

# AQUA Extended: Evaluation Framework Bridging DestinE Simulations and ML Climate Emulator



Jose González-Abad<sup>2\*</sup>, Fernando Iglesias-Suarez<sup>1</sup>, Sergio Portilla<sup>1</sup>, Marvin Axness-Ferrando<sup>3</sup>, Miguel Castrillo<sup>3</sup>, Francisco Doblas-Reyes<sup>3,4</sup>, Amanda Duarte<sup>3</sup>, Aina Gaya-Avila<sup>3</sup>, Hernan Andres Gonzalez Gongora<sup>3</sup>, Jose Manuel Gutierrez Llorente<sup>2</sup>, Christian Lessig<sup>5</sup>, Sebastian Milinski<sup>5</sup>, Ankit Patnala<sup>6</sup>, Alejandro Peraza<sup>3</sup>, Antonio Perez<sup>1</sup>, Daniel San Martin Segura<sup>1</sup>, Jakob Schloer<sup>5</sup>, Martin Schultz<sup>6</sup>

<sup>1</sup>Predictia Intelligent Data Solutions (Predictia); <sup>2</sup>Consejo Superior de Investigaciones Científicas (CSIC); <sup>3</sup>Barcelona Supercomputing Center (BSC); <sup>4</sup>ICREA; <sup>5</sup>European Centre for Medium-Range Weather Forecasts (ECMWF); <sup>6</sup>ForschungsZentrum Jülich (FZJ); \*gonzabad@ifca.unican.es

## 1. Introduction

This framework extends AQUA (DestinE's evaluation tool for Climate Digital Twins) to bridge the gap between physics-based and machine learning (ML) climate simulations. Built on top of the **AQUA diagnostic library**, the framework allows reproducible evaluations across HPC platforms like MareNostrum5. It integrates both standard ML metrics (e.g., RMSE, ACC) and climate-specific diagnostics (e.g., ENSO, NAO, MJO).

## 2. AQUA

AQUA (Climate DT Applications for QUality Assessment) is a model evaluation framework designed for **running diagnostics on high-resolution climate models**, specifically for Climate DT climate simulations being part of Destination Earth activity. The package provides a flexible and efficient framework to process and analyze large volumes of climate data. With its modular design, AQUA offers seamless integration of core functions and a wide range of diagnostic tools that can be run in parallel.

AQUA efficiently handles high-resolution climate data using **catalogues**, **readers**, **fixers**, and **regridders**. It ensures **consistent data access**, **metadata standardisation**, and **regridding**.

```
sources:
  CESM2_r1i1p1f1:
    description: Earth system model: CESM2 (r1i1p1f1)
    driver: yaml_file_cat
    args:
      path: "[[CATALOG_DIR]]/CESM2_piControl_r1i1p1f1/main.yaml"
    metadata:
      description: Machine-learning based climate emulator (different configurations)
      driver: yaml_file_cat
      args:
        path: "[[CATALOG_DIR]]/emulator_inference/main.yaml"
```

Example of a catalogue entry.

```
from open import Reader
reader = ReaderCatalogue("CESM2", models="emulator_inference", exps="historical", sources="Experiment_4.3a",
                        fixer=None)
data = reader.references
```

Example of data loading with AQUA.

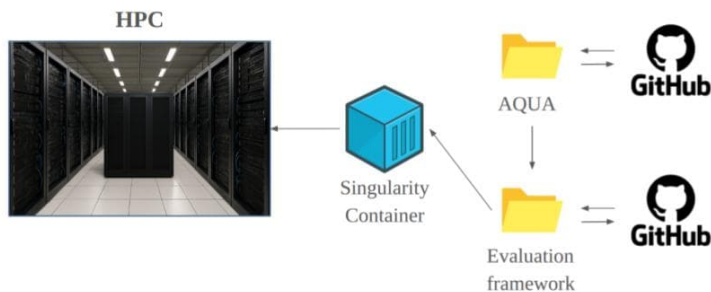
```
plugins:
  source:
    - module: intake_xarray
  sources:
    Experiment_4.3a:
      description: Simulations from the Experiment 4.3a.
      *** Configuration description ***
      driver: netcdf
      args:
        chunks:
          time: 1
        urlpath:
          - <data-path>
        metadata:
          filter_name: <filter-name>
    Experiment_5.1c:
      description: Simulations from the Experiment 5.1c.
      *** Configuration description ***
      driver: netcdf
      args:
        chunks:
          time: 1
        urlpath:
          - <data-path>
        metadata:
          filter_name: <filter-name>
```

Example of the definition of various sources.

