



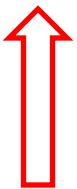
# Towards Safe and Fearless Lossy Compression of Weather and Climate Model Data

**Juniper Tyree<sup>[1,5]</sup>, Sara Faghih-Naini<sup>[2,5]</sup>, Milan Klöwer<sup>[3]</sup>, Tim Reichelt<sup>[3,6]</sup>,  
Peter Dueben<sup>[2,5]</sup>, Karsten Peters-von Gehlen<sup>[4,5]</sup>, and Heikki J. Järvinen<sup>[1,5]</sup>**

[1] INAR, University of Helsinki, [2] ECMWF, [3] University of Oxford, [4] DKRZ, [5] ESiWACE3, [6] Embed2Scale

CONTINENTS Webinar, 28.05.2025





**Feel free to take pictures of  
most slides ...**

**... unless the camera is red**



# Towards exascale weather and climate simulations

Supporting the community of weather and climate modelling in Europe



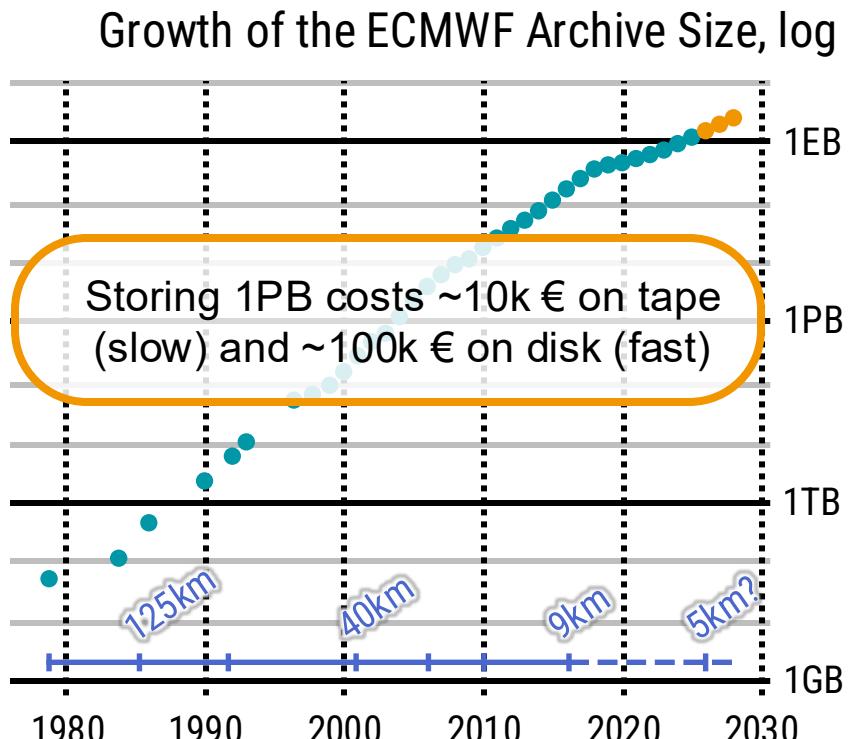
03  
Develop tools to tackle the data challenge of high-resolution models

3.2  
Domain-specific data compression

# **What is Lossy Compression, Why do we need it, and Why is it not used?**



## Problem: Growth in Weather & Climate Model Data Volumes



We need  
lossy  
compression  
  
But losing  
data accuracy  
is scary



## **Problem:** Fear of silently losing problem-specific information

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

**Domain scientist:** "Nice! ... Actually, did we 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<sub>2</sub>**  
...

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:** "Just check the compressed outputs to see if anything broke"



**Problem:** We need lossy compression but are afraid of using it

Lossy Compression is crucial to handle Data Growth

**Unclear safety requirements** and error tolerances for lossy compression in weather and climate science

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



# How do we bridge the gap between lossy compression and user needs to build trust?

# Our roadmap towards Safe, Trusted, and Fearless Lossy Compression



## The Roadmap to *Fearless Lossy Compression*



### **Open Science Online Laboratory to Try it out and Convince Yourself**

*We don't even know what error tolerance we have!*

**Specify clear safety requirements**

*What compressor should I use?*

**Benchmark compressors for safety and performance**

*But my safety requirements are far stricter than yours!*

**Create safeguards to use any compressor without fear**



# Compression should be *Easy & Fearless*

**“We don’t want to disrupt the science, but  
disrupt the data challenge”**

Heikki Järvinen, ESiWACE3 kick-off, BOG5 session, 01.02.2023



# What are our safety requirements?

# What errors can we tolerate?

**Juniper Tyree<sup>[1,4]</sup>, Milan Klöwer<sup>[2]</sup>, Sara Faghih-Naini<sup>[3,4]</sup>, and Peter Dueben<sup>[3,4]</sup>**

[1] INAR, University of Helsinki, [2] University of Oxford, [3] ECMWF, [4] ESIWACE3



**Problem:** Unclear which compression errors can be tolerated

**Model outputs are stored with excessive precision**

Trailing mantissa bits contain false information / noise

**Tolerance varies by variable, resolution, use case**

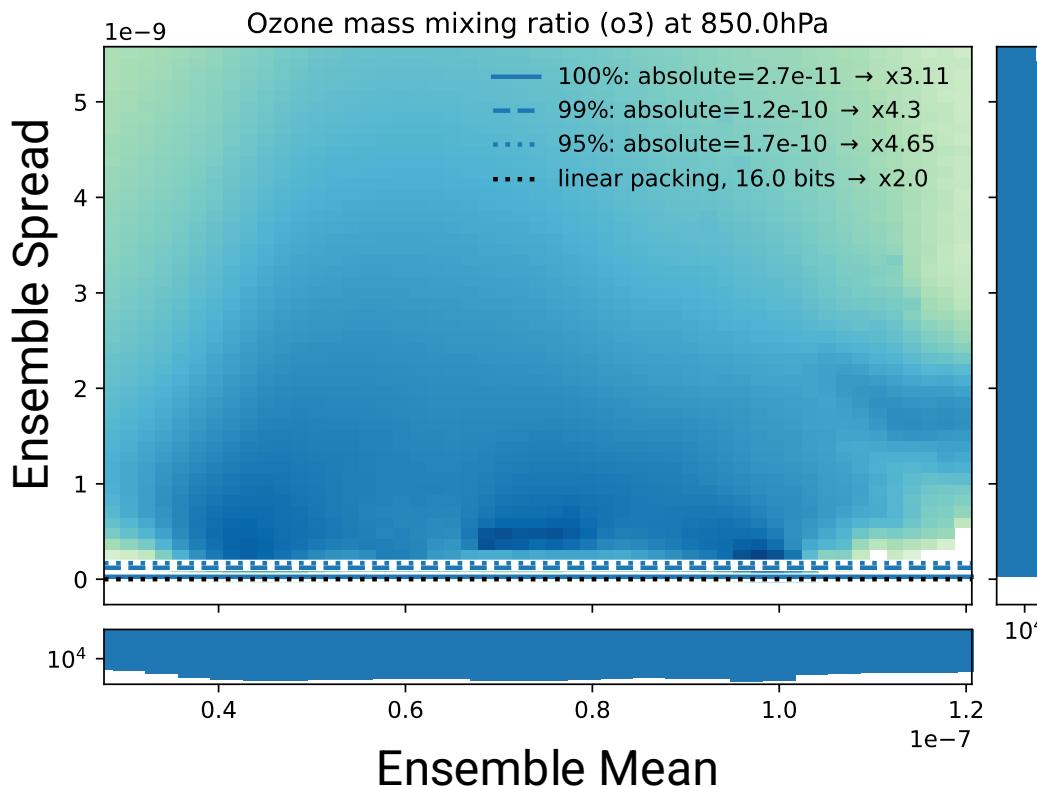
Model output archival must use conservative tolerances

