

---

# **Idcpy**

***Release 0.13***

**Alex Pinard, Allison Baker, Anderson Banihirwe, Dorit Hammerlin**

**Mar 16, 2021**



**CONTENTS:**

<b>1</b>	<b>Installation using Conda (recommended)</b>	<b>3</b>
<b>2</b>	<b>Alternative Installation</b>	<b>5</b>
<b>3</b>	<b>Accessing the tutorial</b>	<b>7</b>
<b>4</b>	<b>Re-create notebooks with Pangeo Binder</b>	<b>9</b>
<b>5</b>	<b>Documentation</b>	<b>11</b>
5.1	Installation . . . . .	11
5.2	Development . . . . .	12
5.3	API Reference . . . . .	13
5.4	Examples . . . . .	21
<b>6</b>	<b>Indices and tables</b>	<b>109</b>
	<b>Index</b>	<b>111</b>





ldcpy is a utility for gathering and plotting metrics from NetCDF or Zarr files using the Pangeo stack. It also contains a number of statistical and visual tools for gathering metrics and comparing Earth System Model data files.

**AUTHORS** Alex Pinard, Allison Baker, Anderson Banihirwe, Dorit Hammerling

**COPYRIGHT** 2020 University Corporation for Atmospheric Research

**LICENSE** Apache 2.0

Documentation and usage examples are available [here](#).



## INSTALLATION USING CONDA (RECOMMENDED)

Ensure conda is up to date and create a clean Python (3.6+) environment:

```
conda update conda
conda create --name ldcpy python=3.8
conda activate ldcpy
```

Now install ldcpy:

```
conda install -c conda-forge ldcpy
```



## ALTERNATIVE INSTALLATION

Ensure pip is up to date, and your version of python is at least 3.6:

```
pip install --upgrade pip  
python --version
```

Install cartopy using the instructions provided at <https://scitools.org.uk/cartopy/docs/latest/installing.html>.

Then install ldcpy:

```
pip install ldcpy
```



## ACCESSING THE TUTORIAL

If you want access to the tutorial notebook, clone the repository (this will create a local repository in the current directory):

```
git clone https://github.com/NCAR/ldcpy.git
```

Start by enabling Hinterland for code completion and code hinting in Jupyter Notebook and then opening the tutorial notebook:

```
jupyter nbextension enable hinterland/hinterland  
jupyter notebook
```

The tutorial notebook can be found in docs/source/notebooks/TutorialNotebook.ipynb, feel free to gather your own metrics or create your own plots in this notebook!

Other example notebooks that use the sample data in this repository include PopData.ipynb and MetricsNotebook.ipynb.

The AWSDataNotebook grabs data from AWS, so can be run on a laptop with the caveat that the files are large.

The following notebooks assume that you are using NCAR's JupyterHub (<https://jupyterhub.ucar.edu>): Large-DataGladenotebook.ipynb, CompressionSamples.ipynb, and error\_bias.ipynb





## **RE-CREATE NOTEBOOKS WITH PANGEO BINDER**

Try the notebooks hosted in this repo on Pangeo Binder. Note that the session is ephemeral. Your home directory will not persist, so remember to download your notebooks if you make changes that you need to use at a later time!

Note: All example notebooks are in docs/source/notebooks (the easiest ones to use in binder first are TutorialNotebook.ipynb and PopData.ipynb)



## 5.1 Installation

### 5.1.1 Installation using Conda (recommended)

Ensure conda is up to date and create a clean Python (3.6+) environment:

```
conda update conda
conda create --name ldcpy python=3.8
conda activate ldcpy
```

Now install ldcpy:

```
conda install -c conda-forge ldcpy
```

### 5.1.2 Alternative Installation

Ensure pip is up to date, and your version of python is at least 3.6:

```
pip install --upgrade pip
python --version
```

Install cartopy using the instructions provided at <https://scitools.org.uk/cartopy/docs/latest/installing.html>.

Then install ldcpy:

```
pip install ldcpy
```

### 5.1.3 Accessing the tutorial (for users)

If you want access to the tutorial notebook, clone the repository (this will create a local repository in the current directory):

```
git clone https://github.com/NCAR/ldcpy.git
```

Start by activating the ldcpy environment, enabling Hinterland for code completion in Jupyter Notebook and then starting the notebook server:

```
conda activate ldcpy
conda install -c conda-forge jupyter_contrib_nbextensions
jupyter contrib nbextension install --user
jupyter nbextension enable hinterland/hinterland
jupyter notebook
```

The tutorial notebook can be found in docs/source/notebooks/TutorialNotebook.ipynb, feel free to gather your own metrics or create your own plots in this notebook!

## 5.2 Development

### 5.2.1 Installation for Developers

First, clone the repository and cd into the root of the repository:

```
git clone https://github.com/NCAR/ldcpy.git
cd ldcpy
```

For a development install, do the following in the ldcpy repository directory:

```
conda env update -f environment.yml
conda activate ldcpy
python -m pip install -e .
```

Then install the pre-commit script and git hooks for code style checking:

```
pre-commit install
```

This code block enables optional extensions for code completion, code hinting and minimizing tracebacks in Jupyter. Then start the jupyter notebook server in your browser (at localhost:8888):

```
jupyter nbextension enable hinterland/hinterland
jupyter nbextension enable skip-traceback/main

conda activate ldcpy
jupyter notebook
```

### 5.2.2 Instructions and Tips for Contributing

For viewing changes to documentation in the repo, do the following:

```
pip install -r docs/requirements.txt
sphinx-reload docs/
```

This starts and opens a local version of the documentation in your browser (at localhost:5500/index.html) and keeps it up to date with any changes made. Note that changes to docstrings in the code will not trigger an update, only changes to the .rst files in the docs/ folder.

If you have added a feature or fixed a bug, add new tests to the appropriate file in the tests/ directory to test that the feature is working or that the bug is fixed. Before committing changes to the code, run all tests from the project root directory to ensure they are passing.

```
pytest -n 4
```

Additionally, rerun the TutorialNotebook in Jupyter (Kernel -> Restart & Run All). Check that no unexpected behavior is encountered in these plots.

Now you are ready to commit your code. pre-commit should automatically run black, flake8, and isort to enforce style guidelines. If changes are made, the first commit will fail and you will need to stage the changes that have been made before committing again. If, for some reason, pre-commit fails to make changes to your files, you should be able to run the following to clean the files manually:

```
black --skip-string-normalization --line-length=100 .
flake8 .
isort .
```

### 5.2.3 Adding new package dependencies to ldcpy

- 1) Adding new package dependencies requires updating the code in the following four places:

/ci/environment.yml /ci/environment-dev.yml /ci/upstream-dev-environment.yml /requirements.txt

If the package dependency is specifically used for documentation, instead of adding it to /requirements.txt, add it to:

/docs/source/requirements.txt

If this package is only used for documentation, skip the remaining steps.

- 2) If the package is one that includes C code (such as numpy or scipy), update the autodoc\_mock\_imports list in /docs/source/conf.py. The latest build of the documentation can be found at (<https://readthedocs.org/projects/ldcpy/builds/>), if the build fails and the error message indicates a problem with the newest package - try adding it to autodoc\_mock\_imports.
- 3) Finally, update the ldcpy-feedstock repository (git clone <https://github.com/conda-forge/ldcpy-feedstock.git>), or manually create a branch and add the dependency in the browser. Name the branch add-<new\_dependency\_name>. In the file /recipe/meta.yaml, in the “requirements” section, under “run”, add your dependency to the list.
- 4) If the CI build encounters errors after adding a dependency, check the status of the CI workflow at (<https://github.com/NCAR/ldcpy/actions?query=workflow%3ACI>) to determine if the error is related to the new package.

## 5.3 API Reference

This page provides an auto-generated summary of ldcpy’s API. For more details and examples, refer to the relevant chapters in the main part of the documentation.

### 5.3.1 ldcpy Util (ldcpy.util)

```
ldcpy.util.check_metrics(ds, varname, set1, set2, ks_tol=0.05, pcc_tol=0.99999, spre_tol=5.0,
                        ssim_tol=0.99995, **metrics_kwargs)
```

Check the K-S, Pearson Correlation, and Spatial Relative Error metrics

#### Parameters

- **ds** (*xarray.Dataset*) – An xarray dataset containing multiple netCDF files concatenated across a ‘collection’ dimension
- **varname** (*str*) – The variable of interest in the dataset

- **set1** (*str*) – The collection label of the “control” data
- **set2** (*str*) – The collection label of the (1st) data to compare
- **ks\_tol** (*float*, *optional*) – The p-value threshold (significance level) for the K-S test (default = .05)
- **pcc\_tol** (*float*, *optional*) – The default Pearson correlation coefficient (default = .99999)
- **spre\_tol** (*float*, *optional*) – The percentage threshold for failing grid points in the spatial relative error test (default = 5.0).
- **ssim\_tol** (*float*, *optional*) – The threshold for the ssim test (default = .999950)
- **\*\*metrics\_kwargs** – Additional keyword arguments passed through to the DatasetMetrics instance.

**Returns out**

**Return type** Number of failing metrics

## Notes

Check the K-S, Pearson Correlation, and Spatial Relative Error metrics from:

A. H. Baker, H. Xu, D. M. Hammerling, S. Li, and J. Clyne, “Toward a Multi-method Approach: Lossy Data Compression for Climate Simulation Data”, in J.M. Kunkel et al. (Eds.): ISC High Performance Workshops 2017, Lecture Notes in Computer Science 10524, pp. 30–42, 2017 (doi:10.1007/978-3-319-67630-2\_3).

Check the SSIM metric from:

A.H. Baker, D.M. Hammerling, and T.L. Turton. “Evaluating image quality measures to assess the impact of lossy data compression applied to climate simulation data”, Computer Graphics Forum 38(3), June 2019, pp. 517-528 (doi:10.1111/cgf.13707).

K-S: fail if p-value < .05 (significance level) Pearson correlation coefficient: fail if coefficient < .99999 Spatial relative error: fail if > 5% of grid points fail relative error SSIM: fail if SSIM < .99995

`ldcpy.util.collect_datasets` (*varnames*, *list\_of\_ds*, *labels*, *\*\*kwargs*)

Concatenate several different xarray datasets across a new “collection” dimension, which can be accessed with the specified labels. Stores them in an xarray dataset which can be passed to the ldcpy plot functions (Call this OR `open_datasets()` before plotting.)

## Parameters

- **varnames** (*list*) – The variable(s) of interest to combine across input files (usually just one)
- **list\_of\_datasets** (*list*) – The datasets to be concatenated into a collection
- **labels** (*list*) – The respective label to access data from each dataset (also used in plotting fns)
- **\*\*kwargs** : (optional) – Additional arguments passed on to `xarray.concat()`. A list of available arguments can be found here: <https://xarray-test.readthedocs.io/en/latest/generated/xarray.concat.html>

**Returns out** – a collection containing all the data from the list datasets

**Return type** `xarray.Dataset`

`ldcpy.util.compare_stats(ds, varname: str, set1: str, set2: str, significant_digits: int = 5, include_ssim_metric: bool = False, **metrics_kwargs)`

Print error summary statistics of two DataArrays

#### Parameters

- **ds** (*xarray.Dataset*) – An xarray dataset containing multiple netCDF files concatenated across a ‘collection’ dimension
- **varname** (*str*) – The variable of interest in the dataset
- **set1** (*str*) – The collection label of the “control” data
- **set2** (*str*) – The collection label of the (1st) data to compare
- **significant\_digits** (*int, optional*) – The number of significant digits to use when printing stats, (default 5)
- **include\_ssim\_metric** (*bool, optional*) – Whether or not to compute the ssim metric, (default: False)
- **\*\*metrics\_kwargs** – Additional keyword arguments passed through to the DatasetMetrics instance.

#### Returns out

**Return type** *None*

`ldcpy.util.open_datasets(varnames, list_of_files, labels, **kwargs)`

Open several different netCDF files, concatenate across a new ‘collection’ dimension, which can be accessed with the specified labels. Stores them in an xarray dataset which can be passed to the ldcpy plot functions.

