# An Online (Compression) Laboratory for Climate Science and Meteorology

## Exploring Data Compression and more from the Comfort of your web browser

**ESiWACE3 WP3.2: Juniper Tyree[1], Sara Faghih-Naini[2], Peter Dueben[2], Karsten Peters-von Gehlen[3], Heikki J. Järvinen[1]**

[1] University of Helsinki, [2] ECMWF, [3] DKRZ

CASTIEL2 Code of the Month, 19.02.2025

Why do we need but avoid to use Lossy Compression?

The Compression Laboratory

An Online Open CliMet Science Laboratory

Future Outlook

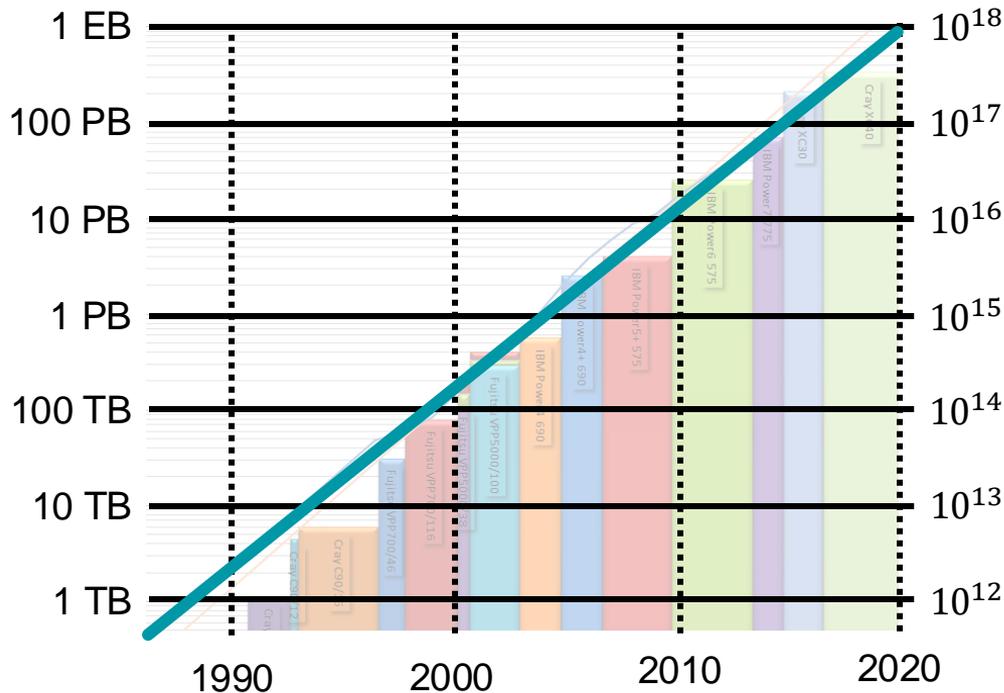Why do we need but avoid to use Lossy Compression?

The Compression Laboratory

An Online Open CliMet Science Laboratory

Future Outlook

# The Problem − Exponential data growth

### ECMWF Archive Size

Weather and climate science has reached the exascale for data storage

Data needs to be produced, stored, analyzed, transferred, published, and archived, which all have I/O costs

Cost per stored petabyte (PB) is ~10k € for cold (tapes) and ~100k € for hot (disk) storage

https://climate.copernicus.eu/sites/default/files/2019-11/Martin%20Palkovic.pdf

# The Solution? – Lossy compression can provide huge savings

## Compressing atmospheric data into its real information content

Milan Klöwer ✉, Miha Razinger, Juan J. Dominguez, Peter D. Düben & Tim N. Palmer

13k Accesses | 63 Altmetric | Metrics

ⓘ A preprint version of the article is available at Research Square.

### Abstract

Hundreds of petabytes are produced annually at weather and climate forecast centers worldwide. Compression is essential to reduce storage and to facilitate data sharing. Current techniques do not distinguish the real from the false information in data, leaving the level of meaningful precision unassessed. Here we define the bitwise real information content from information theory for the Copernicus Atmospheric Monitoring Service (CAMS). Most variables contain fewer than 7 bits of real information per value and are highly compressible due to spatio-temporal correlation. Rounding bits without real information to zero facilitates lossless compression algorithms and encodes the uncertainty within the data itself. All CAMS data are 17× compressed relative to 64-bit floats, while preserving 99% of real information. Combined with four-dimensional compression, factors beyond 60× are achieved. A data compression Turing test is proposed to optimize compressibility while minimizing information loss for the end use of weather and climate forecast data.