**Community needs recommendations and standards**

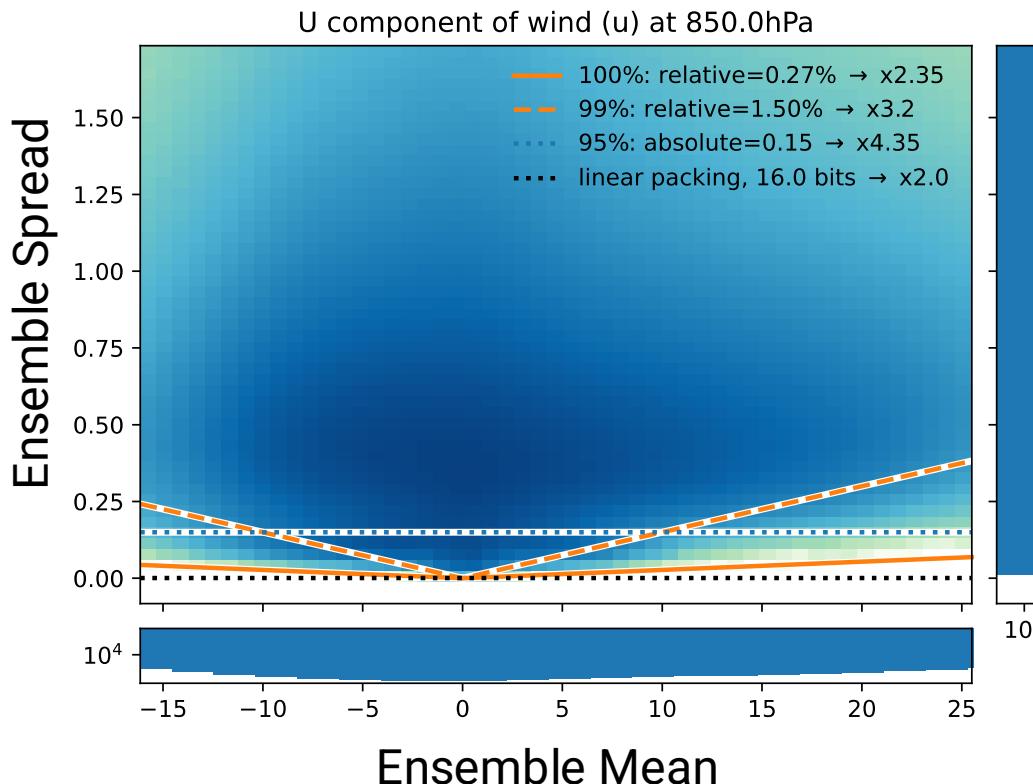
GRIB already uses lossy linear packing, improve on that