#### Parameters

- **varnames** (*list*) – The variable(s) of interest to combine across input files (usually just one)
- **list\_of\_files** (*list*) – The file paths for the netCDF file(s) to be opened
- **labels** (*list*) – The respective label to access data from each netCDF file (also used in plotting fcn)
- **\*\*kwargs** – (optional) – Additional arguments passed on to `xarray.open_mfdataset()`. A list of available arguments can be found here: [http://xarray.pydata.org/en/stable/generated/xarray.open\\_dataset.html](http://xarray.pydata.org/en/stable/generated/xarray.open_dataset.html)

**Returns out** – a collection containing all the data from the list of files

**Return type** *xarray.Dataset*

`ldcpy.util.subset_data(ds, subset=None, lat=None, lon=None, lev=None, start=None, end=None, time_dim_name='time', vertical_dim_name=None, lat_coord_name=None, lon_coord_name=None)`

Get a subset of the given dataArray, returns a dataArray

### 5.3.2 ldcpy Plot (ldcpy.plot)

```
class ldcpy.plot.MetricsPlot (ds, varname, metric, sets, group_by=None, scale='linear',
                               metric_type='raw', plot_type='spatial', transform='none',
                               subset=None, approx_lat=None, approx_lon=None, lev=0,
                               color='coolwarm', standardized_err=False, quantile=None,
                               calc_ssim=False, contour_levs=24, short_title=False,
                               axes_symmetric=False, legend_loc='upper right', vert_plot=False,
                               tex_format=False, legend_offset=None)
```

This class contains code to plot metrics in an xarray Dataset that has either 'lat' and 'lon' dimensions, or a 'time' dimension.

```
time_series_plot (da_sets, titles)
    time series plot
```

```
ldcpy.plot.plot (ds, varname, calc, sets, group_by=None, scale='linear', calc_type='raw',
                 plot_type='spatial', transform='none', subset=None, lat=None, lon=None, lev=0,
                 color='coolwarm', quantile=None, start=None, end=None, calc_ssim=False,
                 short_title=False, axes_symmetric=False, legend_loc='upper right', vert_plot=False,
                 tex_format=False, legend_offset=None)
```

Plots the data given an xarray dataset

#### Parameters

- **ds** (*xarray.Dataset*) – The input dataset
- **varname** (*str*) – The name of the variable to be plotted
- **calc** (*str*) – The name of the metric to be plotted (must match a property name in the DatasetMetrics class in ldcpy.plot, for more information about the available metrics see ldcpy.DatasetMetrics) Acceptable values include:
  - ns\_con\_var
  - ew\_con\_var
  - mean
  - std
  - variance
  - prob\_positive
  - prob\_negative
  - odds\_positive
  - zscore
  - mean\_abs
  - mean\_squared
  - rms
  - sum
  - sum\_squared
  - corr\_lag1
  - quantile
  - lag1



- `standardized_mean`
- `annual_harmonic_relative_ratio`
- `pooled_variance_ratio`
- **sets** (*list* *<str>*) – The labels of the dataset to gather metrics from
- **group\_by** (*str*) – how to group the data in time series plots. Valid groupings:
  - `time.day`
  - `time.dayofyear`
  - `time.month`
  - `time.year`
- **scale** (*str*, *optional*) – time-series y-axis plot transformation. (default “linear”) Valid options:
  - `linear`
  - `log`
- **calc\_type** (*str*, *optional*) – The type of operation to be performed on the metrics. (default ‘raw’) Valid options:
  - `raw`: the unaltered metric values
  - `diff`: the difference between the metric values in the first set and every other set
  - `ratio`: the ratio of the metric values in (2nd, 3rd, 4th... sets/1st set)
  - `metric_of_diff`: the metric value computed on the difference between the first set and every other set
- **plot\_type** (*str* , *optional*) – The type of plot to be created. (default ‘spatial’) Valid options:
  - `spatial`: a plot of the world with values at each lat and lon point (takes the mean across the time dimension)
  - `time-series`: A time-series plot of the data (computed by taking the mean across the lat and lon dimensions)
  - `histogram`: A histogram of the time-series data
- **transform** (*str*, *optional*) – data transformation. (default ‘none’) Valid options:
  - `none`
  - `log`
- **subset** (*str*, *optional*) – subset of the data to gather metrics on (default None). Valid options:
  - `first5`: the first 5 days of data
  - `winter`: data from the months December, January, February
  - `spring`: data from the months March, April, May
  - `summer`: data from the months June, July, August
  - `autumn`: data from the months September, October, November
- **lat** (*float*, *optional*) – The latitude of the data to gather metrics on (default None).

- **lon** (*float* , *optional*) – The longitude of the data to gather metrics on (default None).
- **lev** (*float*, *optional*) – The level of the data to gather metrics on (used if plotting from a 3d data set), (default 0).
- **color** (*str*, *optional*) – The color scheme for spatial plots, (default ‘coolwarm’). see [https://matplotlib.org/3.1.1/gallery/color/colormap\\_reference.html](https://matplotlib.org/3.1.1/gallery/color/colormap_reference.html) for more options
- **quantile** (*float*, *optional*) – A value between 0 and 1 required if metric=“quantile”, corresponding to the desired quantile to gather, (default 0.5).
- **start** (*int*, *optional*) – A value between 0 and the number of time slices indicating the start time of a subset, (default None).
- **end** (*int*, *optional*) – A value between 0 and the number of time slices indicating the end time of a subset, (default None)
- **calc\_ssim** (*bool*, *optional*) – Whether or not to calculate the ssim (structural similarity index) between two plots (only applies to plot\_type = ‘spatial’), (default False)
- **short\_title** (*bool*, *optional*) – If True, use a shortened title in the plot output (default False).
- **axes\_symmetric** (*bool*, *optional*) – Whether or not to make the colorbar axes symmetric about zero (used in a spatial plot) (default False)
- **legend\_loc** (*str*, *optional*) – The location to put the legend in a time-series plot in single-column format (plot\_type = “time\_series”, vert\_plot=True) (default “upper right”)
- **vert\_plot** (*bool*, *optional*) – If true, forces plots into a single column format and enlarges text. (default False)
- **tex\_format** (*bool*, *optional*) – Whether to interpret all plot output strings as latex formatting (default False)
- **legend\_offset** (*2-tuple*, *optional*) – The x- and y- offset of the legend. Moves the corner of the legend specified by legend\_loc to the specified location specified (where (0,0) is the bottom left corner of the plot and (1,1) is the top right corner). Only affects time-series, histogram, and periodogram plots.

**Returns out**

**Return type** *None*

`ldcpy.plot.tex_escape(text)`

**Parameters** **text** – a plain text message

**Returns** the message escaped to appear correctly in LaTeX

### 5.3.3 Idcpy Metrics (ldcpy.metrics)

```
class ldcpy.metrics.DatasetMetrics(ds: xarray.core.dataarray.DataArray, aggregate_dims:
    list, time_dim_name: str = 'time', lat_dim_name: Optional[str] = None, lon_dim_name: Optional[str] = None,
    vert_dim_name: Optional[str] = None, lat_coord_name: Optional[str] = None, lon_coord_name: Optional[str] =
    None, q: float = 0.5, spre_tol: float = 0.0001)
```

This class contains metrics for each point of a dataset after aggregating across one or more dimensions, and a method to access these metrics. Expects a DataArray.

**get\_metric** (*name: str, q: Optional[int] = 0.5, grouping: Optional[str] = None, ddof=1*)

Gets a metric aggregated across one or more dimensions of the dataset

**Parameters**

- **name** (*str*) – The name of the metric (must be identical to a property name)
- **q** (*float, optional*) – (default 0.5)

**Returns out** – A DataArray of the same size and dimensions the original dataarray, minus those dimensions that were aggregated across.

**Return type** `xarray.DataArray`

**get\_single\_metric** (*name: str*)

Gets a metric consisting of a single float value

**Parameters name** (*str*) – the name of the metric (must be identical to a property name)

**Returns out** – The metric value

**Return type** `float`

**property annual\_harmonic\_relative\_ratio**

The annual harmonic relative to the average periodogram value in a neighborhood of 50 frequencies around the annual frequency NOTE: This assumes the values along the “time” dimension are equally spaced. NOTE: This metric returns a lat-lon array regardless of aggregate dimensions, so can only be used in a spatial plot.

**property annual\_harmonic\_relative\_ratio\_pct\_sig**

The percentage of points past the significance cutoff (p value  $\leq 0.01$ ) for the annual harmonic relative to the average periodogram value in a neighborhood of 50 frequencies around the annual frequency

**property ew\_con\_var**

The East-West Contrast Variance averaged along the aggregate dimensions

**property lag1**

The deseasonalized lag-1 autocorrelation value by day of year NOTE: This metric returns an array of spatial values as the data set regardless of aggregate dimensions, so can only be plotted in a spatial plot.

**property lag1\_first\_difference**

The deseasonalized lag-1 autocorrelation value of the first difference of the data by day of year NOTE: This metric returns an array of spatial values as the data set regardless of aggregate dimensions, so can only be plotted in a spatial plot.

**property mae\_day\_max**

The day of maximum mean absolute value at the point. NOTE: only available in spatial and spatial comparison plots

**property mean**

The mean along the aggregate dimensions

**property mean\_abs**

The mean of the absolute errors along the aggregate dimensions

**property mean\_squared**

The absolute value of the mean along the aggregate dimensions

**property ns\_con\_var**

The North-South Contrast Variance averaged along the aggregate dimensions

**property odds\_positive**

The odds that a point is positive =  $\text{prob\_positive}/(1-\text{prob\_positive})$

**property pooled\_variance**

The overall variance of the dataset

**property pooled\_variance\_ratio**

The pooled variance along the aggregate dimensions

**property prob\_negative**

The probability that a point is negative

**property prob\_positive**

The probability that a point is positive

**property root\_mean\_squared**

The absolute value of the mean along the aggregate dimensions

**property standardized\_mean**

The mean at each point along the aggregate dimensions divided by the standard deviation NOTE: will always be 0 if aggregating over all dimensions

**property std**

The standard deviation along the aggregate dimensions

**property variance**

The variance along the aggregate dimensions

**property zscore**

The z-score of a point averaged along the aggregate dimensions under the null hypothesis that the true mean is zero. NOTE: currently assumes we are aggregating along the time dimension so is only suitable for a spatial plot.

**property zscore\_cutoff**

The Z-Score cutoff for a point to be considered significant

**property zscore\_percent\_significant**

The percent of points where the zscore is considered significant

```
class ldcpy.metrics.DiffMetrics(ds1: xarray.core.dataarray.DataArray, ds2: xarray.core.dataarray.DataArray, aggregate_dims: Optional[list] = None, **metrics_kwargs)
```

This class contains metrics on the overall dataset that require more than one input dataset to compute

**get\_diff\_metric** (*name: str*)

Gets a metric on the dataset that requires more than one input dataset

**Parameters** *name* (*str*) – The name of the metric (must be identical to a property name)

**Returns** *out*

**Return type** *float*

**property covariance**

The covariance between the two datasets

**property ks\_p\_value**

The Kolmogorov-Smirnov p-value

**property max\_spatial\_rel\_error**

We compute the relative error at each grid point and return the maximum.

**property normalized\_max\_pointwise\_error**

The absolute value of the maximum pointwise difference, normalized by the range of values for the first set

**property normalized\_root\_mean\_squared**

The absolute value of the mean along the aggregate dimensions, normalized by the range of values for the first set

**property pearson\_correlation\_coefficient**

returns the pearson correlation coefficient between the two datasets

**property spatial\_rel\_error**

At each grid point, we compute the relative error. Then we report the percentage of grid point whose relative error is above the specified tolerance (1e-4 by default).

**property ssim\_value**

We compute the SSIM (structural similarity index) on the visualization of the spatial data.

**property ssim\_value\_fp**

We compute the SSIM (structural similarity index) on the spatial data - using the data itself (we do not create an image).

Here we scale from [0,1] - then quantize to 256 bins

**property ssim\_value\_fp\_old**

To mimic what zchecker does - the ssim on the fp data with original constants and no scaling. This will return Nan on POP data.

## 5.4 Examples

### 5.4.1 Tutorial

ldcpy is a utility for gathering and plotting metrics from NetCDF or Zarr files using the Pangeo stack. This tutorial notebook targets comparing CESM data in its original form to CESM data that has undergone lossy compression (meaning that the reconstructed file is not exactly equivalent to the original file). The tools provided in ldcpy are intended to highlight differences due to compression artifacts in order to assist scientist in evaluating the amount of lossy compression to apply to their data.

The CESM data used here are NetCDF files in “timeseries” file format, meaning that each NetCDF file contains one (major) output variable (e.g., surface temperature or precipitation rate) and spans multiple timesteps (daily, monthly, 6-hourly, etc.). CESM timeseries files are regularly released in large public datasets.

```
[1]: # Add ldcpy root to system path
import sys

sys.path.insert(0, '../..../')

# Import ldcpy package
# Autoreloads package everytime the package is called, so changes to code will be
# reflected in the notebook if the above sys.path.insert(...) line is uncommented.
%load_ext autoreload
%autoreload 2

# suppress all of the divide by zero warnings
import warnings

warnings.filterwarnings("ignore")

import ldcpy
```

(continues on next page)

(continued from previous page)

```
# display the plots in this notebook
%matplotlib inline
```

## Overview

This notebook demonstrates the use of ldcpy on the sample data included with this package. It explains how to open datasets (and view metadata), display basic statistics about the data, and create both time-series and spatial plots of the datasets and related metrics. Plot examples start out with the essential arguments, and subsequent examples explore the additional plotting options that are available.

For information about installation, see [these instructions](#), and for information about usage, see the API reference [here](#).

## Loading Datasets and Viewing Metadata

The first step in comparing the data is to load the data from the files that we are interested into a “collection” for ldcpy to use. To do this, we use `ldcpy.open_datasets()`. This function requires the following three arguments:

- *varnames* : the variable(s) of interest to combine across files (typically the timeseries file variable name)
- *list\_of\_files* : a list of full file paths (either relative or absolute)
- *labels* : a corresponding list of names (or labels) for each file in the collection

Note: This function is a wrapper for `xarray.open_mfdatasets()`, and any additional key/value pairs passed in as a dictionary are used as arguments to `xarray.open_mfdatasets()`. For example, specifying the chunk size (“chunks”) will be important for large data (see `LargeDataGladeNotebook.ipynb` for more information and an example).

We setup three different collections of timeseries datasets in these examples:

- *col\_ts* contains daily surface temperature (TS) data (2D data) for 100 days
- *col\_prect* contains daily precipitation rate (PRECT) data (2D data) for 60 days
- *col\_t* contains monthly temperature (T) data (3D data) for 3 months

These datasets are collections of variable data from several different netCDF files, which are given labels in the third parameter to the `ldcpy.open_datasets()` function. These names/labels can be whatever you want (e.g., “orig”, “control”, “bob”, ...), but they should be informative because the names will be used to select the appropriate dataset later and as part of the plot titles.

In this example, in each dataset collection we include a file with the original (uncompressed) data as well as additional file(s) with the same data subject to different levels of lossy compression.

