# PYCMOR

Empowering scientists in CMORization

**Pavan Siligam**, *January 23, 2026*

Download PDF

# What is CMORization?

The process of converting climate model output data into a standardized format, in accordance with rules set by the Climate Model Intercomparison Projects (CMIP) and CF Conventions, ensuring consistency in variable-names, units, metadata, and file structure for easier analysis and comparison across different models.

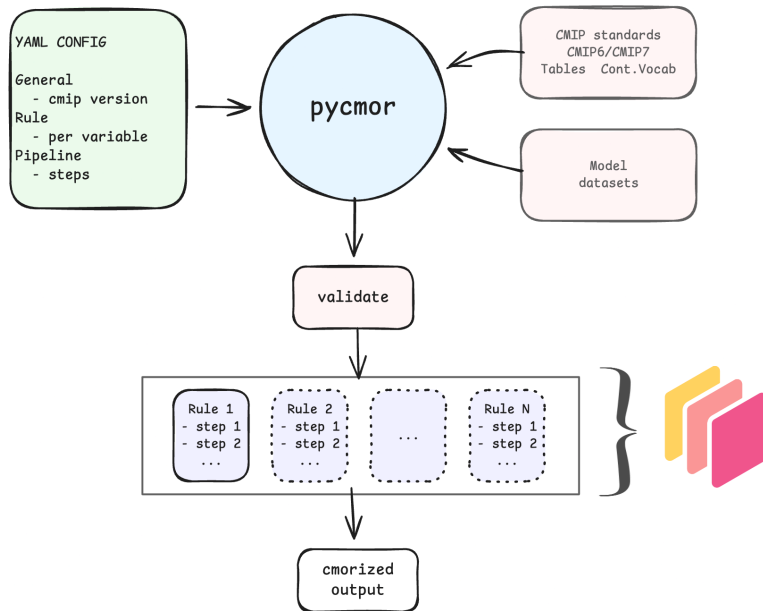For `pycmor` it means to provide compliance for the following

**CF Conventions**

- Ensure metadata is set
  - coordinate attributes
  - variable attributes
  - global attributes

**CMIP6/CMIP7 specification + Controlled Vocabulary**

- Ensure file structure
  - file(s) naming
  - directory structure
- match variable units
- match temporal resolution
- match dimensions
- bounds (time_bnds, lat_bnds, …)

# pycmor workflow

YAML CONFIG

General
 - cmip version
Rule
 - per variable
Pipeline
 - steps

pycmor

CMIP standards
CMIP6/CMIP7
Tables  Cont.Vocab

Model
datasets

validate

| Rule 1 | Rule 2 | | Rule N |
| - step 1 | - step 1 | ... | - step 1 |
| - step 2 | - step 2 | | - step 2 |
| ... | ... | | ... |

cmorized
output

# yaml configuration

Entry point

```
1   general:
2     name: "AWI-ESM-1-1-lr PI Control"
3     description: "CMOR configuration for AWIESM 1.1 LR"
4     maintainer: "pgierz"
5     email: "pgierz@awi.de"
6     cmor_version: "CMIP7"
7   pycmor:
8     # parallel: True
9     warn_on_no_rule: False
10    use_flox: True
11    dask_cluster: "slurm"
12    dask_cluster_scaling_mode: fixed
13    fixed_jobs: 1
14    # minimum_jobs: 8
15    # maximum_jobs: 30
16    # You can add your own path to the dimensionless mapping t
17    # If nothing is specified here, it will use the built-in o
18  inherit:
19    # Common attributes shared by all rules
20    activity_id: CMIP
21    institution_id: AWI
22    source id: AWI-CM-1-1-HR
```

| section | used for |
|---|---|
| general | cmor_version, CV_Dir, CMIP_Table_Dir, … |
| pycmor | application settings |
| inherit | common attributes shared by all rules |
| rules | **per variable settings**: cmor_variable, compound_name, model_variable, data_path, … |
| pipelines | steps to be executed for each variable (uses **Default pipeline** if not provided) |
| jobqueue | SLURM specific settings for launching workers |

# Pipeline Steps

## Default Steps in Pipeline

PyCMOR provides these standard processing steps.

```
1   STEPS = (
2       "pycmor.core.gather_inputs.load_mfdataset",
3       "pycmor.std_lib.generic.get_variable",
4       "pycmor.std_lib.add_vertical_bounds",
5       "pycmor.std_lib.timeaverage.timeavg",
6       "pycmor.std_lib.units.handle_unit_conversion",
7       "pycmor.std_lib.attributes.set_global",
8       "pycmor.std_lib.attributes.set_variable",
9       "pycmor.std_lib.attributes.set_coordinates",
10      "pycmor.std_lib.dimensions.map_dimensions",
11      "pycmor.core.caching.manual_checkpoint",
12      "pycmor.std_lib.generic.trigger_compute",
13      "pycmor.std_lib.generic.show_data",
14      "pycmor.std_lib.files.save_dataset",
15  )
```