## COMPRESSING MULTIDIMENSIONAL WEATHER AND CLIMATE DATA INTO NEURAL NETWORKS

**Langwen Huang**
Department of Computer Science
ETH Zurich
langwen.huang@inf.ethz.ch

**Torsten Hoefler**
Department of Computer Science
ETH Zurich
torsten.hoefler@inf.ethz.ch

### ABSTRACT

Weather and climate simulations produce petabytes of high-resolution data that are later analyzed by researchers in order to understand climate change or severe weather. We propose a new method of compressing this multidimensional weather and climate data: a coordinate-based neural network is trained to overfit the data, and the resulting parameters are taken as a compact representation of the original grid-based data. While compression ratios range from 300x to more than 3,000x, our method outperforms the state-of-the-art compressor SZ3 in terms of weighted RMSE, MAE. It can faithfully preserve important large scale atmosphere structures and does not introduce artifacts. When using the resulting neural network as a 790x compressed dataloader to train the WeatherBench forecasting model, its RMSE increases by less than 2%. The three orders of magnitude compression democratizes access to high-resolution climate data and enables numerous new research directions.

# The Problem – Lossy compression is scary

**Computer scientist:** "I can give you great savings by using lossy compression."

**Domain scientist:** "Nice! … Actually, did you remember to keep the …

- **logarithmic error** norms small for **specific humidity**
- $L_2$ **error** small for **temperature**
- **budgets** correct when **integrating** over long climate trajectories
- $L_\infty$ **norm** small for **extreme precipitation events**
- **integral in the vertical** correct even though it sums across **several orders of magnitude**
- delta at high precision for $CO_2$
  **…**

And how does lossy data compression change the **enthalpy budget** when cold rain is entering a warm ocean and changing the vertical layering of the vertical column in an ocean model???"

**Computer scientist:** "Goodbye."

# The Problem – We need lossy compression but avoid to use it

Lossy Compression is crucial to handle Data Growth

Safety Requirements for lossy compression
in weather and climate science are vague

**Lack of Trust** in Lossy Compression
to uphold the safety requirements
for everyday domain-scientist users

# The Roadmap to *Fearless* Lossy Compression

Lossy Compression is crucial to work with Data Growth

Establish clear Safety Requirements

Simplify building **Trust** in Lossy Compression:
(1) Benchmark existing methods for safety and performance
(2) Provide safeguards to make any compressor safe
(3) **Convince yourself** in an Online Compression Laboratory

Why do we need but avoid to use Lossy Compression?

The Compression Laboratory

An Online Open CliMet Science Laboratory

Future Outlook

# A brief overview of the Compression Laboratory Notebooks

How to compress?

What compression?

Which datasets?

Examples case studies



Example: Introduction to loading and compressing data
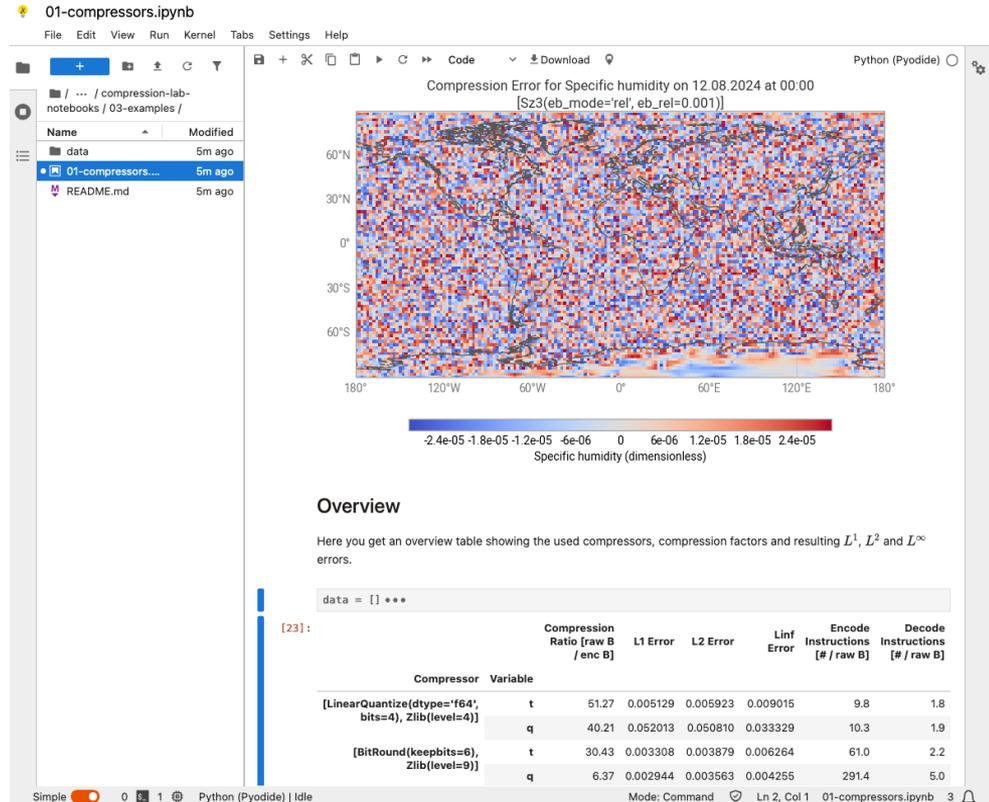
# A brief overview of the Compression Laboratory Notebooks

