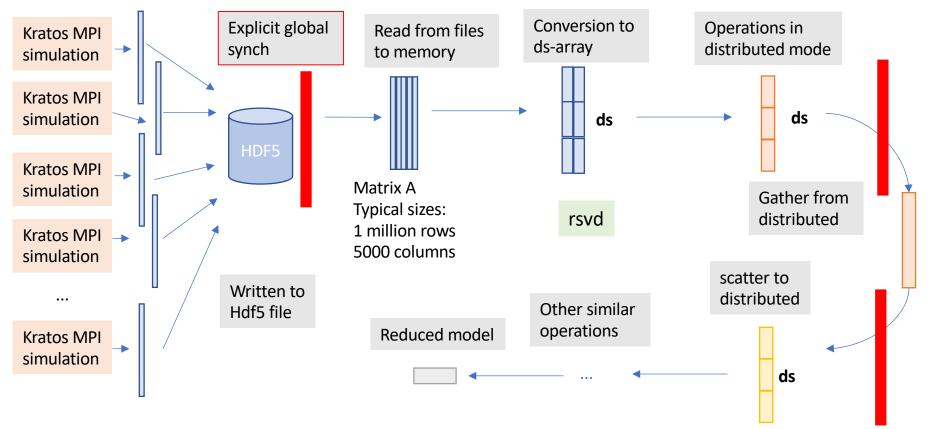# Integration of HPC and data analytics in a single workflow

```python
import dislib as ds

@constraint(computingUnits=8)
@task(Y_blocks={Type: COLLECTION_IN, Depth: 2},
def my_qr(Y_blocks):
    Y = np.block(Y_blocks)
    Q,R = np.linalg.qr(Y, mode='reduced')
    return Q,R


def rsvd(A, desired_rank):
    k = desired_rank
    ...
    Y = A @ Omega
    Q,R = my_qr(Y._blocks)
    Q=load_blocks_rechunk([Q], ...)
    ...
```
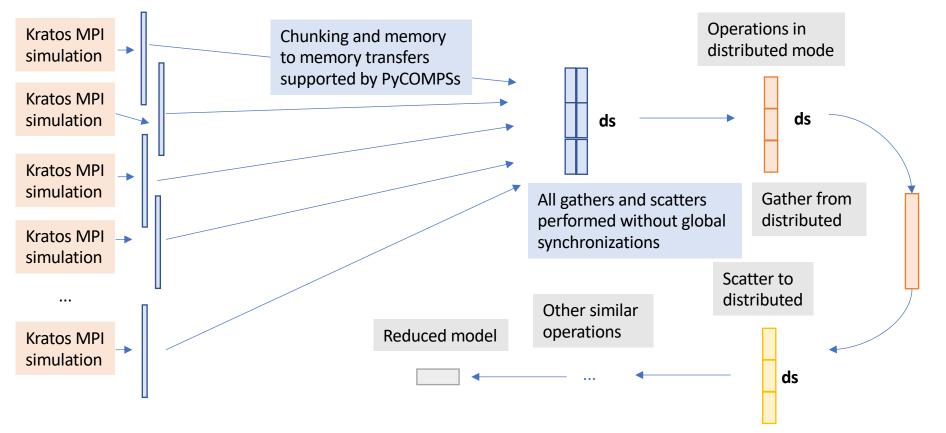
```python
@constraint(computing_units="8")
@mpi(runner="mpirun", processes="16")
@task(returns = np.array)
def ExecuteInstance_Task(
    ...
    current_parameters =
    ...
    simulation = GetTrain
    simulation.Run()
    return simulation.GetSnapshotsMatrix()
```

```python
....
for instance in range (0,TotalNumberOFCases):
    blocks.append(ExecuteInstance_Task(model, parameters, pars, inst
    ...
U, s = rsvd(A, desired_rank)
    ...
```

# Pillar I sample workflow – initial version

eFlows4HPC

Kratos MPI simulation

Kratos MPI simulation

Kratos MPI simulation

Kratos MPI simulation

...

