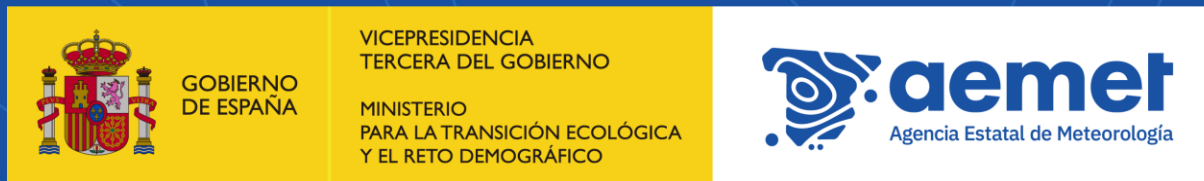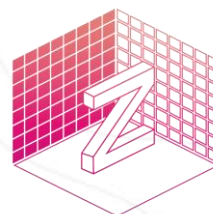# Unlocking Large Datasets on Cloud-Native Data Workflows

*Lessons on using Icechunk, Virtualizar, and Kerchunk to enable Single-Point Access to Distributed Cloud Files through xarray*

# *Accessing NWC SAF CF buckets with kerchunk, virtualizarr & icechunk.*

**Short guide for:** https://gitlab.aemet.es/jllisov/vzarr.git

GOBIERNO DE ESPAÑA

VICEPRESIDENCIA TERCERA DEL GOBIERNO

MINISTERIO PARA LA TRANSICIÓN ECOLÓGICA Y EL RETO DEMOGRÁFICO

aemet — Agencia Estatal de Meteorología

EUMETSAT NWC SAF

Explore                                                                Sign in

Search or go to...

jllisov@aemet.es  /  **Vzarr**

**Project**

V  Vzarr

88  Manage               >

📅  Plan                  >

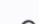</>  Code                 >

🚀  Build                 >

Deploy                   >

Operate                  >

Monitor                  >

Analyze                  >

# Vzarr

## Accessing NWC SAF buckets with kerchunk/virtualizarr/icechunk

## Description

This project is a collection of notebooks that demonstrate how to build virtual data cubes from NWC SAF data buckets. It focuses specifically on the bucket containing NWC SAF GEO outputs hosted on the European Weather Cloud (EWC). The bucket used is a precursor to the future MTG format of NWC SAF, containing files with a time dimension. This bucket is publicly accessible within the EWC infrastructure. The bucket name is:

`nwc-saf.0-degree.level-2-cf`

⚠️ Important: These notebooks are designed to be executed inside the EWC environment. Pleas~ these notebook should be executed inside the EWC.

## ⚙️ Requirements

- EWC access (https://www.europeanweather.cloud/)
- Conda (Anaconda/Miniconda)

## Installation

1. Clone the repository
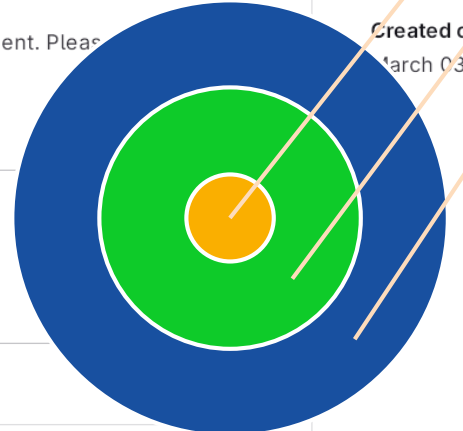
**Project information**

⚬ 45 Commits

⑂ 1 Branch

🏷️ 0 Tags

📄 README

⚖️ MIT License

Auto DevOps enabled

Created on
March 03, 2025

**kerchunk**

**virtualizarr**

**icechunk**

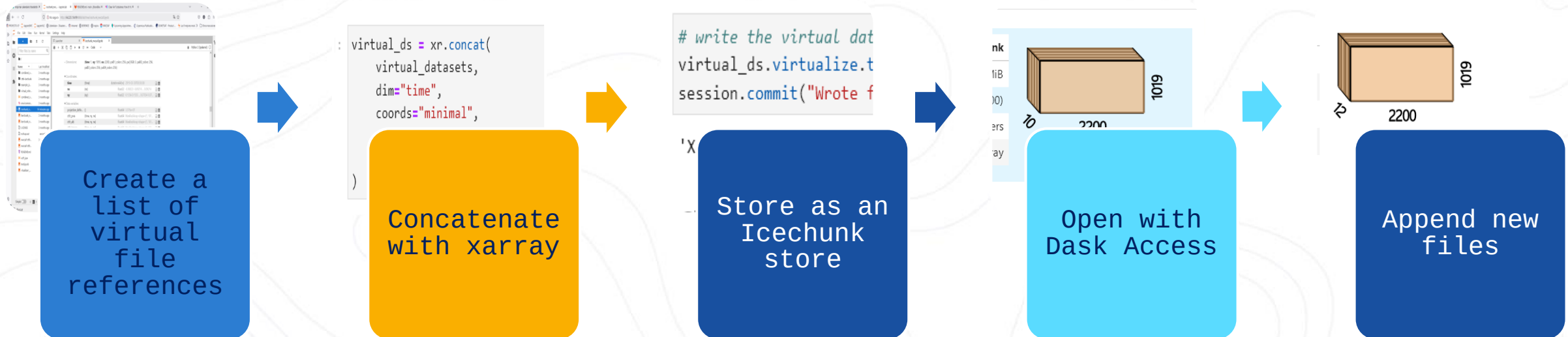The project is a set of notebooks teaching how to build <u>virtual data cubes</u> from NWC SAF CF buckets.

Virtual access implies not duplicating the files to build the data cube.

The format of the files (for the lessons) is the precursor for MTG format of the NWC SAF.
This bucket is public inside the EWC (***"nwc-saf.0-degree.level-2-cf"***)

The full lessons can only be executed on EWC.

**5**

**Retrieving from the store**

Appending to the store

# Append

Here we follow the same previous steps to create a virtual dataset, but we add an append_dim argument to the to_icechunk function. This will allow to expand the reference in the store.

```python
[129]:
open_dataset_options = {"chunks": {}}  # opens passed to xarray
so = dict(endpoint_url=S3_ENDPOINT_URL, anon=True, default_fill_cache=False, default_cache_type="none")

virtual_datasets_a = [
    open_virtual_dataset(url, indexes={}, reader_options={"storage_options": so,  "open_dataset_options":open_dataset_options},
                         loadable_variables=['ny', 'nx', 'time'], decode_times=True)
    for url in ctth_files[10:12]
]
```

```python
[130]:
virtual_ds_a = xr.concat(
    virtual_datasets_a,
    dim="time",
    coords="minimal",
    compat="override",
    combine_attrs="override",
)
```

```python
[132]: append_session = repo.writable_session("main")
```

```python
[133]: virtual_ds_a.virtualize.to_icechunk(append_session.store, append_dim="time")
```

```python
[134]: append_session.commit("wrote 2 more slots of data")
```

# Features

The file colection is accessed via xarray.

Data Access:

xarray datasets

- *Dask array*

Not file duplication

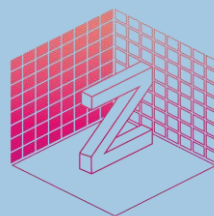An index file is build (json/parquet)

# Lessons learned

The time dimension to stack the files is « needed ». Use CF time dimension

The spatial dimensions should match, this should be forced

icechunch offers the best performance & is more pythonic

Use parquet for big collections

jllisov@aemet.es