*Note: If you do not need to get the data from files (e.g., you have already used `xarray.open_dataset()`), then use `ldcpy.collect_datasets()` instead of `ldcpy.open_datasets` (see example in `AWSDataNotebook.ipynb`).*

```
[2]: # col_ts is a collection containing TS data
col_ts = ldcpy.open_datasets(
    ["TS"],
    [
        "../.../data/cam-fv/orig.TS.100days.nc",
        "../.../data/cam-fv/zfp1.0.TS.100days.nc",
        "../.../data/cam-fv/zfp1e-1.TS.100days.nc",
    ],
    ["orig", "zfpA1.0", "zfpA1e-1"],
)
# col_prect contains PRECT data
```

(continues on next page)

(continued from previous page)

```

col_prect = ldcpy.open_datasets(
    ["PRECT"],
    [
        "../.../data/cam-fv/orig.PRECT.60days.nc",
        "../.../data/cam-fv/zfp1e-7.PRECT.60days.nc",
        "../.../data/cam-fv/zfp1e-11.PRECT.60days.nc",
    ],
    ["orig", "zfp1e-7", "zfp1e-11"],
)
# col_t contains 3D T data (here we specify the chunk to be a single timeslice)
col_t = ldcpy.open_datasets(
    ["T"],
    [
        "../.../data/cam-fv/cam-fv.T.3months.nc",
        "../.../data/cam-fv/c.fzip.cam-fv.T.3months.nc",
    ],
    ["orig", "comp"],
    chunks={"time": 1},
)

dataset size in GB 0.07

dataset size in GB 0.04

dataset size in GB 0.04

```

Note that running the `open_datasets` function (as above) prints out the size of each dataset collection. For `col_prect`, the `chunks` parameter is used by DASK (which is further explained in [LargeDataGladeNotebook.ipynb](#)).

Printing a dataset collection reveals the dimension names, sizes, datatypes and values, among other metadata. The dimensions and the length of each dimension are listed at the top of the output. Coordinates list the dimensions vertically, along with the data type of each dimension and the coordinate values of the dimension (for example, we can see that the 192 latitude data points are spaced evenly between -90 and 90). Data variables lists all the variables available in the dataset. For these timeseries files, only the one (major) variable will be of interest. For `col_t`, that variable is temperature (T), which was specified in the first argument of the `open_datasets()` call. The so-called major variable will have the required “lat”, “lon”, and “time” dimensions. If the variable is 3D (as in this example), a “lev” dimension will indicate that the dataset contains values at multiple altitudes (here, `lev=30`). Finally, a “collection” dimension indicates that we concatenated several datasets together. (In this `col_t` example, we concatenated 2 files together.)

```

[3]: # print information about col_t
col_t

[3]: <xarray.Dataset>
Dimensions:      (collection: 2, lat: 192, lev: 30, lon: 288, time: 3)
Coordinates:
  * lat          (lat) float64 -90.0 -89.06 -88.12 -87.17 ... 88.12 89.06 90.0
  * lev          (lev) float64 3.643 7.595 14.36 24.61 ... 957.5 976.3 992.6
  * lon          (lon) float64 0.0 1.25 2.5 3.75 5.0 ... 355.0 356.2 357.5 358.8
  * time         (time) object 1920-02-01 00:00:00 ... 1920-04-01 00:00:00
  * collection   (collection) <U4 'orig' 'comp'
Data variables:
  T              (collection, time, lev, lat, lon) float32 dask.array<chunksize=(1, 1, ↵
    ↵30, 192, 288), meta=np.ndarray>
Attributes:

```

(continues on next page)

(continued from previous page)

```

Conventions:      CF-1.0
source:           CAM
case:             b.e11.B20TRC5CNBDRD.f09_g16.031
title:           UNSET
logname:          mickelso
host:             ys0219
Version:          $Name$
revision_Id:      $Id$
initial_file:     b.e11.B20TRC5CNBDRD.f09_g16.001.cam.i.1920-01-01-00000.nc
topography_file:  /glade/p/cesmdata/cseg/inputdata/atm/cam/topo/USGS-gtop...
history:          Thu Jul 9 14:15:11 2020: ncks -d time,0,2,1 cam-fv.T.6...
NCO:             netCDF Operators version 4.7.9 (Homepage = http://nco.s...

```

## Comparing Summary Statistics

The `compare_stats` function can be used to compute and compare the overall statistics for a **single** timeslice in two datasets. To use this function, four arguments are required. In order, they are:

- *ds* - a single time slice in a collection of datasets read in from `ldcpy.open_datasets()` or `ldcpy.collect_datasets()`
- *varname* - the variable name we want to get statistics for (in this case ‘TS’ is the variable in our dataset collection *col\_ts* )
- *set1* - the label of one particular dataset in the collection we are interested in (‘orig’)
- *set2* - and the label of another dataset that we want to compare it with (‘zfpA1.0’).

Additionally, two optional arguments can be specified:

- *significant\_digits* - the number of significant digits to print (default = 5)
- *include\_ssim\_metric* - include the ssim metric (default = False. Note this takes a bit of time for 3D vars)

```

[4]: # print 'TS' statistics about 'orig', 'zfpA1.0', and diff between the two datasets
      # for time slice = 0
      ds = col_ts.isel(time=0)
      ldcpy.compare_stats(ds, "TS", "orig", "zfpA1.0", include_ssim_metric=True)

```

```

mean orig                : 274.71
mean zfpA1.0             : 274.71
mean diff                 : 0.0057673

variance orig            : 534.01
variance zfpA1.0         : 533.68

standard deviation orig   : 23.109
standard deviation zfpA1.0 : 23.102

max value orig           : 315.58
max value zfpA1.0        : 315.57
min value orig           : 216.74
min value zfpA1.0        : 216.82

max abs diff             : 0.40588
min abs diff             : 0
mean abs diff            : 0.05852
mean squared diff        : 3.3262e-05

```

(continues on next page)



(continued from previous page)

```

root mean squared diff      : 0.075273
normalized root mean squared diff : 0.00076154
normalized max pointwise error : 0.0041064
pearson correlation coefficient : 0.99999
ks p-value                  : 1
spatial relative error(% > 0.0001) : 68.958
max spatial relative error   : 0.001474
ssim                        : 0.99845
ssim_fp                     : 0.9822

```

We can also generate metrics on a particular dataset. While this is done “behind the scenes” with the plotting functions, we first demonstrate here how the user can access this data without creating a plot.

We use an object of type `ldcpy.DatasetMetrics` to gather metrics on a dataset. To create an `ldcpy.DatasetMetrics` object, we first grab the particular dataset from our collection that we are interested in (in the usual xarray manner). For example, the following will grab the data for the TS variable labeled ‘orig’ in the `col_ts` dataset that we created:

```
[5]: # get the orig dataset
my_data = col_ts["TS"].sel(collection="orig")
```

Then we create a `DatasetMetrics` object using the data and a list of dimensions that we want to aggregate the data along. We usually want to aggregate all of the timeslices (“time”) or spatial points (“lat” and “lon”):

```
[6]: ds_metrics_across_time = ldcpy.DatasetMetrics(my_data, ["time"])
ds_metrics_across_space = ldcpy.DatasetMetrics(my_data, ["lat", "lon"])
```

Now when we call the `get_metric()` method on this class, a metric will be computed across each of these specified dimensions. For example, below we compute the “mean” across time.

```
[7]: my_data_mean_across_time = ds_metrics_across_time.get_metric("mean")

# trigger computation
my_data_mean_across_time.load()

[7]: <xarray.DataArray 'TS' (lat: 192, lon: 288)>
array([[228.94805908, 229.39446915, 229.42184357, ..., 229.37858994,
        228.95278442, 229.4450769 ],
       [229.3337999 , 229.38760071, 229.32789291, ..., 228.73066544,
        229.06293549, 229.02659225],
       [229.74119156, 229.81762497, 229.76844055, ..., 229.55835098,
        229.75506577, 229.57969193],
       ...,
       [237.00477859, 237.03717865, 237.06771225, ..., 236.91639221,
        236.94580048, 236.97352493],
       [236.69006699, 236.70527969, 236.7207486 , ..., 236.63906601,
        236.65667786, 236.6738382 ],
       [236.29601639, 236.29747787, 236.29879272, ..., 236.29072723,
        236.29267105, 236.29443573]])
Coordinates:
  * lat      (lat) float64 -90.0 -89.06 -88.12 -87.17 ... 88.12 89.06 90.0
  * lon      (lon) float64 0.0 1.25 2.5 3.75 5.0 ... 355.0 356.2 357.5 358.8
    collection <U8 'orig'
Attributes:
  units:      K
  long_name:  Surface temperature (radiative)
  cell_methods: time: mean
```

```
[8]: # Here just ask for the spatial mean at the first time step
my_data_mean_across_space = ds_metrics_across_space.get_metric("mean").isel(time=0)
# trigger computation
my_data_mean_across_space.load()

[8]: <xarray.DataArray 'TS' ()>
array(274.71370277)
Coordinates:
  time          object 1920-01-01 00:00:00
  collection    <U8 'orig'
Attributes:
  units:        K
  long_name:    Surface temperature (radiative)
  cell_methods: time: mean
```

There are many currently available metrics to choose from. A complete list of metrics that can be passed to `get_metric()` is available [here](#).

## Spatial Plots

### A Basic Spatial Plot

First we demonstrate the most basic usage of the `ldcpy.plot()` function. In its simplest form, `ldcpy.plot()` plots a single spatial plot of a dataset (i.e., no comparison) and requires:

- *ds* - the collection of datasets read in from `ldcpy.open_datasets()` (or `ldcpy.collect_datasets()`)
- *varname* - the variable name we want to get statistics for
- *sets* - a list of the labels of the datasets in the collection that we are interested in
- *calc* - the name of the metric to be plotted

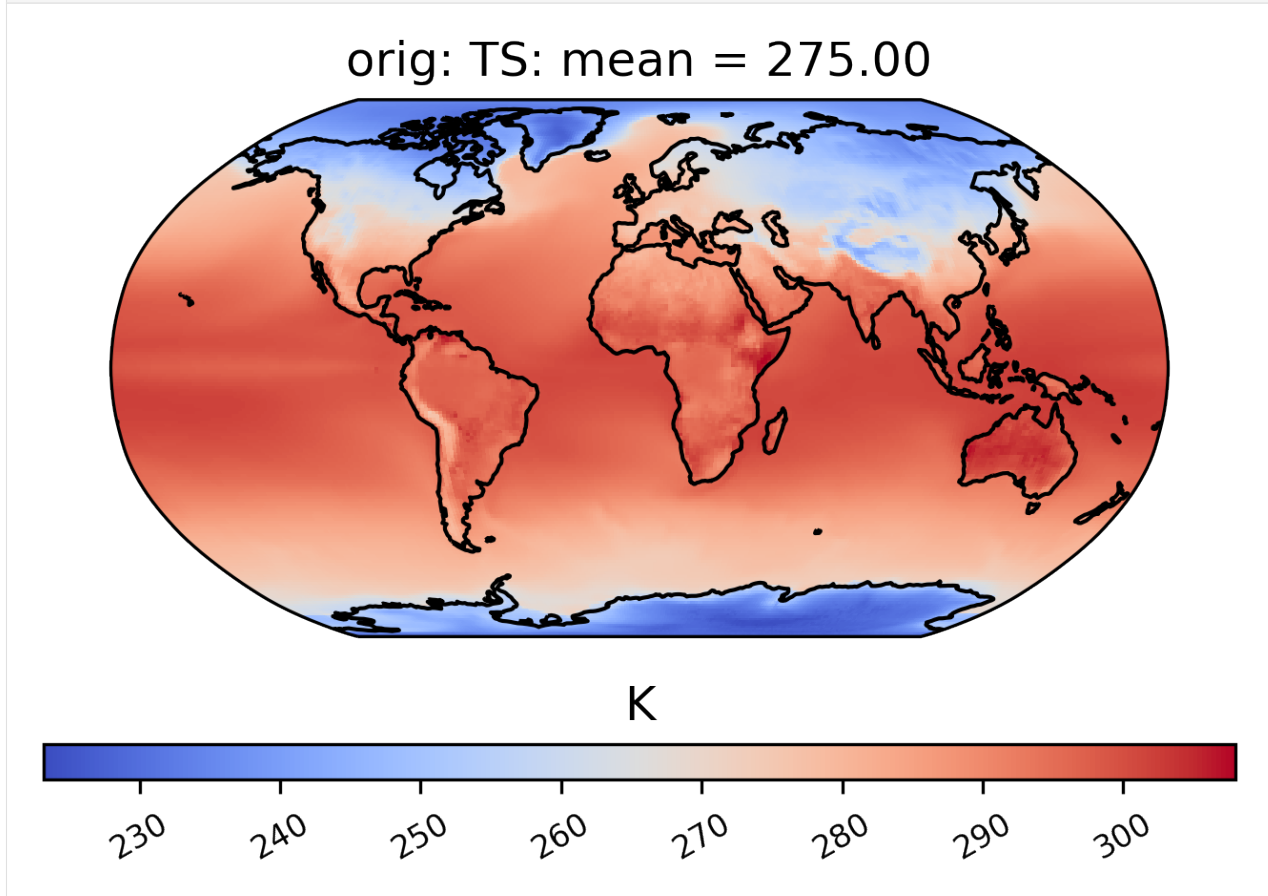
There are a number of optional arguments as well, that will be demonstrated in subsequent plots. These options include:

- *group\_by*
- *scale*
- *calc\_type*
- *plot\_type*
- *transform*
- *subset*
- *lat*
- *lon*
- *lev*
- *color*
- *quantile*
- *start*
- *end*

A full explanation of each optional argument is available in the documentation [here](#) - as well as all available options. By default, a spatial plot of this data is created. The title of the plot contains the name of the dataset, the variable of interest, and the metric that is plotted.