## Results: Variables with absolute error bounds, e.g. O3, temp

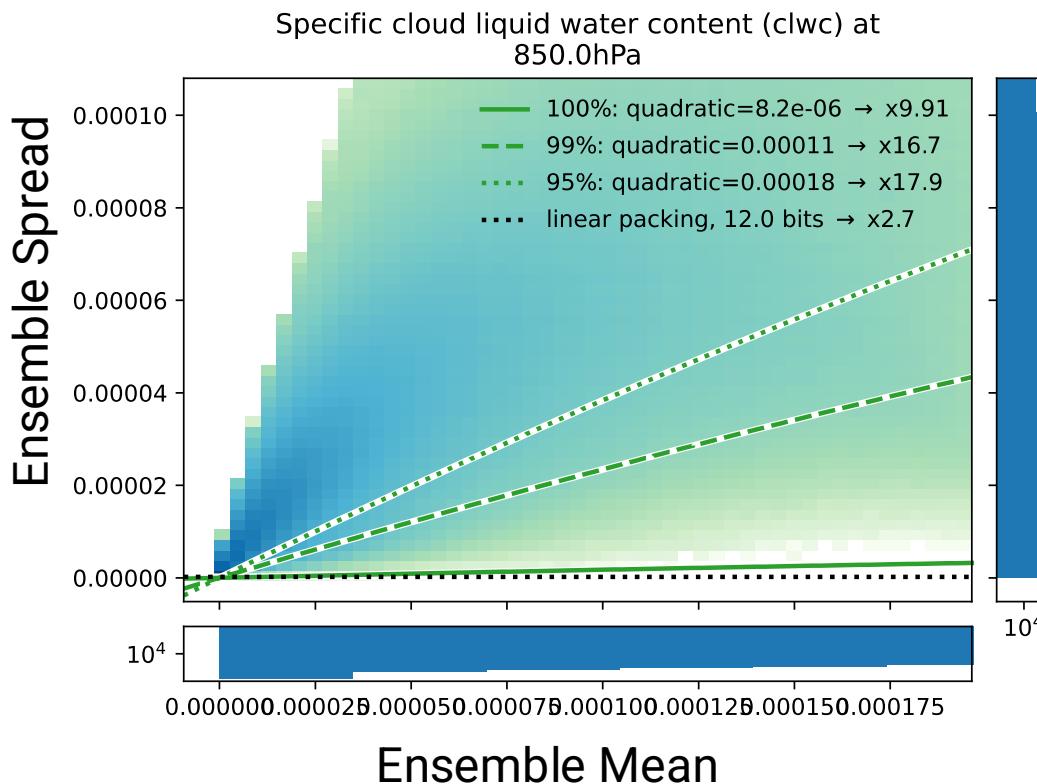


## Results: Variables with relative error bounds, e.g. u/v wind



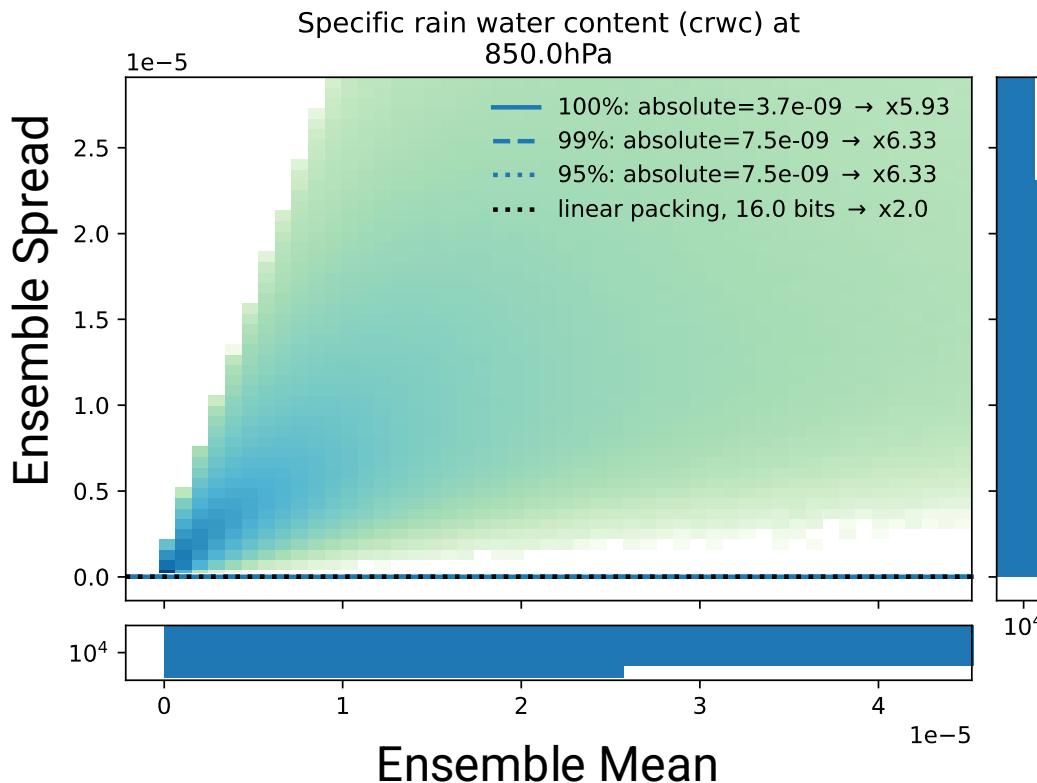


## Results: Variables with relative error bounds, e.g. cloud water

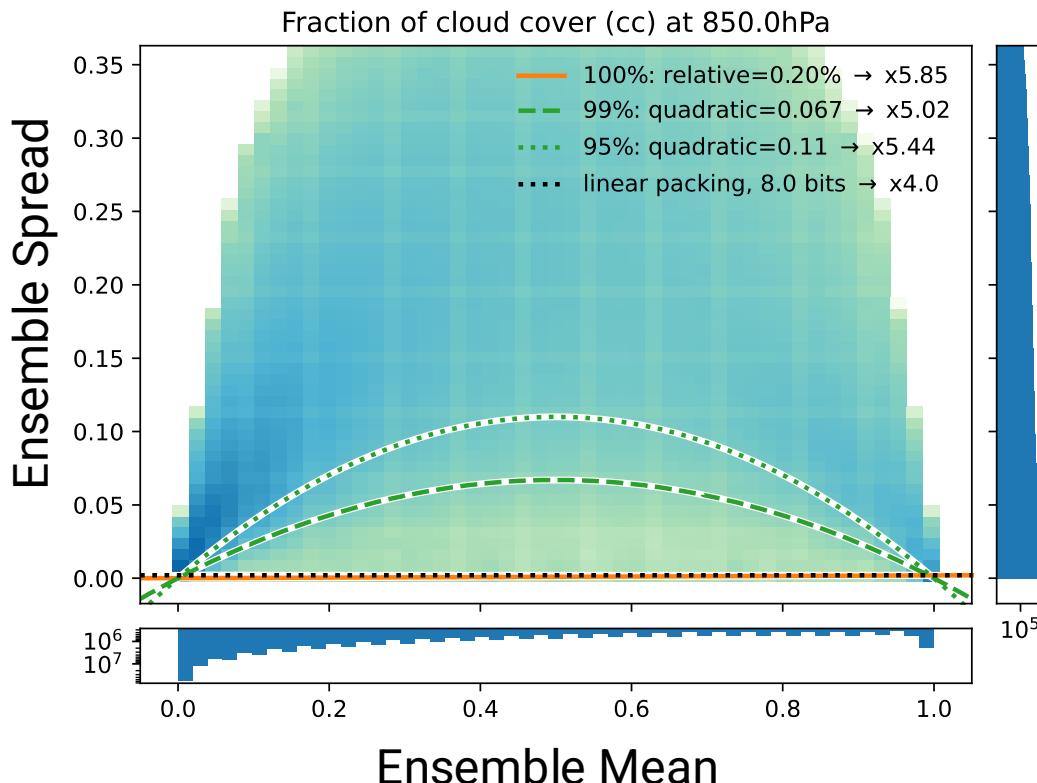




## Results: Variables with unclear error bounds, e.g. rainwater



## Results: Variables with quadratic error bounds, e.g. cloud cover





## Outlook: Next steps, publication, and future work

-Q2

### **Comparing expert and automatic bounds**

Experience vs data analysis

Q2-Q3

### **Initial publication and table of error bounds**

Three suggested levels of error bounds per variable

...

### **Outreach, Feedback, Refinement, Standards**

Further feedback and analyses, replace linear packing



# How do existing compressors perform? Are they safe to use?



# *ClimateBenchPress:*

# A Benchmark for Compression of Climate Data

Tim Reichelt<sup>[1,7]</sup>, Juniper Tyree<sup>[2,8]</sup>, Milan Klöwer<sup>[1]</sup>, Peter Dueben<sup>[3,8]</sup>, Bryan Lawrence<sup>[4]</sup>,  
Dorit Hammerling<sup>[5]</sup>, Allison Baker<sup>[6]</sup>, Sara Faghih-Naini<sup>[3,8]</sup>, and Philip Stier<sup>[1,7]</sup>

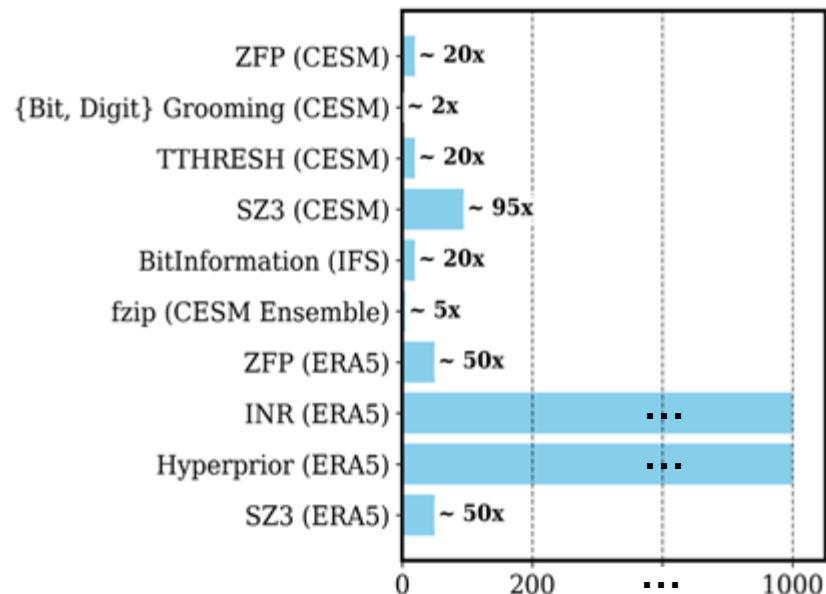
[1] University of Oxford, [2] INAR, University of Helsinki, [3] ECMWF, [4] NCAS, University of Reading,