How to compress?

What compression?

Which datasets?

Examples case studies



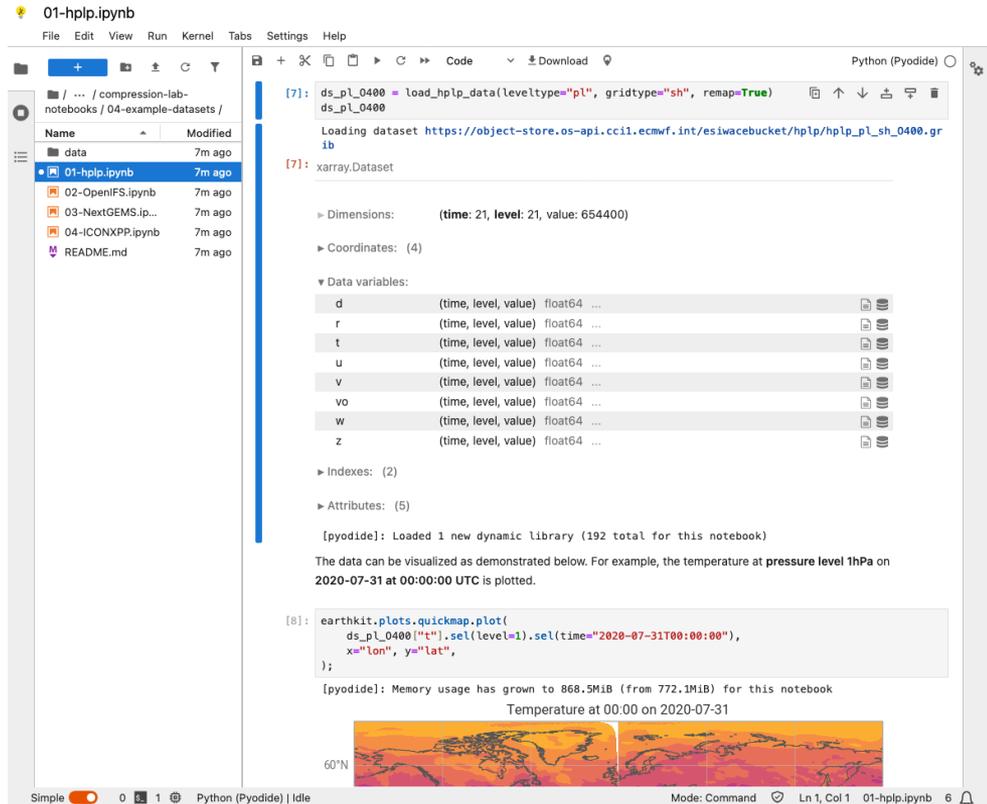Example: Comparing lossy compressor performance

# A brief overview of the Compression Laboratory Notebooks

How to compress?

What compression?

Which datasets?

Examples case studies



Example: Opening large remote datasets

# A brief overview of the Compression Laboratory Notebooks

How to compress?

What compression?

Which datasets?