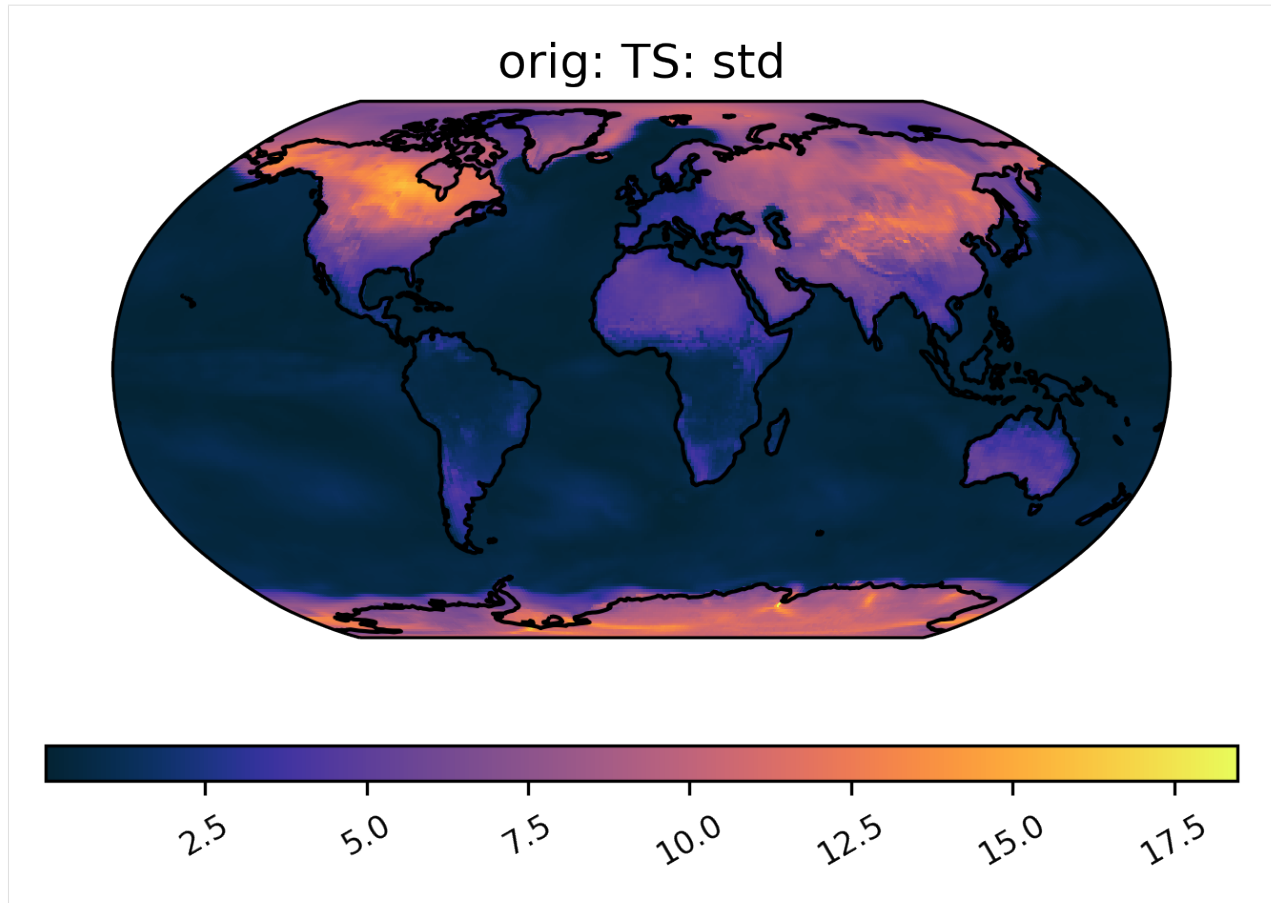
The following command generates a plot showing the mean TS (surface temperature) value from our dataset collection *col\_ts* over time at each point in the dataset labeled 'orig'. Note that plots can be saved by using the matplotlib function `savefig`:

```
[9]: ldcpy.plot(col_ts, "TS", sets=["orig"], calc="mean", plot_type="spatial")
# Uncomment to save this image to the filesystem
# import matplotlib.pyplot as plt
# plt.savefig(f"MYFIGURE.png", bbox_inches='tight')
```



We can also plot metrics other than the mean, such as the standard deviation at each grid point over time. We can also change the color scheme (for a full list of metrics and color schemes, see the [documentation](#)). Here is an example of a plot of the same dataset from *col\_ts* as above but using a different color scheme and the standard deviation metric. Notice that we do not specify the 'plot\_type' this time, because it defaults to 'spatial':

```
[10]: # plot of the standard deviation of TS values in the col_ds "orig" dataset
ldcpy.plot(col_ts, "TS", sets=["orig"], calc="std", color="cmo.thermal")
```

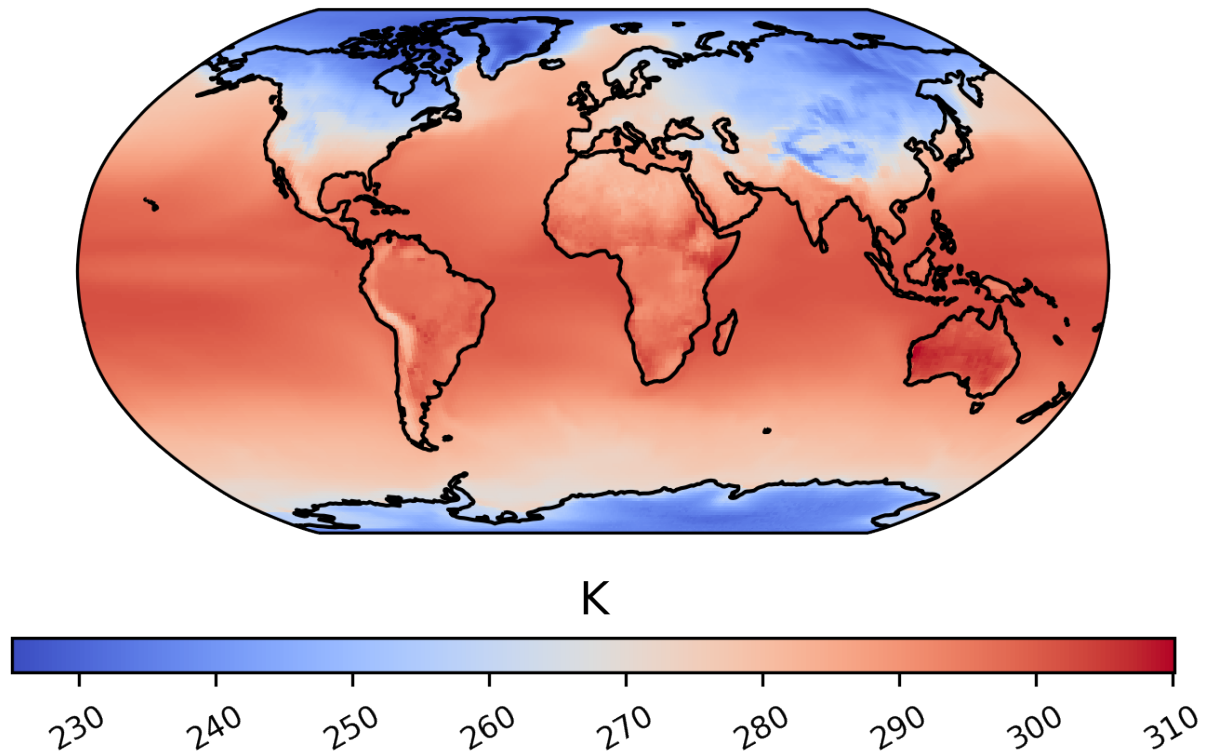


### Spatial Subsetting

Plotting the metric of a subset of the data (e.g., not all of the time slices in the dataset) is possible using the `subset` keyword. In the plot below, we just look at “winter” data (Dec., Jan., Feb.). Other options for `subset` are available [here](#).

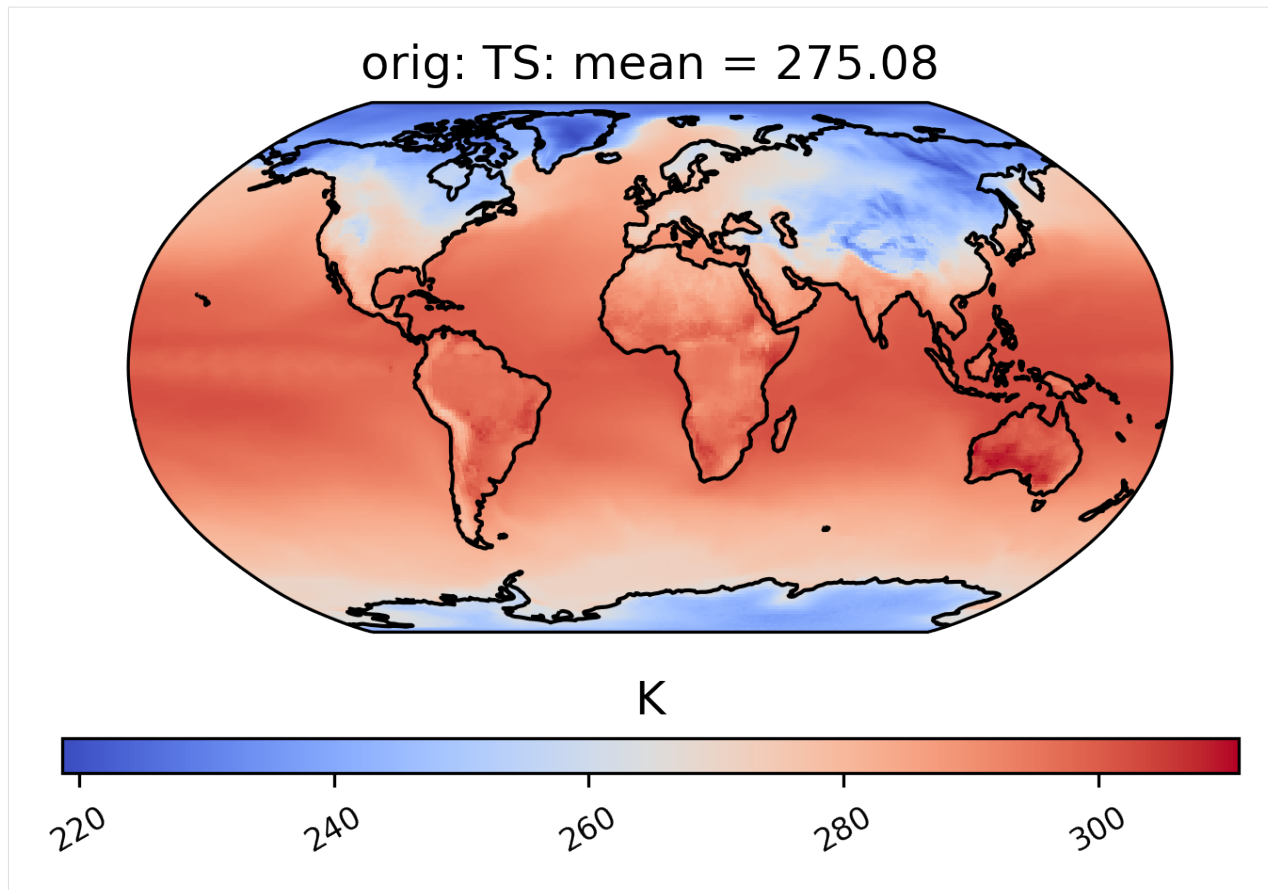
```
[11]: # plot of mean winter TS values in the col_ts "orig" dataset
ldcpy.plot(col_ts, "TS", sets=["orig"], calc="mean", subset="winter")
```