Kratos MPI simulation

Explicit global synch

HDF5

Written to Hdf5 file

Read from files to memory

Matrix A
Typical sizes:
1 million rows
5000 columns

Conversion to ds-array

**ds**

rsvd

Operations in distributed mode

**ds**

Gather from distributed

scatter to distributed

**ds**

Reduced model

Other similar operations

...

# Pillar I sample workflow – optimized version

eFlows4HPC

Kratos MPI simulation

Kratos MPI simulation

Chunking and memory to memory transfers supported by PyCOMPSs

Kratos MPI simulation

Kratos MPI simulation

...

Kratos MPI simulation

**ds**

All gathers and scatters performed without global synchronizations

Operations in distributed mode

**ds**

Gather from distributed

Scatter to distributed

Reduced model

Other similar operations

...

**ds**

# Pillar II: Climate



ESM Ensemble run over HPC

Data Analytics

Dynamic (AI-assisted) workflow

HPDA & ML/DL

- Perform climate predictions: temperature, precipitation or wind speed

- AI-assisted pruning of the ESM workflow

- Study of Tropical Cyclones (TC) in the North Pacific, with in-situ analytics
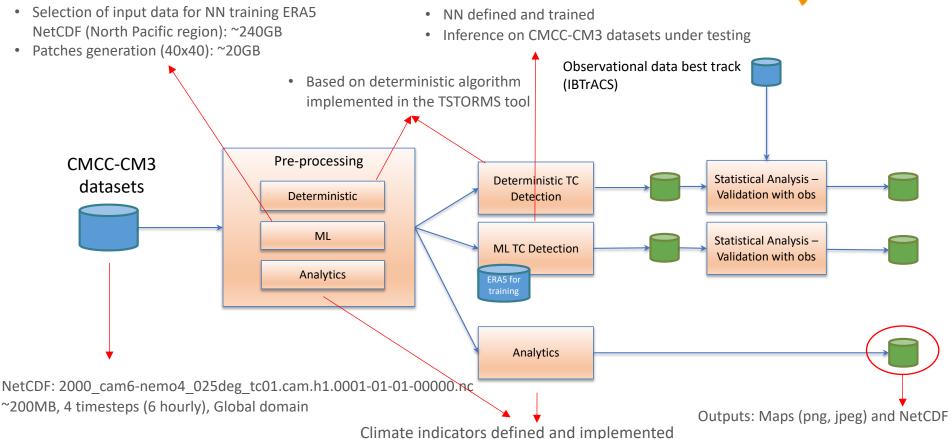
11/07/2022

# ESM Dynamic Workflow

- Ensemble of FESOM simulations

- Intermediate simulation results stored persistently to Hecuba

- Pruning mechanism cancels given simulations based on dynamic anaysis of data stored in Hecuba
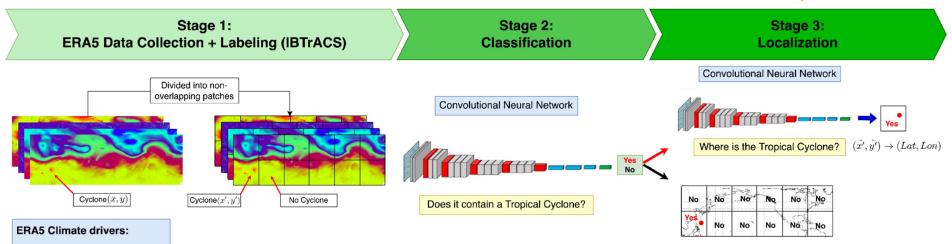
# Use case 1 – CMCC-CM3 based



- Selection of input data for NN training ERA5 NetCDF (North Pacific region): ~240GB
- Patches generation (40x40): ~20GB

- NN defined and trained
- Inference on CMCC-CM3 datasets under testing

- Based on deterministic algorithm implemented in the TSTORMS tool

Observational data best track (IBTrACS)