Examples case studies

*Work in Progress:*

Compressing precipitation

with Clément Bouvier (UH), Joona Cornér (UH), Peter Dueben (ECMWF), and Milan Klöwer (Oxford)

Compressing model states during runtime

with Madeleine Ekblom (FMI) and Milan Klöwer (Oxford)

demo.ipynb

File  Edit  View  Run  Kernel  Tabs  Settings  Help

/ 02-19-2025-05-41-34-9e598c0b /
climet-lab-demo /

Filter files by name

| Name | Modified |
|---|---|
| • 🖼 demo.ipynb | now |
| 📄 LICENSE.txt | 2s ago |
| 📄 t2m.nc | now |
| 🐍 utils.py | 17m ago |

Code

Download

Python (Pyodide)

```
[4]: t = xr.open_dataset("t2m.nc").t2m.rename(dict(valid_time="time")); t
```

```
[4]: xarray.DataArray  't2m'  (time: 1, latitude: 721, longitude: 1440)
```

```
[5]: chart = earthkit.plots.quickmap.plot(t, units="degC")
     chart.title("Uncompressed {variable_name} on {time:%d.%m.%Y at %H:%M}");
```

Uncompressed 2 metre temperature on 01.12.2012 at 14:00

2 metre temperature (°C)

Simple  0  $_  1  ⚙  Python (Pyodide) | Idle

Mode: Command    Ln 1, Col 1    demo.ipynb    3

File  Edit  View  Run  Kernel  Tabs  Settings  Help

Python (Pyodide) ○

▣ / ···
/ 02-19-2025-05-41-34-9e598c0b /
climet-lab-demo /

Filter files by name

| Name ▲ | Modified |
|---|---|
| • ▣ demo.ipynb | now |
| 📄 LICENSE.txt | 2s ago |
| 📄 t2m.nc | now |
| 🐍 utils.py | 17m ago |

Code ▾   ⬇ Download   ◉

```python
[6]: compressor = [fcbench.codecs.BitRound(keepbits=5), fcbench.codecs.Zlib(level=9)]

stats = []
t_c = fcbench.compressor.compress_decompress(t, compressor, measurements=stats)

chart = earthkit.plots.quickmap.plot(t_c, units="degC")
chart.title("Compressed {variable_name} on {time:%d.%m.%Y at %H:%M}");
```

Compressed 2 metre temperature on 01.12.2012 at 14:00



2 metre temperature (°C)

Simple 🔘  0 ⬛ 1 ⚙ Python (Pyodide) | Idle

Mode: Command  🛡  Ln 1, Col 1  demo.ipynb  3 🔔

File   Edit   View   Run   Kernel   Tabs   Settings   Help

Code   Download   Python (Pyodide)

/ 02-19-2025-05-41-34-9e598c0b /
climet-lab-demo /

Filter files by name

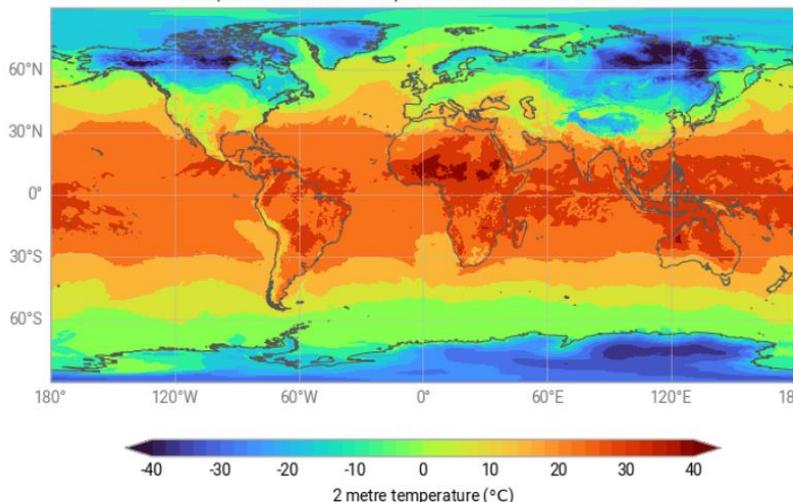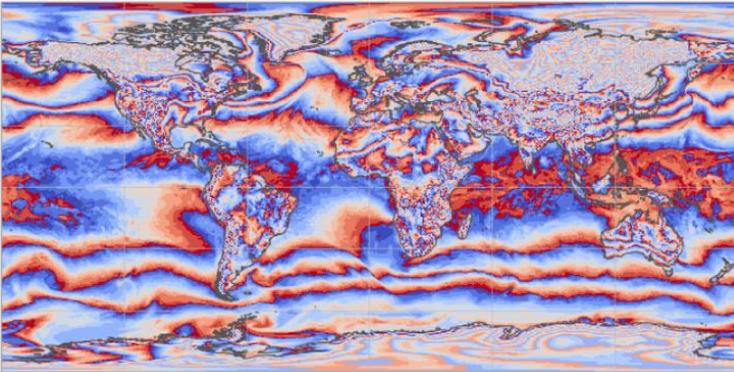| Name | Modified |
|---|---|
| • demo.ipynb | now |
| LICENSE.txt | 2s ago |
| t2m.nc | now |
| utils.py | 17m ago |

```python
[7]: utils.format_compress_stats(compressor, stats)
```

```
[7]:
```