orig: TS: mean = 275.02 subset:winter



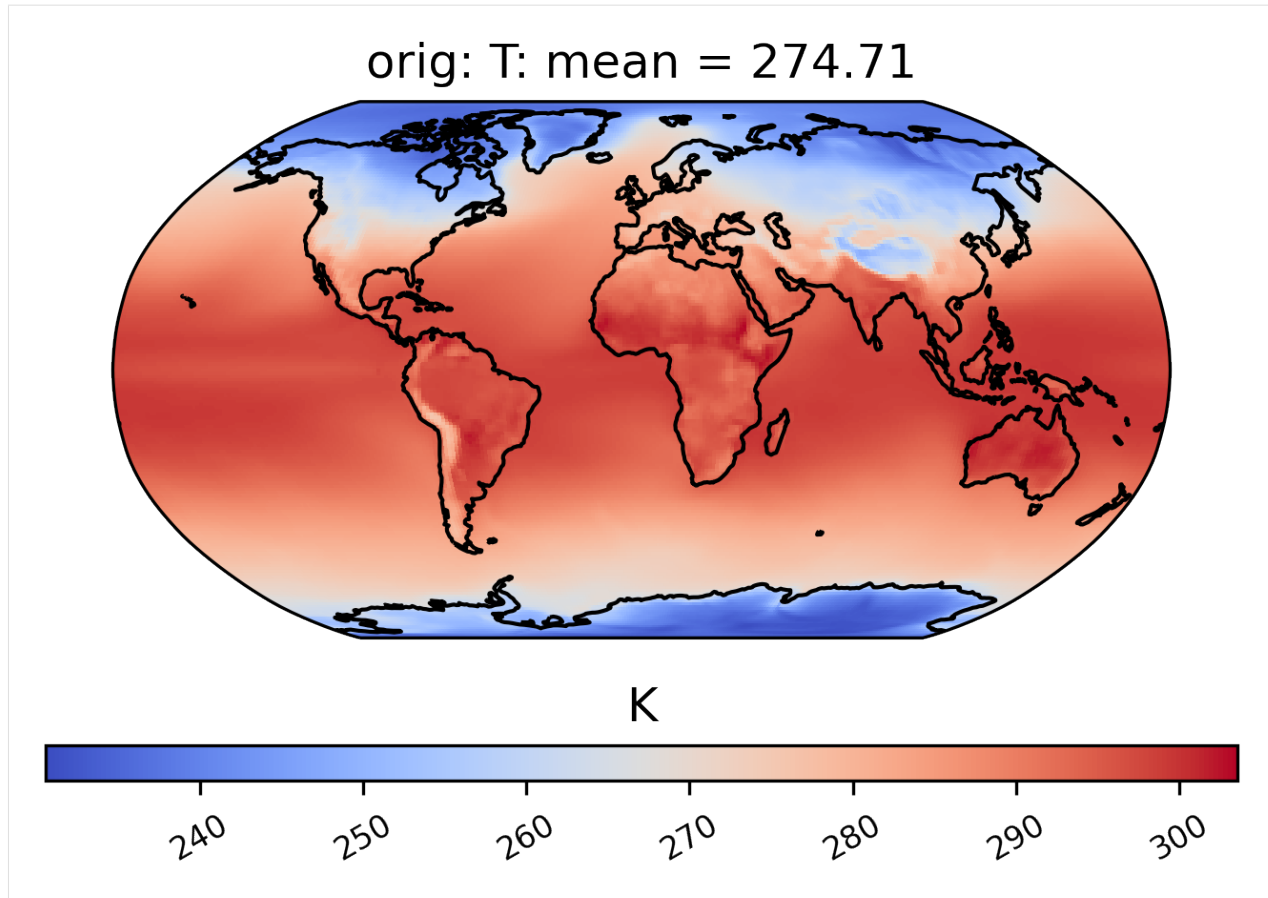
It is also possible to plot metrics for a subset of the time slices by specifying the start and end indices of the time we are interested in. This command creates a spatial plot of the mean TS values in “orig” for the first five days of data:

```
[12]: # plot of the first 5 TS values in the ds orig dataset
ldcpy.plot(col_ts, "TS", sets=["orig"], calc="mean", start=0, end=5)
```



Finally, for a 3D dataset, we can specify which vertical level to view using the “lev” keyword. Note that “lev” is a dimension in our dataset *col\_t* (see printed output for *col\_t* above), and in this case lev=30, meaning that lev ranges from 0 and 29, where 0 is at the surface (by default, lev=0):

```
[13]: # plot of T values at lev=29 in the col_t "orig" dataset
ldcpy.plot(col_t, "T", sets=["orig"], calc="mean", lev=29)
```



### Spatial Comparisons and Diff Plots

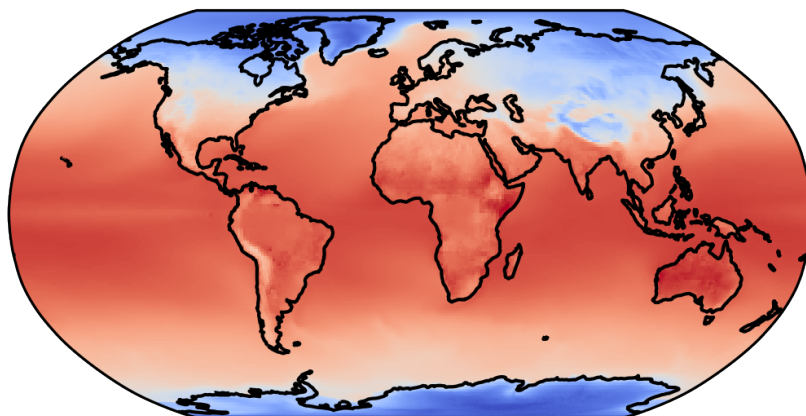
If we want a side-by-side comparison of two datasets, we need to specify an additional dataset name in the `sets` argument. The plot below shows the mean TS value over time at each point in the ‘orig’ (original), ‘zfpA1.0’ (compressed with zfp, absolute error tolerance 1.0) and ‘zfpA1e-1’ (compressed with zfp, absolute error tolerance 0.1) datasets in collection `col_ts`. Note that the “`calc_ssim`” argument indicates whether to calculate the SSIM (structural similarity index) comparing the first plot (SSIM is between 0.0 and 1.0, where 1.0 means the plots are identical) to the second, third, etc... The “`vert_plot`” argument indicates that the arrangement of the subplots should be vertical.

```
[14]: # comparison between mean TS values in col_ts for "orig" and "zfpA1.0" datasets
ldcpy.plot(
    col_ts,
    "TS",
    sets=["orig", "zfpA1.0", "zfpA1e-1"],
    calc="mean",
    calc_ssim=True,
    vert_plot=True,
)

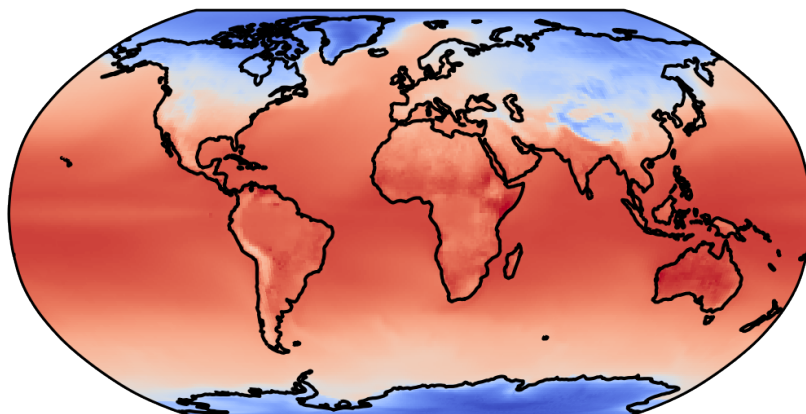
SSIM 1 & 2 = 0.99719

SSIM 1 & 3 = 0.99974
```

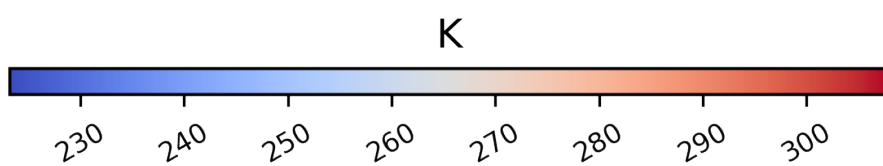
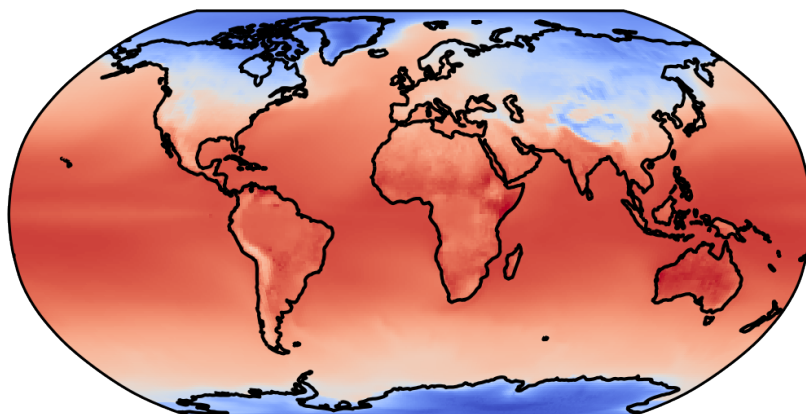
orig: TS: mean = 275.00



zfpA1.0: TS: mean = 274.99



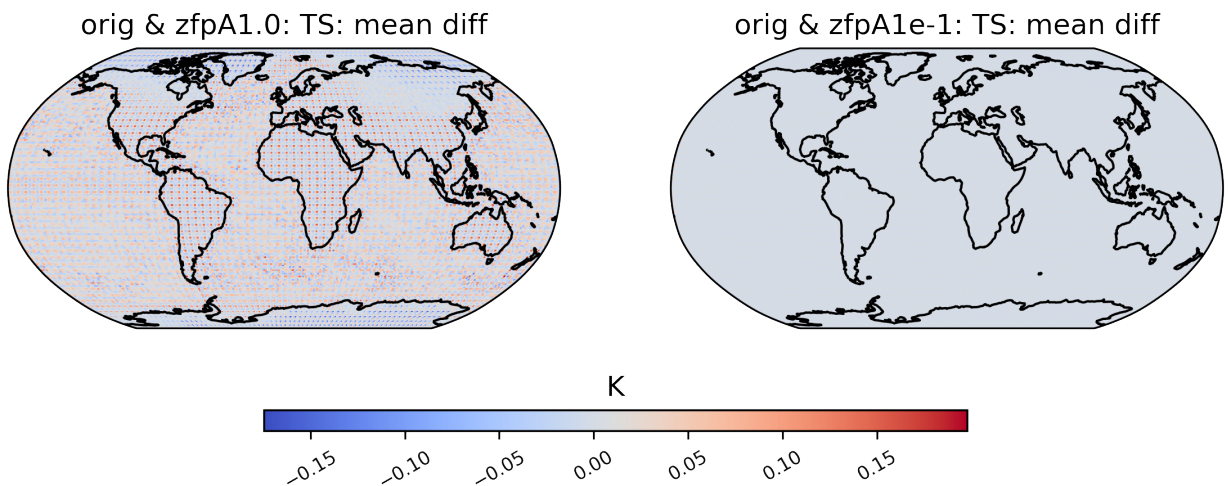
zfpA1e-1: TS: mean = 274.99





To the eye, these plots look identical. This is because the effects of compression are small compared to the magnitude of the data. We can view the compression artifacts more clearly by plotting the difference between two plots. This can be done by setting the `calc_type` to 'diff':

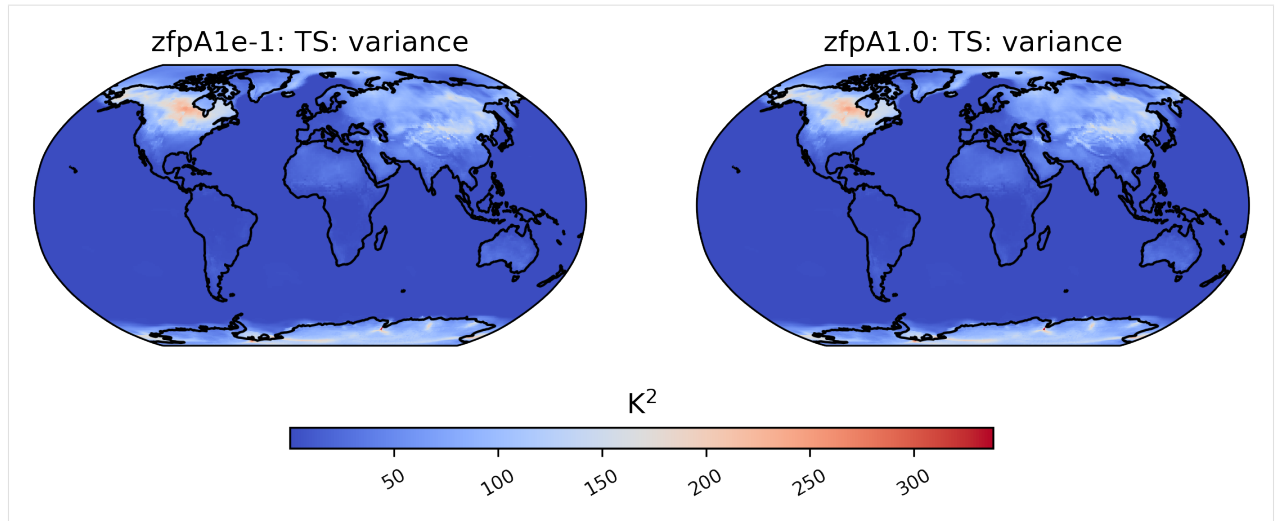
```
[15]: # diff between mean TS values in col_ds "orig" and "zfpA1.0" datasets
ldcpy.plot(
    col_ts,
    "TS",
    sets=["orig", "zfpA1.0", "zfpA1e-1"],
    calc="mean",
    calc_type="diff",
)
```