CMCC-CM3 datasets

Pre-processing

Deterministic

ML

Analytics

ERA5 for training

Deterministic TC Detection

ML TC Detection

Analytics

Statistical Analysis – Validation with obs

Statistical Analysis – Validation with obs

NetCDF: 2000_cam6-nemo4_025deg_tc01.cam.h1.0001-01-01-00000.nc
~200MB, 4 timesteps (6 hourly), Global domain

Climate indicators defined and implemented

Outputs: Maps (png, jpeg) and NetCDF

# Tropical Cyclones Detection ML Workflow: Training Phase

eFlows4HPC



**ERA5 Climate drivers:**
- 10 m wind gust
- Temperature at 500 hPa
- Temperature at 300 hPa
- Mean sea level pressure

- Three main stages:

  - STAGE 1: Gathering of ERA5 climatic maps and generation of patches containing at most 1 tropical Cyclone (TC) each

  - STAGE 2: Classification of TC presence/absence inside the patch

  - STAGE 3: Localization of TC center coordinates in the patches in which the TC was previously classified as present

# Pillar III: Urgent computing for natural hazards



Pillar III explores the modelling of natural catastrophes:

- Earthquakes and their associated tsunamis shortly after such an event is recorded

- Use of AI to estimate intensity maps

- Use of DA and AI tools to enhance event diagnostics

- Areas: Mediterranean basin, Mexico, Iceland and Chile

*Tsunami-HySEA GPU-based code*

# UCIS4EQ workflow: http services as tasks



Urgent Computing Integrated Services for
Earthquakes: UCIS4EQ

```python
@http(request="POST", resource="SalvusRun", ...')
@task(returns=1)
def run_salvus(event_id, trial, input, resources):
    """
    """
    pass
@http(request="POST", resource="cmt",...)
@task(returns=1)
def calculate_cmt(alert, event_id, domain, precmt):
    """
    """
    pass
...
```

```python
for alert in event['alerts']:
    cmts = calculate_cmt(alert, eid, domain, precmt)
    for cmt in cmts.keys():
        for slip in range(1, region['GPSetup']['trials']+1)
            ...
            rupture = compute_graves_pitarka(eid, alert, ...)
            inputs = build_input_parameters( eid, alert, ...)
            salvus_inputs = build_salvus_parameters( eid, path, ...)
            result = run_salvus( eid, path, ...)
            all_results.append(run_salvus_post(eid, result, ...))
result = run_salvus_plots(eid, basename, domain, resources)
```
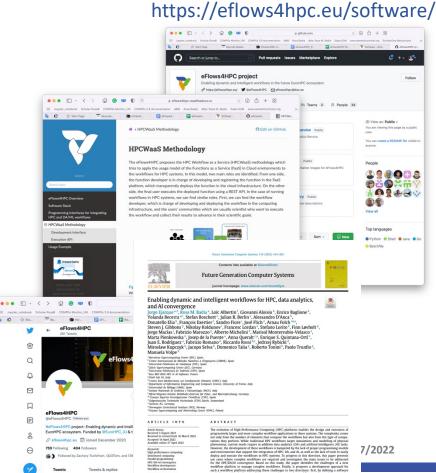
# Project main achievements

- Requirements and software architecture

- Definition and implementation of abstractions to support the integration of different stack components

- Design and development of a minimal workflow

- Design and first version of the HPCWaaS methodology

- Design and implementation of the Data Catalogue

- Definition of first version of Pillars' workflows. Implementation in progress

- First release of project software and documentation available

- Set of internal trainings about software stack components and HPCWaaS

- Good visibility: articles, keynote presentations, media

# Conclusions

- There is a need for providing tools for the development of complex workflows that include HPC modeling and simulation, artificial intelligence components and big data

- eFlows4HPC aims at providing a software stack that supports the development, deployment and execution of complex and dynamic workflows

- The HPCWaaS aims to provide a functionality similar for FaaS in cloud for complex workflows in HPC to make it easier the adoption of HPC technologies

# Project partners

**www.eFlows4HPC.eu**

@eFlows4HPC          eFlows4HPC Project

# Project WP Structure

# Project RoadMap