[5] Colorado School of Mines, [6] NCAR, [7] Embed2Scale, [8] ESIWACE3



## Problem: Why do we need a benchmark?

Compression ratios reported in  
the literature ...



are **difficult to compare** because  
they use different

- Datasets and **resolutions**
- **Variables** and dimensions
- Compression **error bounds**
- Measurement techniques



## Aims: What do we want from a good benchmark?

### **Fair and Comparable Results**

Same data, same error bounds (3 levels), same tests

### **Interpretable Results from Representative Sources**

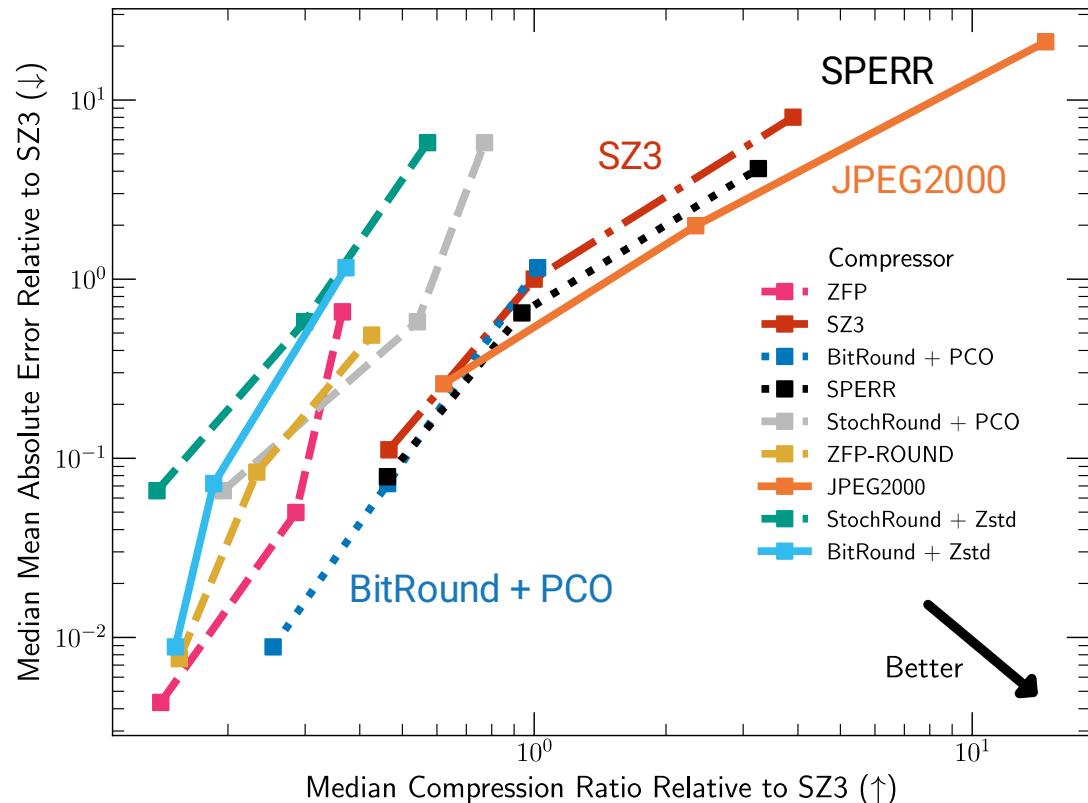
2D & 3D from CMIP6, ERA5, NextGEMS, CAMS, Satellite

### **Ease of use and expansion**

Modular and extensible, several dataset sizes



## Results: Preliminary Insights, averaged across several variables



- **SZ3, SPERR, and JPEG2000:** highest compression
  - JPEG2000 only has global error bounds (not pointwise)
- **Bit Rounding + PCO** is competitive (but is quite conservative)



## Outlook: Next steps, publication, and future work

Q1-Q2

### **Preparing the small version of the benchmark**

Initial subset of data, compressors, and metrics

Q3-Q4

### **Initial publication and recommendations**

Per-variable compressor scoreboard with trade-offs

2026-

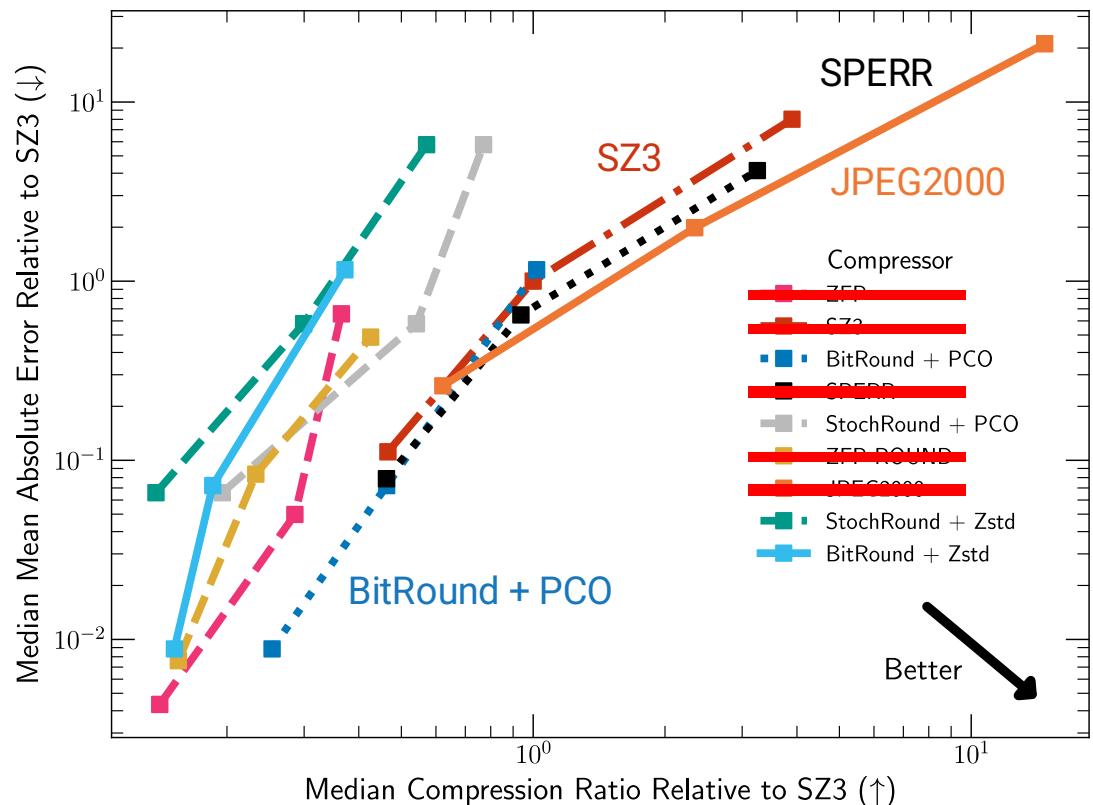
### **Expanded benchmark, renew recommendations**