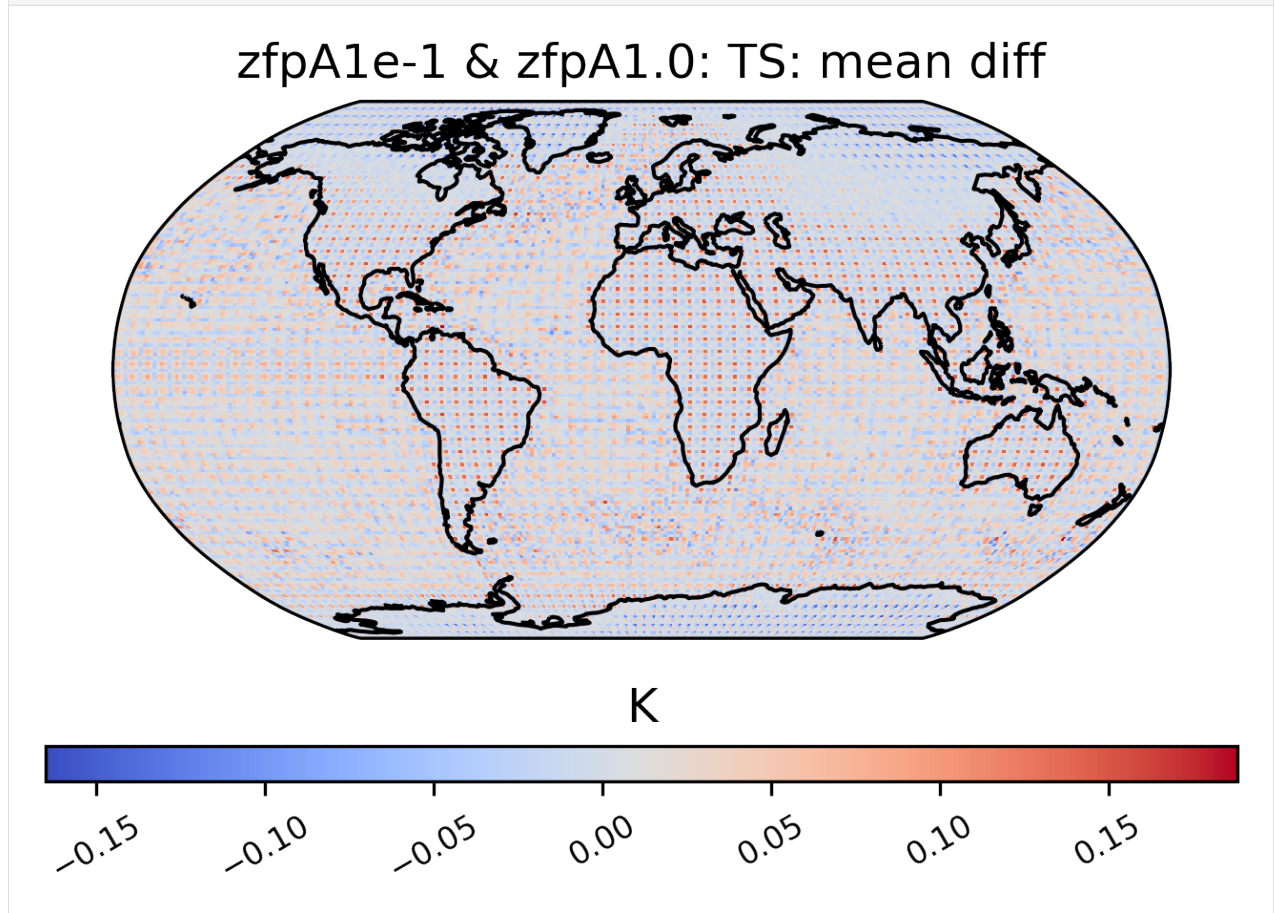
We are not limited to comparing side-by-side plots of the original and compressed data. It is also possible to compare two different compressed datasets side-by-side as well, by using a different dataset name for the first element in `sets`:

```
[16]: # comparison between variance of TS values in col_ts for "zfpA1e-1" and "zfpA1.0"
      ↪ datasets"
ldcpy.plot(
    col_ts,
    "TS",
    sets=["zfpA1e-1", "zfpA1.0"],
    calc="variance",
    calc_ssim=True,
)
```

SSIM 1 & 2 = 0.99972

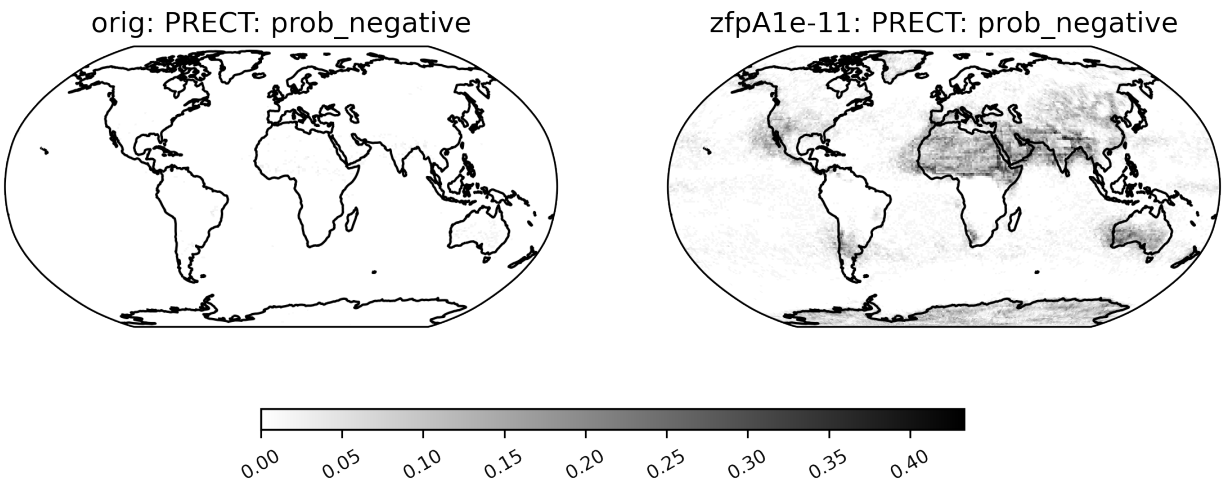


```
[17]: # diff between mean TS values in col_ts for "zfpA1e-1" and "zfpA1.0" datasets
ldcpy.plot(
    col_ts,
    "TS",
    sets=["zfpA1e-1", "zfpA1.0"],
    calc="mean",
    calc_type="diff",
)
```



Sometimes comparison plots can look strikingly different, indicating a potential problem with the compression. This plot shows the probability of negative rainfall (PRECT). We would expect this metric to be zero everywhere on the globe (because negative rainfall does not make sense!), but the compressed output shows regions where the probability is significantly higher than zero:

```
[18]: ldcpy.plot(
    col_prect,
    "PRECT",
    sets=["orig", "zfpA1e-11"],
    calc="prob_negative",
    color="binary",
)
```



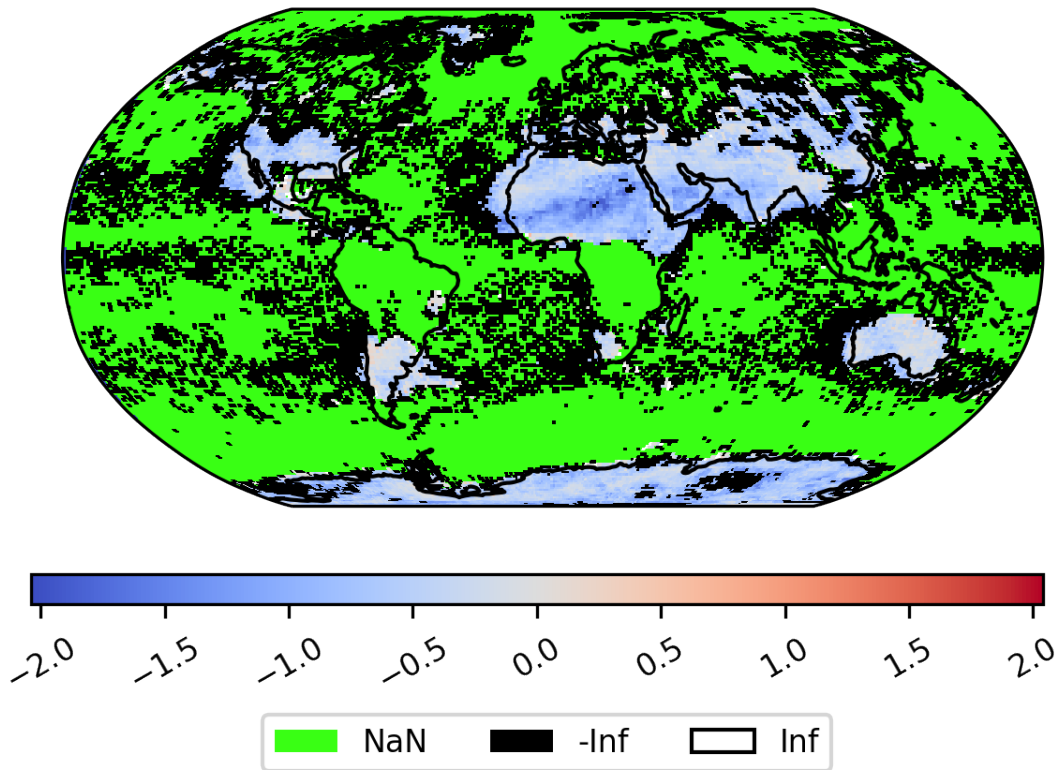
### Unusual Values (inf, -inf, NaN)

Some metrics result in values that are +/- infinity, or NaN (likely resulting from operations like 0/0 or inf/inf). NaN values are plotted in neon green, infinity is plotted in white, and negative infinity is plotted in black (regardless of color scheme). If infinite values are present in the plot data, arrows on either side of the colorbar are shown to indicate the color for +/- infinity. This plot shows the log of the ratio of the odds of positive rainfall over time in the compressed and original output,  $\log(\text{odds\_positive compressed} / \text{odds\_positive original})$ . Here we are suppressing all of the divide by zero warnings for aesthetic reasons.

The following plot showcases some interesting plot features. We can see sections of the map that take NaN values, and other sections that are black because taking the log transform has resulted in many points with the value -inf:

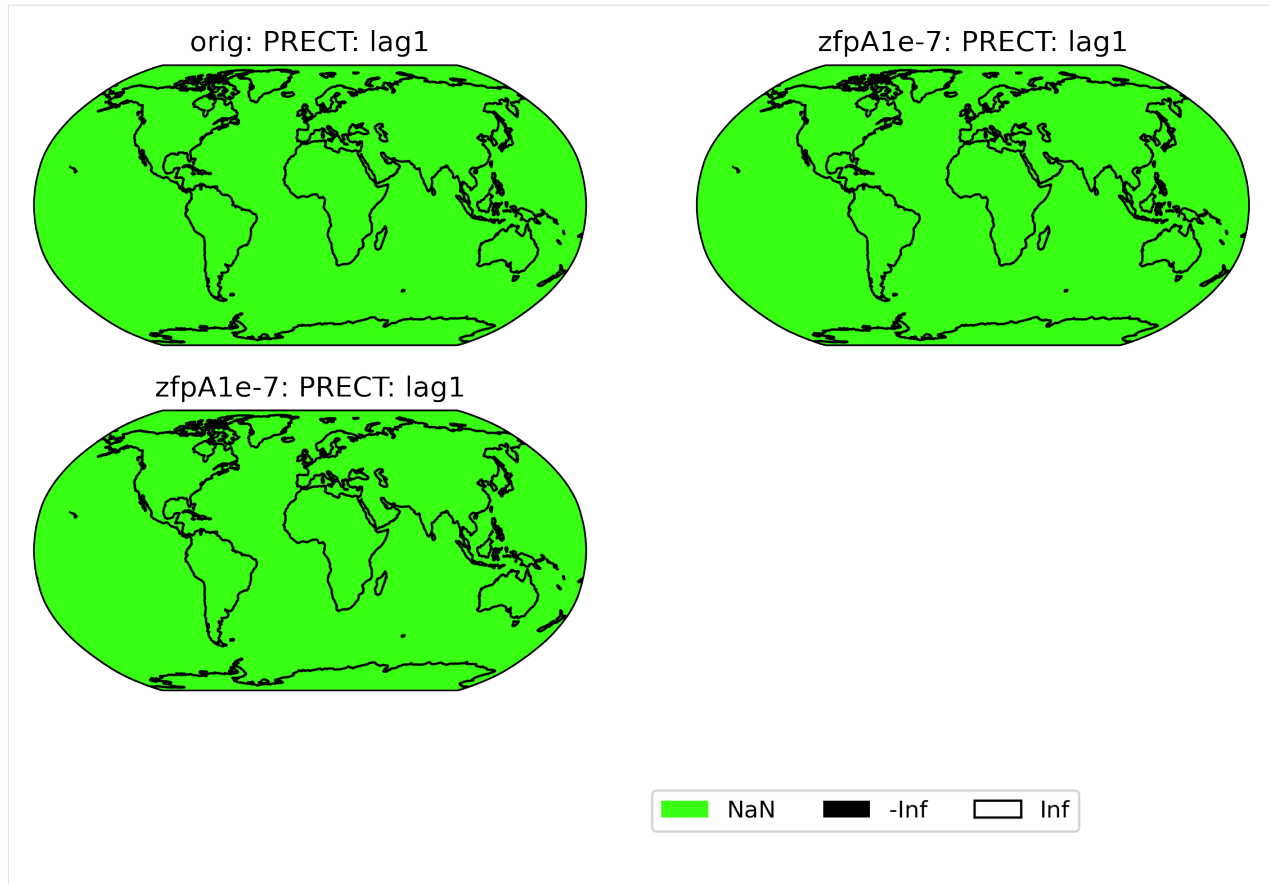
```
[19]: # plot of the ratio of odds positive PRECT values in col_prect zfpA1.0 dataset / col_
      ↪prect orig dataset (log transform)
ldcpy.plot(
    col_prect,
    "PRECT",
    sets=["orig", "zfpA1e-11"],
    calc="odds_positive",
    calc_type="ratio",
    transform="log",
    axes_symmetric=True,
    vert_plot=True,
)
```

orig & zfpA1e-11: PRECT: log10 odds\_positive ratio



If all values are NaN, then the colorbar is not shown. Instead, a legend is shown indicating the green(!) color of NaN values, and the whole plot is colored gray. (If all values are infinite, then the plot is displayed normally with all values either black or white). Because the example dataset only contains 60 days of data, the deseasonalized lag-1 values and their variances are all 0, and so calculating the correlation of the lag-1 values will involve computing  $0/0 = \text{NaN}$ :

```
[20]: # plot of lag-1 correlation of PRECT values in col_prect orig dataset
ldcpy.plot(col_prect, "PRECT", sets=["orig", "zfpA1e-7", "zfpA1e-7"], calc="lag1")
```