| Codec | compression ratio [raw B / enc B] | encode instructions [#/B] | decode instructions [#/B] |
|---|---|---|---|
| BitRound(keepbits=5) | 1.00 | 0.3 | 0.1 |
| Zlib(level=9) | 74.92 | 27.5 | 1.7 |

```python
[8]: with xr.set_options(keep_attrs=True):
         chart = earthkit.plots.quickmap.plot(t_c - t, style=earthkit.plots.styles.Style(
             colors="coolwarm", levels=earthkit.plots.styles.levels.Levels(divergence_point=0.0),
         )); chart.title("Compression Error for {variable_name} on {time:%d.%m.%Y at %H:%M}");
```

Compression Error for 2 metre temperature on 01.12.2012 at 14:00



Simple   0   1   Python (Pyodide) | Idle        Mode: Command   Ln 1, Col 1   demo.ipynb   3

Why do we need but avoid to use Lossy Compression?

The Compression Laboratory

An Online Open CliMet Science Laboratory

Future Outlook

# Why build an Online Laboratory? – Try and Convince Yourself

Trust in lossy compression requires **convincing yourself**

Convincing yourself requires **trying things out yourself**, with your own data and analyses

Trying out code has a setup/installation time cost
Reduce this friction to **engage** more stakeholders

Why build an Online Laboratory? – Try and Convince Yourself

# Trying out compression ... the *easy* way

## Compression should be easy to try out

Let's reduce the barrier to experiment and convince yourself,
a serverless WebAssembly container within your browser:

[compression.lab.climet.eu](compression.lab.climet.eu)

# Online Laboratory for Data Compression in Climate Science and Meteorology

Welcome to the **Online Laboratory for Data Compression in Climate Science and Meteorology**!

If you are familiar with JupyterLab, you should feel right at home with the user interface of this lab. You can use the JupyterLab interface at /lab and a REPL interface at /repl.

In fact, this laboratory is built using JupyterLite, "a JupyterLab distribution that runs entirely in [your] browser" by leveraging WebAssembly. In other words, while you typically need to install JupyterLab on your own machine or connect to a server that executes your code, JupyterLite runs installation-free in your webbrowser and allows your code, data, and information to stay entirely on your machine. To run Python code within your browser, JupyterLite uses Pyodide, "a Python distribution for the browser [...] based on WebAssembly".

While Pyodide already supports an extensive list of scientific Python packages, which we have contributed to, this laboratory comes with additional packages that are commonly used in the weather and climate science community, including (but not limited to) `metpy`, `cfgrib`, `earthkit`, and `xeofs`.

## Getting Started

To get started, click the blue `+` button in the top left to open a new launcher and create a new Python notebook from there. After the Python kernel has initialised, you can execute Python code in the cells of the notebook.

> ✏️ **Tip**
>
> While many Python packages can be `import` ed directly, additional pure Python packages can also be loaded by executing the `%pip install <PACKAGE>` magic inside a cell, after which the package can be imported.

> ⓘ **Note**
>
> The online laboratory has only been tested in recent Firefox and Chrome browsers. Some features may not (yet) be supported in Safari browsers.

> 📢 **Attention**
>
> The online laboratory runs with the strict memory constraints of your web browser. It is therefore recommended to only open and execute one or two notebooks at a time. When a notebook is closed, the kernel will automatically shutdown to preserve resources.
>
> If the online lab runs out of memory, you can save your work, close the notebook, and try to restart it. If you are still running low on memory, you should first download a copy of your notebooks, then reload the notebook page, re-upload the notebook, and continue working on them.
>
> If you intend on executing memory intensive workloads, it is best to continue working on the notebooks locally instead. The online laboratory is primarily designed for initial exploration and for sharing codes in a reproducible environment.