Larger datasets, more compressors (ML), more metrics

**Spooky stories  
of what can go wrong  
with lossy compression**



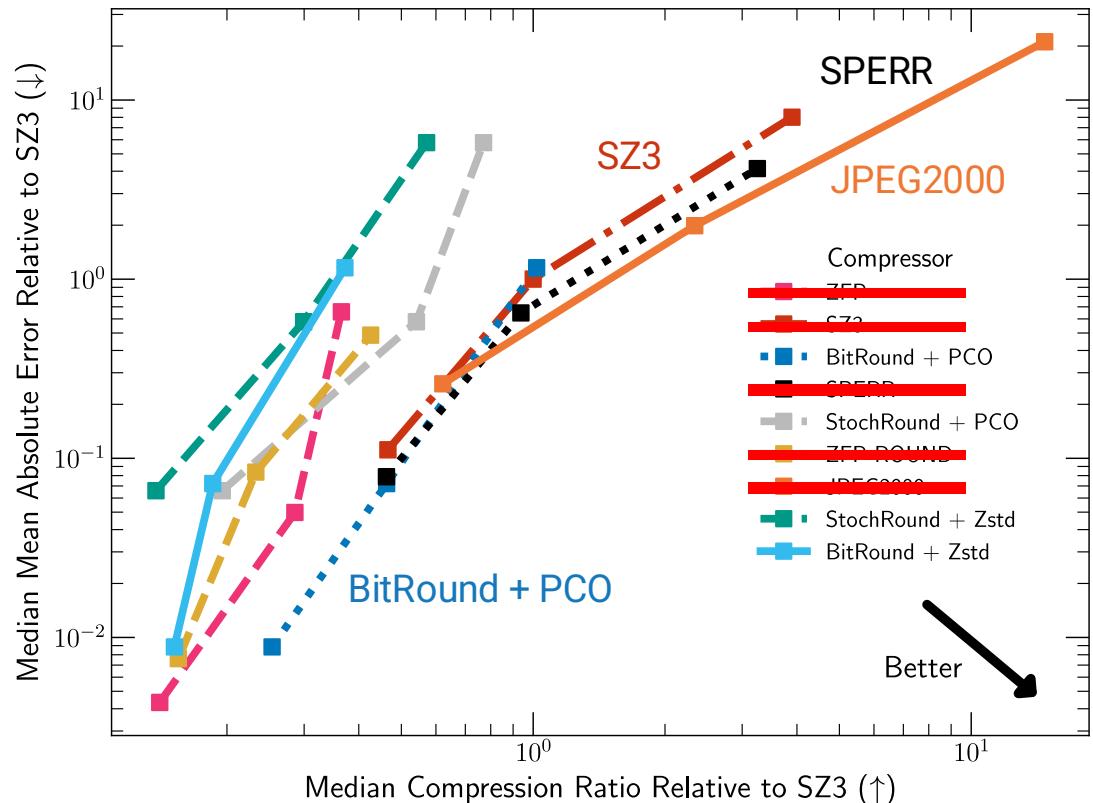
## Problem: Many popular compressors don't handle edge cases



- **ZFP**: error bias, no NaNs, no Inf's
- **ZFP-ROUND**: no NaNs, no Inf's
- **SZ3**: no NaNs (v3.2)
- **SPERR**: no NaNs or big values, loose bounds
- **JPEG2000**: no NaNs, no Inf's, no pointwise error bounds