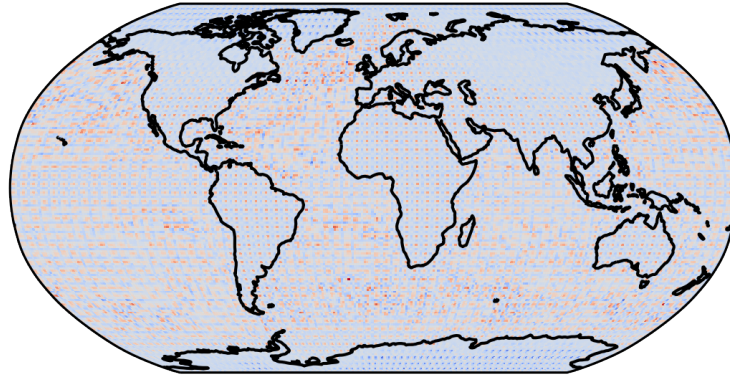
### Other Spatial Plots

Sometimes, we may want to compute a metric on the difference between two datasets. For instance, the `zscore` metric calculates the zscore at each point under the null hypothesis that the true mean is zero, so using the “`metric_of_diff`” `calc_type` calculates the zscore of the diff between two datasets (to find the values that are significantly different between the two datasets). The zscore metric in particular gives additional information about the percentage of significant gridpoints in the plot title:

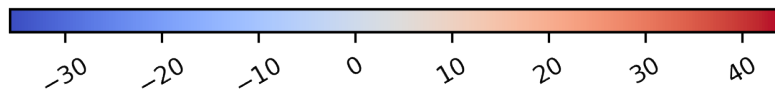
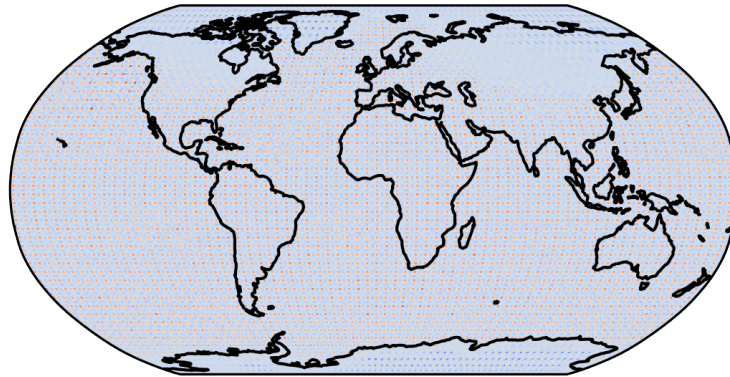
```
[21]: # plot of z-score under null hypothesis that "orig" value= "zfpA1.0" value
ldcpy.plot(
    col_ts,
    "TS",
    sets=["orig", "zfpA1.0", "zfpA1e-1"],
    calc="zscore",
    calc_type="metric_of_diff",
    vert_plot=True,
)
```



orig & zfpA1.0: TS: zscore: cutoff 2.34, % sig: 53.23 metric\_of\_diff



orig & zfpA1e-1: TS: zscore: cutoff 2.40, % sig: 40.96 metric\_of\_diff



## Time-Series Plots

We may want to aggregate the data spatially and look for trends over time. Therefore, we can also create a time-series plot of the metrics by changing the `plot_type` to “time\_series”. For time-series plots, the metric values are on the y-axis and the x-axis represents time. We are also able to group the data by time, as shown below.

### Basic Time-Series Plot

In the example below, we look at the ‘orig’ TS data in collection `col_ts`, and display the spatial mean at each day of the year (our data consists of 100 days).

```
[22]: # Time-series plot of TS mean in ds orig dataset
ldcpy.plot(
    col_ts,
    "TS",
    sets=["orig"],
    calc="mean",
```

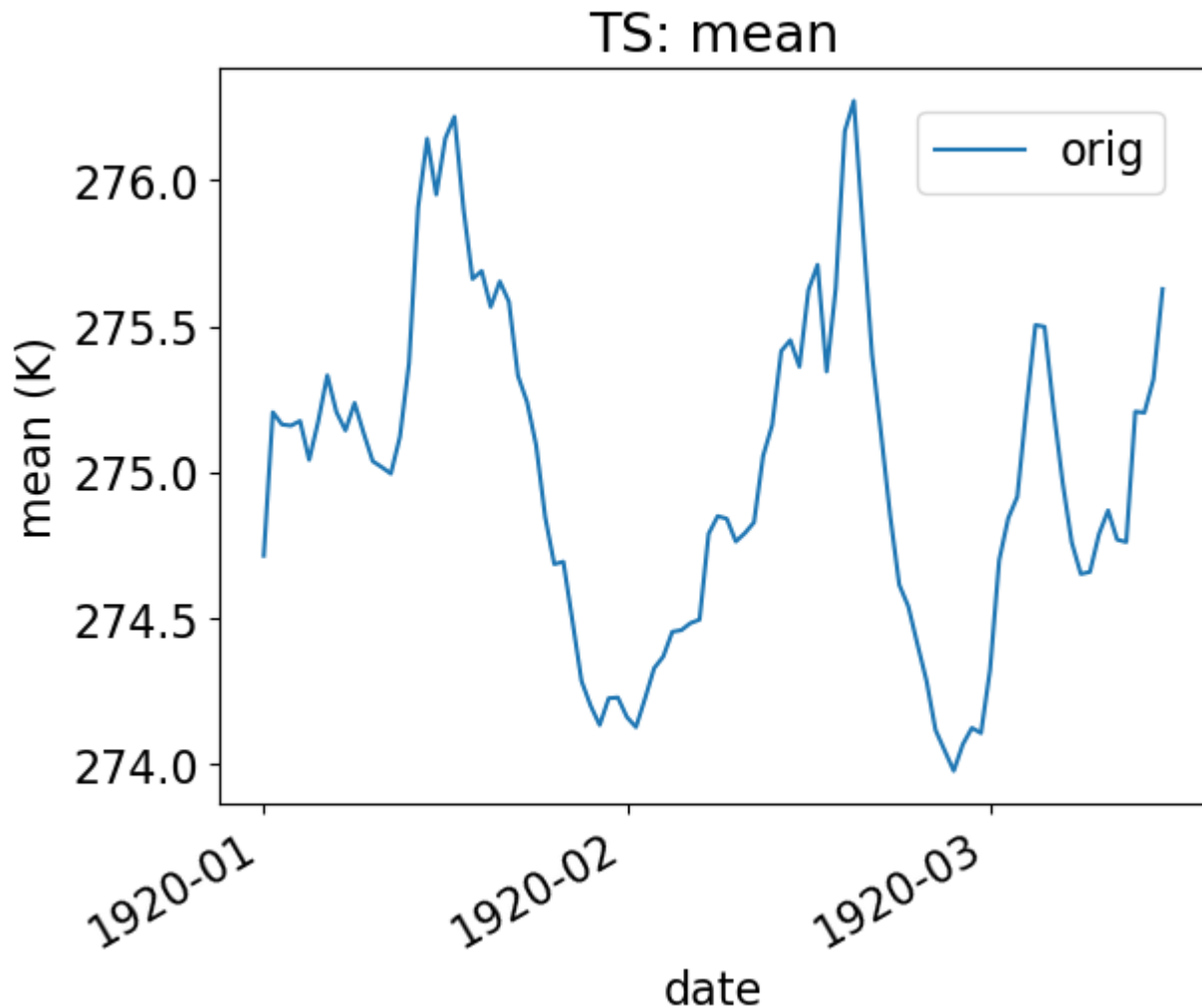
(continues on next page)

(continued from previous page)

```

plot_type="time_series",
vert_plot=True,
legend_loc="best",
)

```



### Using the group\_by keyword

To group the data by time, use the “group\_by” keyword. This plot shows the mean standard deviation over all latitude and longitude points for each month. Note that this plot is not too interesting for our sample data, which has only 100 days of data.

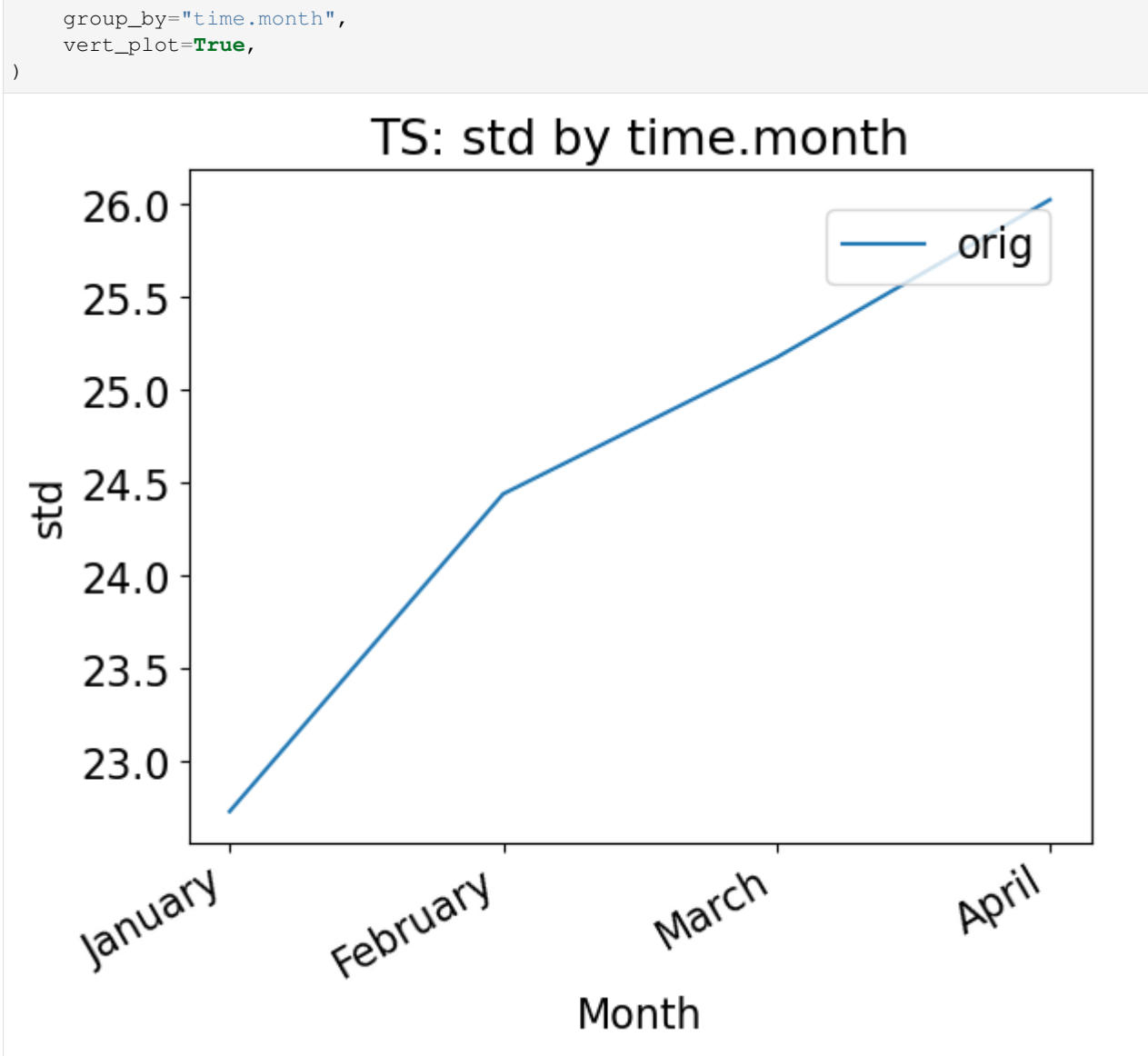
```

[23]: # Time-series plot of TS standard deviation in col_ds "orig" dataset, grouped by month
ldcpy.plot(
    col_ts,
    "TS",
    sets=["orig"],
    calc="std",
    plot_type="time_series",

```

(continues on next page)

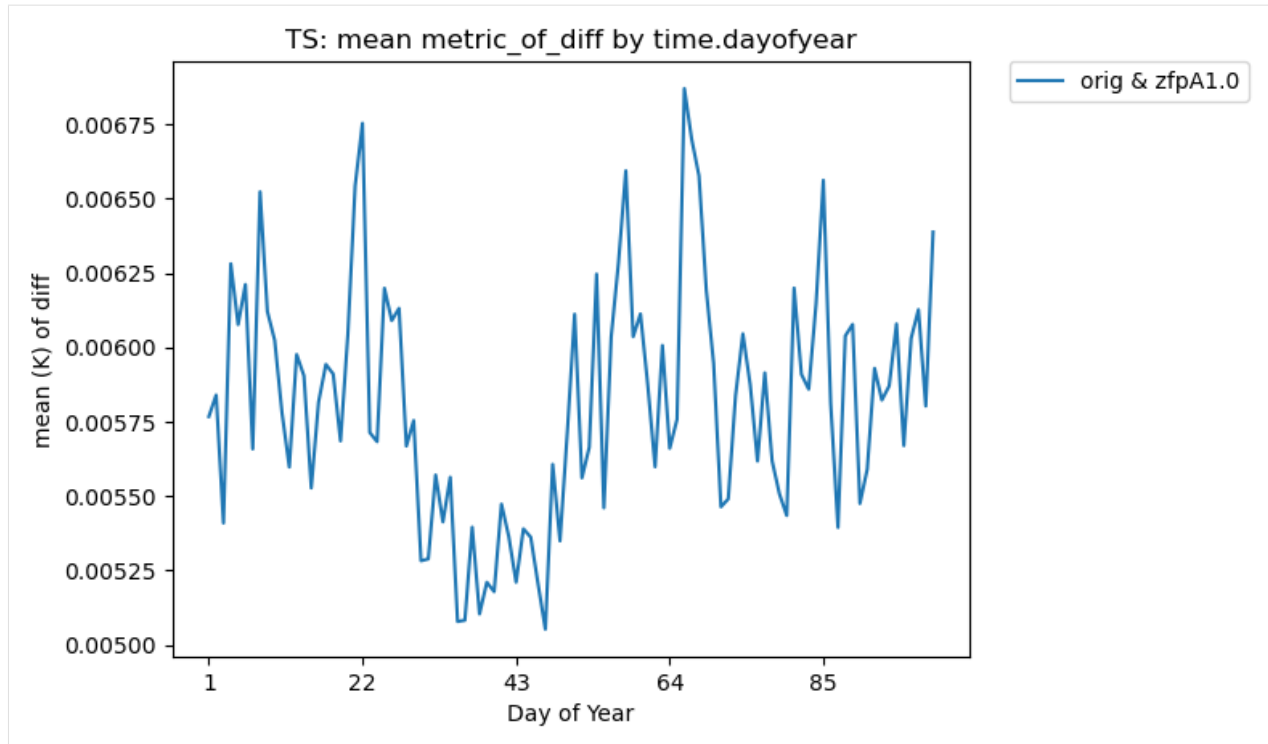
(continued from previous page)



One could also group by days of the year, as below. Again, because we have less than a year of data, this plot looks the same as the previous version. However, this approach would be useful with data containing multiple years.