> ⚠️ **Caution**

## Why build an Online Laboratory? – Ease of Use

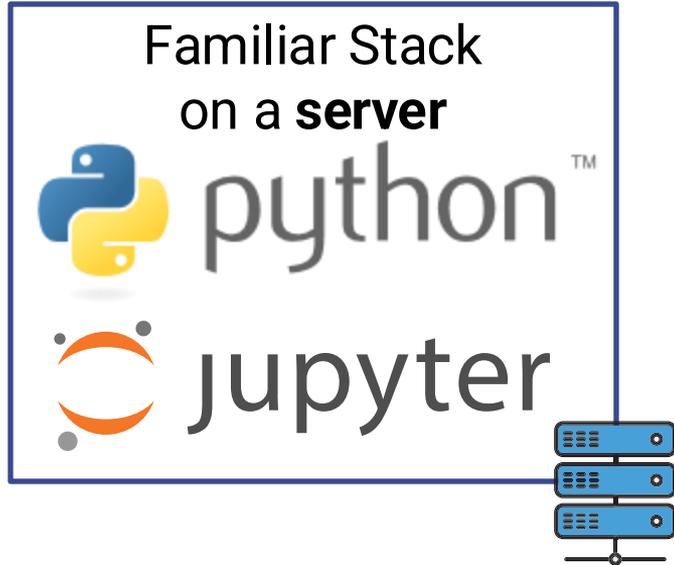**Easy and quick to pick up and go**

Just open a website, **no installation, try it out!**

Open as easily as documentation, but interactive
3x-5x slower Python is sufficient for examples
Motivate and excite first, only then ask for setup

# The technology behind the Online Laboratory

Familiar Stack
on a **server**

python™

jupyter

# The technology behind the Online Laboratory

**Familiar Stack
on a ~~server~~**



*Same* Code
*Same* Interface
*Same* Results

No installation
Easy sharing

**Serverless container**
in your browser



Benefit from the shared **ecosystem** of Python packages:

| **Pure Python** | **Scientific Ports** | **Earth Science Ports** |
|---|---|---|
| (`xarray`, `zarr`, `fsspec`, …) | (`numpy`, `scipy`, `sklearn`, …) | (`eccodes`, `netCDF4`, `cartopy`, …) |
| PyPi | Pyodide Community | Online Laboratory |

# Focus on reducing user friction and working with big data

**Remote streaming access to large datasets (GB, TB, …)**

| **Zarr** | **NetCDF4** | **GRIB** |
|---|---|---|
| natively chunked `fsspec` filesystems | separate metadata JIT with `kerchunk` | requires full scan prior with `gribscan` |

**Downloading smaller datasets:**
Access via `urllib`, `ecmwfapi`, `cdsapi`, `earthkit.data`, …
Load via the native `cfgrib` / `netCDF4` / `h5netcdf` / `zarr` driver

**Local datasets** of any size can be mounted and read in-place

Why build an Online Laboratory? – Community Collaboration

**Service to the greater community**

Locked Python environment with pre-installed packages

Your local notebooks and data never leave your device

Easy **sharing** of code samples and experiments:
e.g. lab.climet.eu/v0.2/github/:org/:repo/:branch/*

# Why build an Online Laboratory? – Interactive Documentation



**Static Code Example**

**Executable Code Example**

< 1 min

https://earthkit-plots.readthedocs.io/en/latest/examples/gallery/grid-types/reduced-gg-point-cloud.html

https://lab.climet.eu/v0.2/raw/github/ecmwf/earthkit-plots/main/docs/examples/gallery/grid-types/reduced-gg-point-cloud.ipynb

Why do we need but avoid to use Lossy Compression?

The Compression Laboratory

An Online Open CliMet Science Laboratory

Future Outlook

# Upcoming Online Laboratory Community Service Features

now
## Open existing repositories or notebooks

05/25
## Customize the locked Python package versions

2025
## Anyone can produce, share, and try out ready-to-go interactive documentation, analyses, or code examples

# The Roadmap to *Fearless* Lossy Compression

**2025**

Establish clear **Safety Requirements**
Benchmark the Safety and Performance of compressors

**Community** Outreach facilitated by the Online Laboratory
**Collaboration** with and **Service** for other projects

**2026**

**Safeguards** to safely use any lossy compressor
**Recommendations** for Operational Lossy Compression

# Thank you very much!

**Thank you also to everyone who has helped with this project:**
Clément Bouvier (UH), Hood Chatham (Pyodide @ MIT), Joona Cornér (UH), Alex Crichton (WASM @ Fermyon),
Milan Klöwer (Oxford), Daniel Köhler (UH), Juha Lento (CSC), Iain Russell (ECMWF),
Jeremy Tuloup (Jupyter @ QuantStack), and James Varndell (ECMWF)

/climet-eu/lab

## CASTIEL2 Code of the Month, 19.02.2025