## Problem: Many popular compressors don't handle edge cases

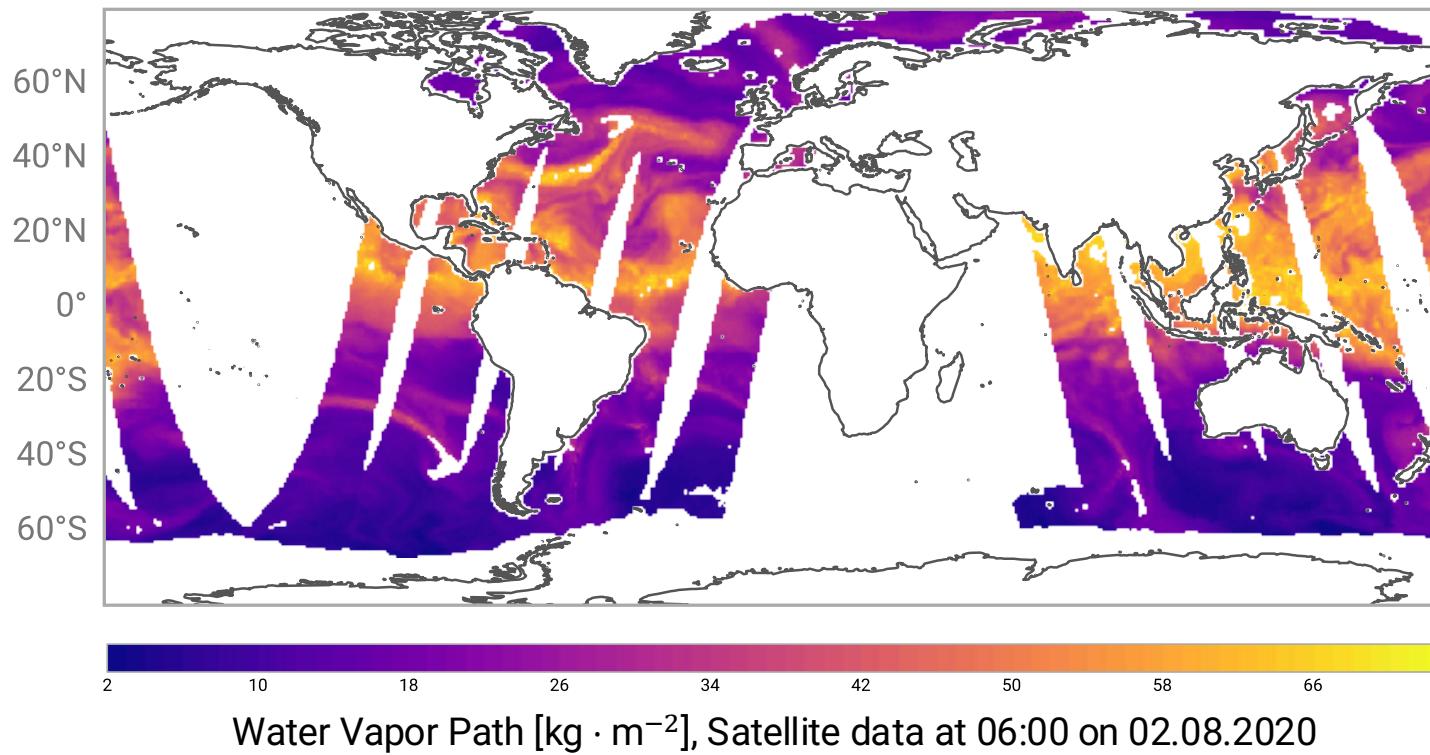


- **ZFP:** error bias, no NaNs, no Inf's
  - **ZFP-ROUND:** no NaNs, no Inf's
  - **SZ3:** no NaNs (v3.2)
  - **SPERR:** no NaNs or big values, loose bounds
  - **JPEG2000:** no NaNs, no Inf's, no pointwise error bounds
- Explicit failure to compress**



## Problem: Hallucination in non-ML lossy compressors

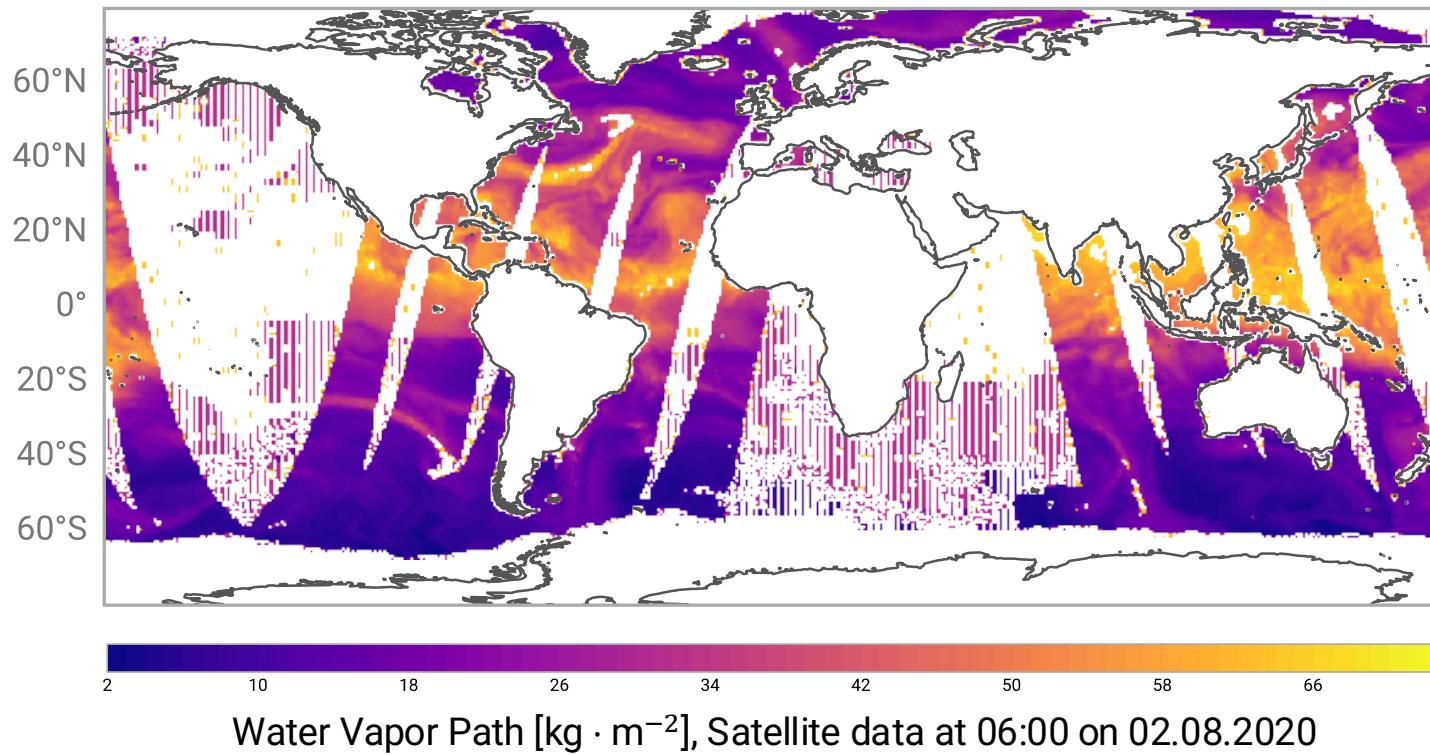
Original data, NaN values on land and from missing coverage





## Problem: Hallucination in non-ML lossy compressors

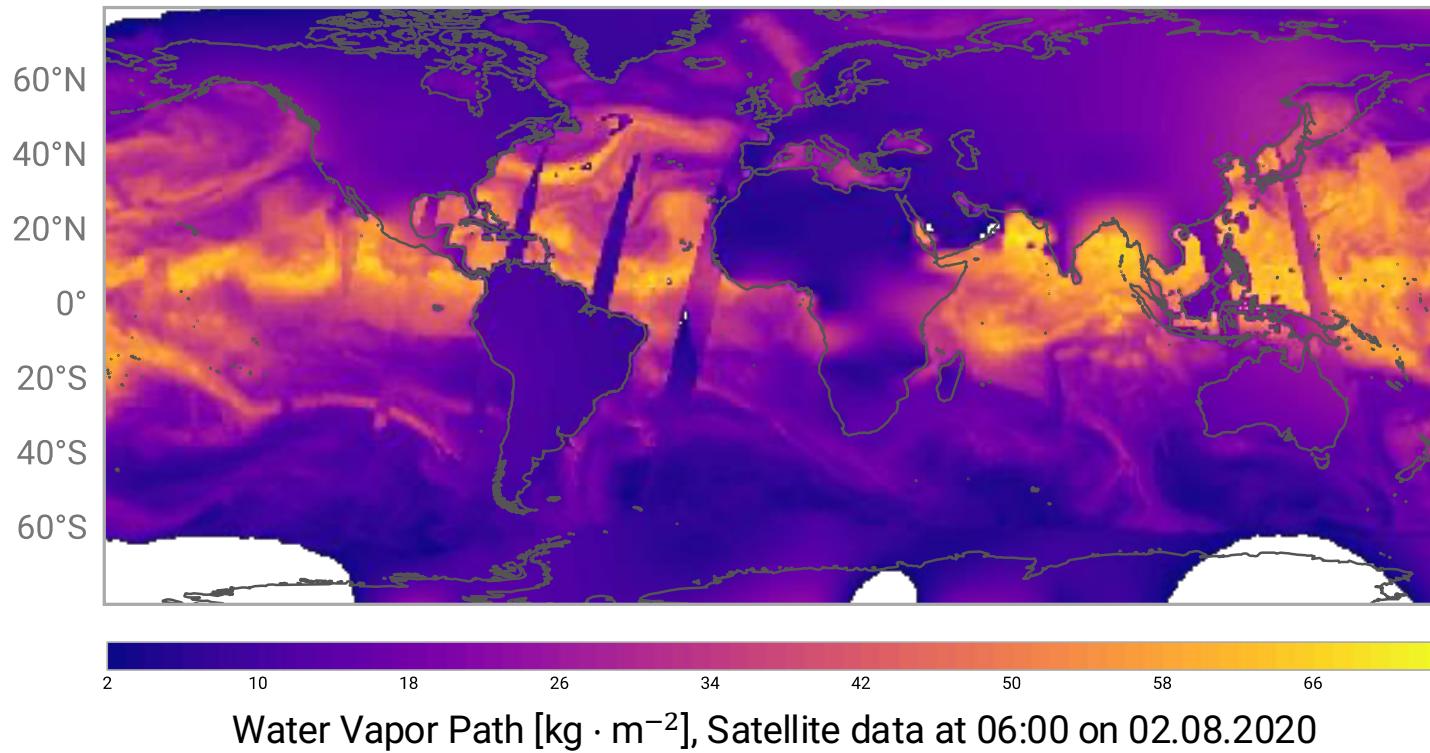
ZFP, absolute error bound  $\epsilon_{\text{abs}} = 0.1 \text{ kg} \cdot \text{m}^{-2}$