## 3. AQUA Diagnostics

AQUA diagnostics are **modular, self-contained components** that evaluate specific aspects of model output. This design allows easy extension and integration with AQUA's data pipeline. Several diagnostics are already available, including tools for major variability modes (ENSO, MJO, NAO), climatological trends, and time series analysis—enabling robust, process-oriented evaluation of emulator performance.

## 4. Evaluation Framework



High-level structure of the framework and its deployment on HPC systems.

The **evaluation framework** provides all necessary tools to perform end-to-end validation of ML-based emulator outputs using AQUA. It includes:

- **Catalogues:** YAML-based mappings of models, experiments, and data sources for scalable and consistent data access.
- **Inference & Alignment Scripts:** Generate emulator outputs from model checkpoints (AnemoI) and align them with reference data.
- **Diagnostic Execution Logic:** Shell scripts and YAML configs to run AQUA diagnostics on selected emulator–reference pairs.
- **Supporting Resources:** Fixer rules and regridding weights for harmonised, comparable inputs across datasets.

## 5. Deployment on HPC Infrastructures

The evaluation framework supports **deployment across multiple HPC systems** and is fully containerised using **Singularity**, a solution well-suited for HPC environments. Singularity containers (provided directly with AQUA) ensure **compatibility with each AQUA release**, simplifying integration. This setup offers key advantages: consistent runtime environments, separation of inference and evaluation workflows, portability across systems, and compliance with strict HPC policies. By encapsulating dependencies and using **AQUA's official containers**, the framework remains **robust, reproducible, and easy to maintain**.

## 6. Extending Diagnostics for ML Emulation

To support more targeted evaluation of ML-based climate emulators, the framework is being **extended with new diagnostics specifically focused on ML performance**. These are developed following AQUA's modular structure, ensuring full compatibility and scalability. Initial additions include:

- **Root Mean Squared Error:** Measures spatial and temporal accuracy across variables.
- **Anomaly Correlation Coefficient:** Evaluates the emulator's ability to reproduce spatial patterns over time.

## 7. Performing the Evaluation

```
Configuration file for the RMSE diagnostic.
```

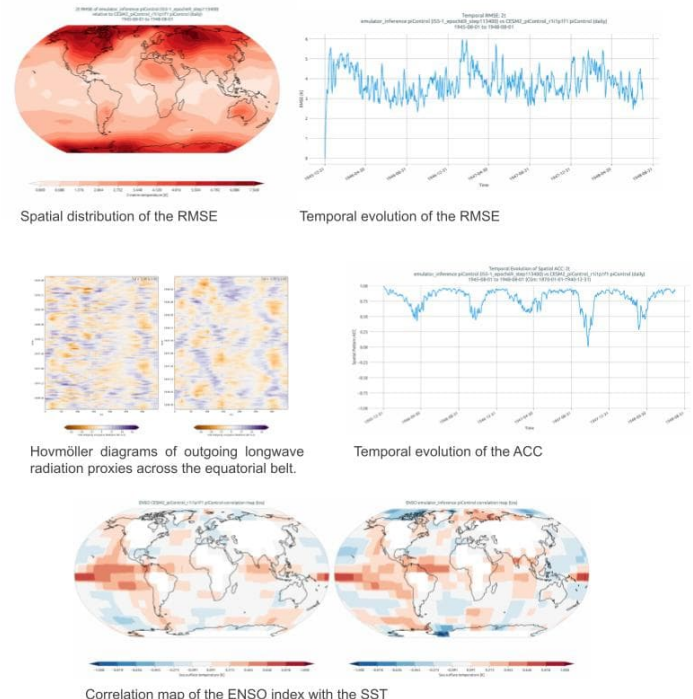
Configuration file for the RMSE diagnostic.

```
Configuration file for the ACC diagnostic.
```

Configuration file for the ACC diagnostic.

In the current version, the evaluation suite includes the following diagnostics:

- Time series comparison
- RMSE
- ACC
- Global biases
- Teleconnection Diagnostics
  - NAO
  - ENSO
  - MJO



## 8. Future Steps

Ongoing work focuses on **scaling the framework to high-resolution data**, adding **new diagnostics** (e.g., SSIM, spectral metrics), and **contributing them back to AQUA**. Usability and accessibility will be enhanced through **curated documentation**, **example configurations**, and **CLI tools**, making the framework easier to adopt and extend by the community.



Funded by  
the European Union

Destination Earth

implemented by