## Adding Custom Steps

### 1. Create Custom Function

```
1   # custom_step.py
2   def custom_step(data, rule):
3       # Your custom processing logic
4       return data
```

### 2. Define Complete Pipeline

**Important**: When adding custom steps, you must specify the **entire pipeline** including both default and custom steps:

```
1   pipelines:
2     - name: my_custom_pipeline
3       steps:
4         - "pycmor.core.gather_inputs.load_mfdataset"
5         - "pycmor.std_lib.generic.get_variable"
6         # ... other default steps ...
7         - "script://path/to/custom_step.py:custom_step"
8         - "pycmor.std_lib.files.save_dataset"
```

### 3. Key Points

- Custom steps require full pipeline specification
- Mix default steps with your custom functions
- Use `script://path/to/file:function_name` syntax

# Coordinate Attributes

Make your coordinates CF-compliant automatically

## Eaxmple: no attribute information

```python
1    import xarray as xr
2    import numpy as np
3
4    # Dataset with coordinates
5    ds = xr.Dataset(
6        { 'tas': (['time', 'lat', 'lon'],
7                  np.random.rand(10, 90, 180)),
8        },
9        coords={
10       'time': np.arange(10),
11       'lat': np.linspace(-89.5, 89.5, 90),
12       'lon': np.linspace(0.5, 359.5, 180),
13   })
14
15   print(ds.lat.attrs)
16   {}
```

`pycmor` provides `set_coordinate_attributes`

In a processing pipeline ( `yaml` config):

```yaml
1    pipelines:
2      steps:
3        - "load_mfdataset"
4        # ... other steps
5        - "set_coordinate_attributes" # Add this step
6        # ... other steps
```

In scripts:

```python
1    from pycmor.std_lib.coordinate_attributes import set_
2
3    # Now coordinates have CF-compliant metadata
4    ds = set_coordinate_attributes(ds, rule)
5    ds['lat'].attrs
6    {
7      'standard_name': 'latitude',
8      'units': 'degrees_north',
9      'axis': 'Y',
10   }
```

# Units Handling in PyCMOR

## 🔬 Chemical units (element-aware)

**Automatic conversion for complex chemical units**

**Example:** `mmolC → kg`

**Previously (manual):**

```
1   mmolC → molC → gC → kgC
2       ÷1e3    ×12.0107    ÷1e3
```

**Now (automatic):**

```
1   # Live conversion log
2   2025-03-13 09:06:37 | INFO  | Converting: mmolC/m2/d → kg m-2 s-1
3   2025-03-13 09:06:37 | DEBUG | Chemical element detected: Carbon
4   2025-03-13 09:06:37 | DEBUG | Registering: molC = 12.0107 * g
```

✓ No manual molecular-weight handling
✓ CMIP-compliant physical units

## 📝 Unit configuration

### Wrong units in source data?

→ Define correct units in `model_unit` field in `yaml` file

### Dimensionless units?

→ Map them explicitly

```
1   # Conversion-only mappings
2   # Output always uses cmor_units
3   so:
4     "0.001": g/kg
5
6   intpp:        # Primary production by phytoplankton
7     "mol m-2 s-1": "molC m-2 s-1"
```

# Frequency Analysis & Processing

PyCMOR analyzes **source data frequency** and processes it to meet **CMIP table requirements**.

**Time averaging** (CMIP-aware):

- Driven by CMIP table frequency string
- **Methods**
  - Instantaneous → `resample( ... ).first()`
  - Mean → `resample( ... ).mean()`
  - Climatology → `groupby( ... ).mean("time")`

**Key principles:**

- ✅ **Downsampling supported** (no artificial data creation)
  - `Daily → Monthly` , `6-hourly → Daily` , `3-hourly → 6-hourly`

```
>>> # Month start dates with day offsets
>>> month_start_day_offset = np.array([
...     "2020-01-01",
...     "2020-02-01",
...     "2020-03-04",  # ← 3-day offset
...     "2020-04-01",
...     "2020-05-01",
... ]).astype("datetime64[s]")
...
>>> xarray.infer_freq(month_start_day_offset)
None
>>> infer_frequency(month_start_day_offset, return_metadata=True)
Inferred Frequency : M
Median Δ (days)    : 30.50
Regular Spacing    : ✅
Status             : valid    # ← below threshold limit
>>>
>>> infer_frequency(month_start_day_offset, return_metadata=True, strict=True)
Inferred Frequency : M
Regular Spacing    : ❌
Strict Mode        : ✅
```