## Problem: Hallucination in non-ML lossy compressors

SZ3, absolute error bound  $\epsilon_{\text{abs}} = 0.1 \text{ kg} \cdot \text{m}^{-2}$





# Compression Safeguards: Fearless Lossy Compression with declarative Safety

Juniper Tyree<sup>[1,4]</sup>, Milan Klöwer<sup>[2]</sup>, Clément Bouvier<sup>[1,4]</sup>, Daniel Köhler<sup>[1]</sup>, Robert Underwood<sup>[3]</sup>,  
Heikki J. Järvinen<sup>[1,4]</sup>, et al.\*

[1] INAR, University of Helsinki, [2] University of Oxford, [3] Argonne National Laboratory, [4] ESiWACE3, \* pending authors



## How can the compression safeguards help You?



**Declare your safety requirements, the safeguards guarantee them**

**Preserve special values (NaN, Inf, 0, missing value, etc.)**

**Regional pointwise error bounds (abs, rel, noa, ratio)**

**Preserve properties, e.g. signs & monotonic sequences**

**Error bounds over quantities of interest ( $x^2$ ,  $\partial X$ , ...)**

**Logical combinations** over all of the above



## Outlook: Next steps, publication, and future work

-Q3

### **Compression Safeguards Software**

Initial integration with numcodecs, possibly Zarr

Q4-

### **Initial publication and evaluation**

Examples of preserving complex safety requirements

2026-

### **Expansion and Operationalization**

Further safeguards, CDF integration, perf optimizations

**Convince Yourself  
through Outreach and an  
Open Science Laboratory**

Trying something out for yourself  
is especially important if you  
have **doubts**

To **convince yourself** that lossy compression  
is safe, you need to **try it out yourself**

Trying something out should be easy  
to **engage earlier** and **delay setup costs**  
after **convincing** users

**Solution for the entire community**

i.e. not just for our project (lossy compression) but also for your project, documentation, demos, ...



## The Online Laboratory for Climate Science and Meteorology

What's easier than opening a URL in your browser?

**Serverless In-Browser**  
Interactive Computing



**PYODIDE**



Code and data stay *local*



No setup, no installation, <1min to start

Reproducible and version-locked  
Ensure your examples keep working

Supports many compiled scientific + especially Earth  
Science Python packages

Extra support for accessing large datasets

**Ease of use: same code, same results**

The screenshot shows a web browser window with the URL [earthkit-plots.readthedocs.io](https://earthkit-plots.readthedocs.io). The page title is "Welcome to the earthkit-plots documentation". A large orange rounded rectangle highlights the main heading "Example: Software documentation is non-interactive". The left sidebar contains navigation links for "EXAMPLES", "DOCUMENTATION", "INSTALLATION", and "PROJECTS". The main content area includes a status bar with "ESEE Foundation Maturity Incubating License Apache 2.0" and a "stable" dropdown. The footer contains copyright information for ECMWF.

# Example: Software documentation is non-interactive

## Welcome to the earthkit-plots documentation

ESEE Foundation Maturity Incubating License Apache 2.0

earthkit-plots leverages the power of the **earthkit** ecosystem to make producing publication-quality scientific graphics as simple and convenient as possible.

It is built on top of popular data science and visualisation tools like **numpy**, **xarray**, **matplotlib** and **cartopy**, but provides a very high-level API enriched with domain-specific knowledge, making it exceptionally easy to use.

Key features include:

- ⚡ Concise, high-level API** Generate high-quality visualisations with minimal code.
- 🧠 Intelligent formatting** Titles and labels automatically adapt based on common metadata standards.
- 🎨 Customisable style libraries** Easily swap styles to match your organisation, project, or personal preferences.
- 🔍 Automatic data styling** Detects metadata like variables and units to optionally apply appropriate formatting and styling.
- 🌐 Complex grids supported out-of-the-box** Visualise grids like HEALPix and reduced gaussian without any extra legwork.

earthkit-plots

stable

Search docs

EXAMPLES

Examples

Gallery

- Domains and projections
- Gridded data
- Grid types
- Ancillary data

DOCUMENTATION

User guide

API Reference

Development

INSTALLATION

Installation

License

PRODUCTS

earthkit plots

37b / 44

# Example: Software documentation is non-interactive

## Gallery

### Domains and projections

Built-in named domains

Using a cartopy CRS

Custom domains

Combining domains

### Gridded data

El Niño

Hatched shading

Model orography

Temperature and pressure

<https://earthkit-plots.readthedocs.io>

Copyright 2022- European Centre for Medium-Range Weather Forecasts (ECMWF). Licensed under the Apache License, Version 2.0.

earthkit-plots

stable

Search docs

EXAMPLES

Examples

Gallery

- Domains and projections
- Gridded data
  - El Niño
  - Hatched shading
  - Model orography
  - Temperature and pressure
  - Time zones
  - Unstructured grids
  - Categorical data
  - EFAS (Lambert Azimuthal Equal Area)
  - RGB composite
- Grid types
- Ancillary data

chart.ocean(color="#444")  
chart.coastlines(color="white")  
chart.gridlines()  
chart.legend(label="model orography (\$m\$)")  
chart.show()

37c / 44

lab.climet.eu

# Demo: Interactive Documentation in the Online Laboratory

File Edit View Run Kernel Code Download Python (Pyodide)

[5]: # Compress the geopotential ( $m^2/s^2$ ) using ZFP

```
z = CodecStack(  
    Zfp(mode="fixed-accuracy", tolerance=100) ← strict error bound  
).encode_decode_data_array(data[0].to_xarray().z)
```

chart = earthkit.plots.Map(domain=[-15, 35, 32, 72]) •••

Simple  Python (Pyodide) | Idle

Mode: Edit  Ln 1, Col 1 model-orography.ipynb 3

<https://lab.climet.eu/v0.3.0/raw/github-tag/ecmwf/earthkit-plots/0.3.1/docs/examples/gallery/gridded-data/model-orography.ipynb>



lab.climet.eu 38b / 44

# Demo: Interactive Documentation in the Online Laboratory

File Edit View Run Kernel Code Download Python (Pyodide)

[7]: # Compress the geopotential ( $m^2/s^2$ ) using ZFP

```
z = CodecStack(  
    Zfp(mode="fixed-accuracy", tolerance=10000) ← loose error bound  
).encode_decode_data_array(data[0].to_xarray().z)
```

chart = earthkit.plots.Map(domain=[-15, 35, 32, 72]) •••

Simple Python (Pyodide) | Idle Mode: Edit Ln 1, Col 1 model-orography.ipynb 4



# Demo: ClimateBenchPress

To summarize the performance of the two compressors across all the error bounds we can make a rate-distortion plot. Interestingly, for this dataset BitRounding seems to perform better than SZ3 because BitRounding produces higher compression ratios for a similar error level.

```
[15]: plot_variable_rd_curve(  
    results,  
    distortion_metric="MAE",  
    outfile=plotspath / "rd_curve.png"  
)  
Image((plotspath / "rd_curve.png").read_bytes(), width=512)
```

[15]:

Compression Ratio [raw B / enc B]	MAE (BitRound + PCO)	MAE (SZ3)
$5 \times 10^0$	$10^{-3}$	$10^{-2.5}$
$10^1$	$10^{-2}$	$10^{-1.5}$
$2 \times 10^1$	$10^{-1.5}$	$10^{-1}$

Mode: Command    Ln 1, Col 1    climatebenchpress.ipynb    2



lab.climet.eu 40 / 44

## Demo: Online Compression Laboratory

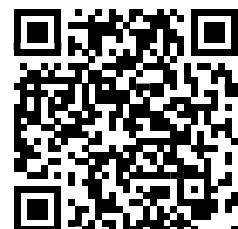
```
[16]: with xr.set_options(keep_attrs=True):
    da_zfp_error = (da_zfp - da).compute()

plot_data(
    da_zfp_error, title_prefix="Compression Error for ",
    title_postfix=f"\n{n{zfp_compressor}}", error=True,
)
```

Compression Error for 2 metre temperature on 01.12.2012 at 14:00  
CodecStack(Asinh(linear\_width=1.0), Zfp(mode='fixed-accuracy', tolerance=0.01))

2 metre temperature (K)

Simple  Python (Pyodide) | Idle Mode: Edit  Ln 1, Col 1 01-intro.ipynb 2



How can **You** benefit from the Online Laboratory today?

[docs.climet.eu](https://docs.climet.eu)



**Interactive and Reproducible Open Science Laboratory**

**Share your existing notebooks from repos / gists / URLs**

**Customize the provided packages and versions**

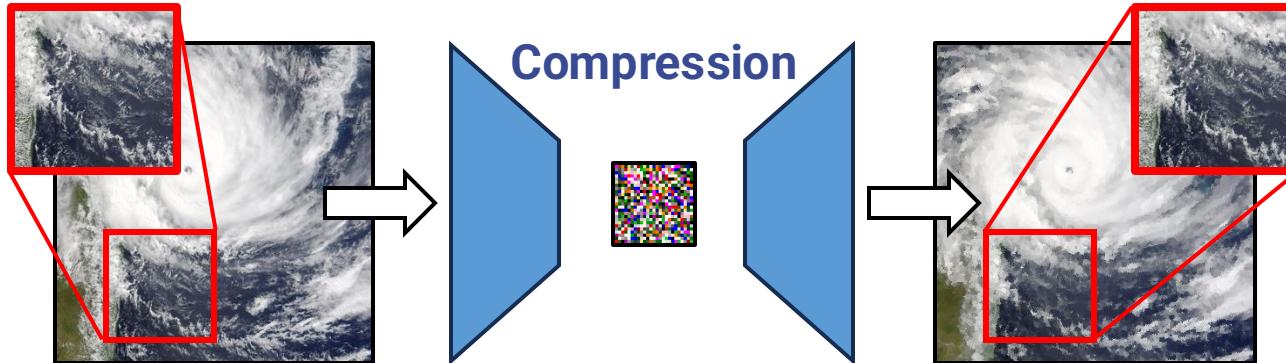
**Ongoing development and support and new packages**

**Explore lossy data compression on climate model data**



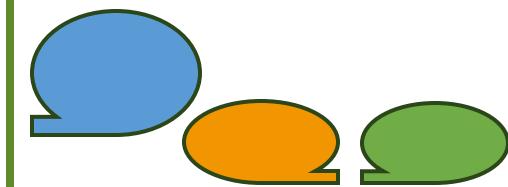
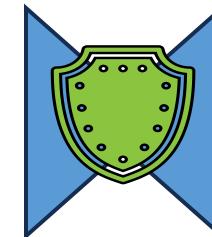
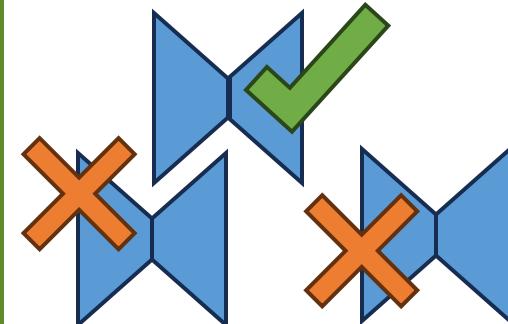
# In summary ...

# Weather & Climate Prediction will soon produce too much data



*Paralyzing fear of compression artifacts or loss of details*

## Safe and Fearless Lossy Data Compression in Climate Science



**Juniper Tyree**  
juniper.tyree@helsinki.fi



# Towards Safe and Fearless Lossy Compression of Weather and Climate Model Data

**Thank You for Your attention!**

 ID 0000-0002-7923-9609

**Juniper Tyree**  
[juniper.tyree@helsinki.fi](mailto:juniper.tyree@helsinki.fi)

 juntyr

Disclaimer: Funded by the European Union. Views and opinions expressed are however those of the author(s) only and do not necessarily reflect those of the European Union or the European High Performance Computing Joint Undertaking (JU). Neither the European Union nor the granting authority can be held responsible for them.

**CONTINENTS Webinar, 28.05.2025**



**EuroHPC**  
Joint Undertaking

*Funded by the European Union. This work has received funding from the European High Performance Computing Joint Undertaking (JU) under grant agreement No 101093054.*



**Co-funded by  
the European Union**