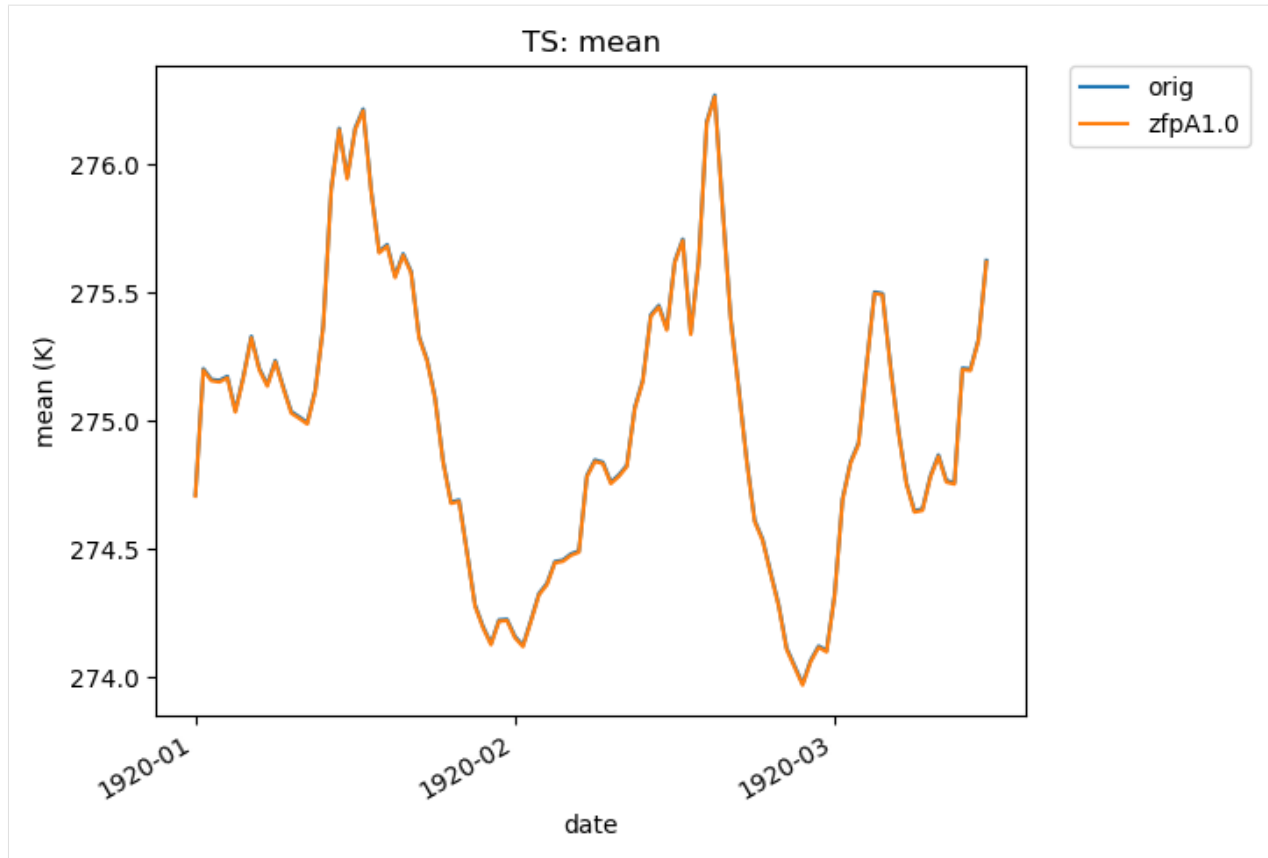
```
[24]: # Time-series plot of TS mean in col_ts "orig" dataset, grouped by day of year
ldcpy.plot(
    col_ts,
    "TS",
    sets=["orig", "zfpA1.0"],
    calc="mean",
    calc_type="metric_of_diff",
    plot_type="time_series",
    group_by="time.dayofyear",
)
```





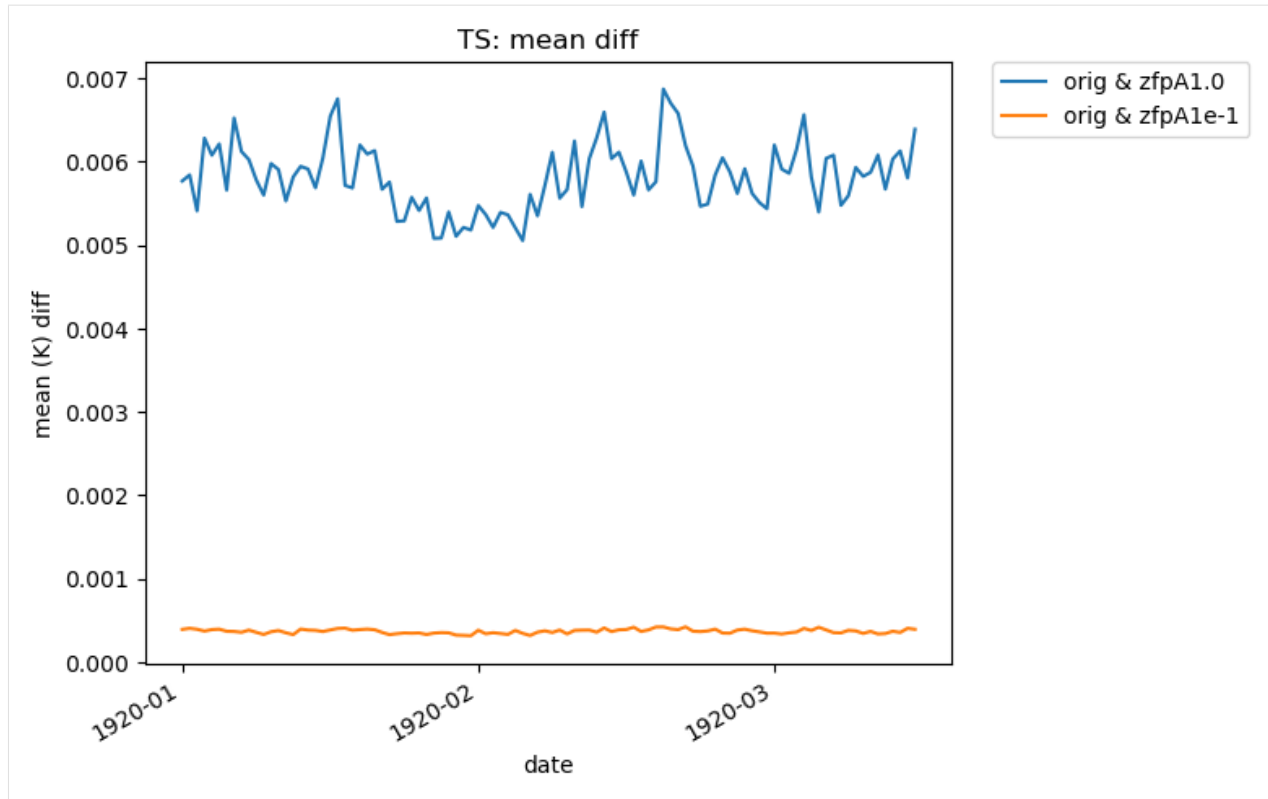
We can also overlay multiple sets of time-series data on the same plot. For instance, we can plot the mean of two datasets over time. Note that the blue and orange lines overlap almost perfectly:

```
[25]: # Time-series plot of TS mean in col_ts 'orig' dataset
ldcpy.plot(
    col_ts,
    "TS",
    sets=["orig", "zfpA1.0"],
    calc="mean",
    plot_type="time_series",
)
```



If we change the `calc_type` to “diff”, “ratio” or “metric\_of\_diff”, the first element in `sets` is compared against subsequent elements in the `sets` list. For example, we can compare the difference in the mean of two compressed datasets to the original dataset like so:

```
[26]: # Time-series plot of TS mean differences (with 'orig') in col_ts orig dataset
ldcpy.plot(
    col_ts,
    "TS",
    sets=["orig", "zfpA1.0", "zfpA1e-1"],
    calc="mean",
    plot_type="time_series",
    calc_type="diff",
)
```

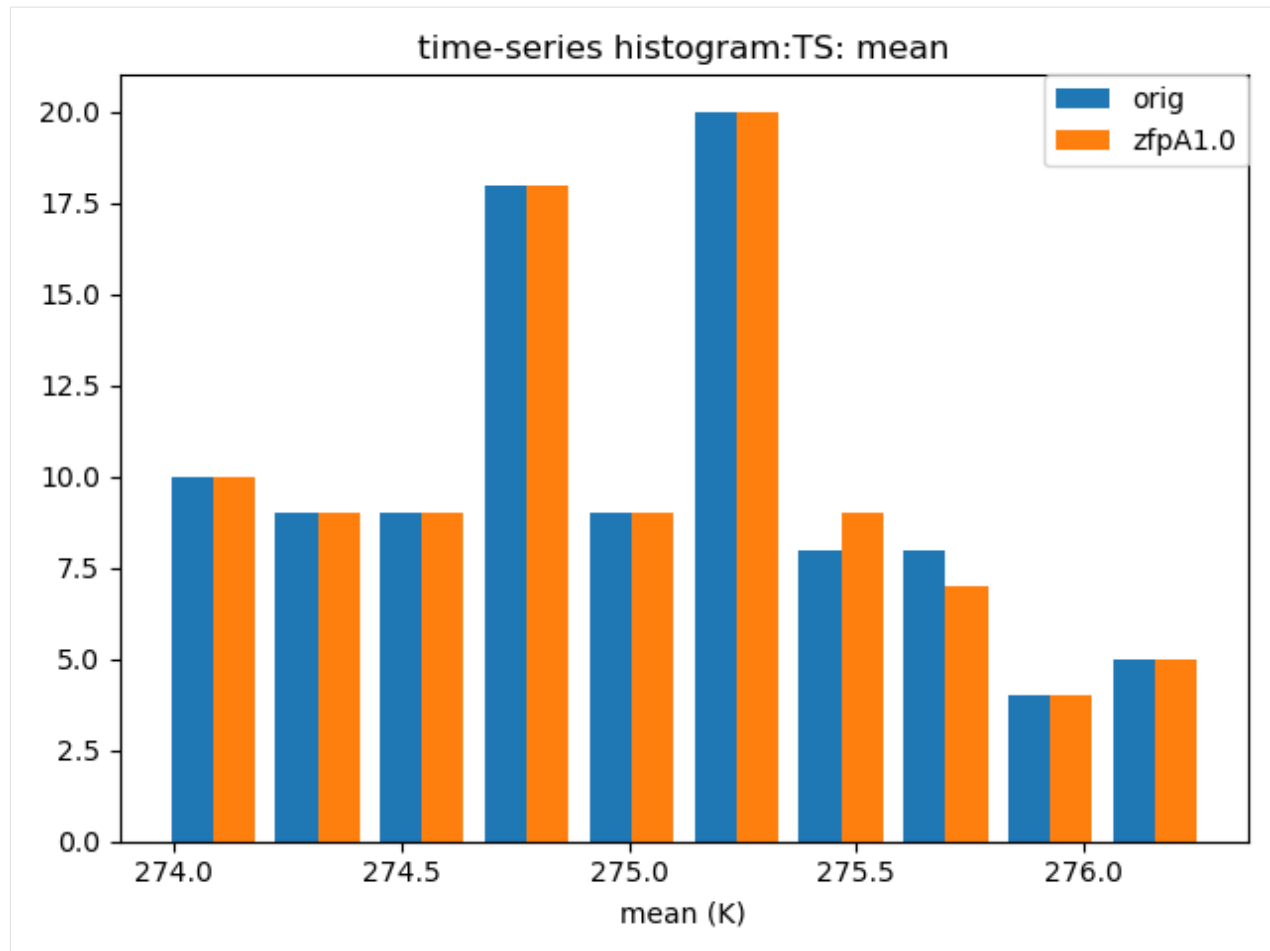


## Histograms

We can also create a histogram of the data by changing the `plot_type` to 'histogram'. Note that these histograms are currently generated from time-series metric values (a histogram of spatial values is not currently available).

The histogram below shows the mean values of TS in the 'orig' and 'zfpA1.0' datasets in our collection `col_ts`. Recall that this dataset contains 100 timeslices.

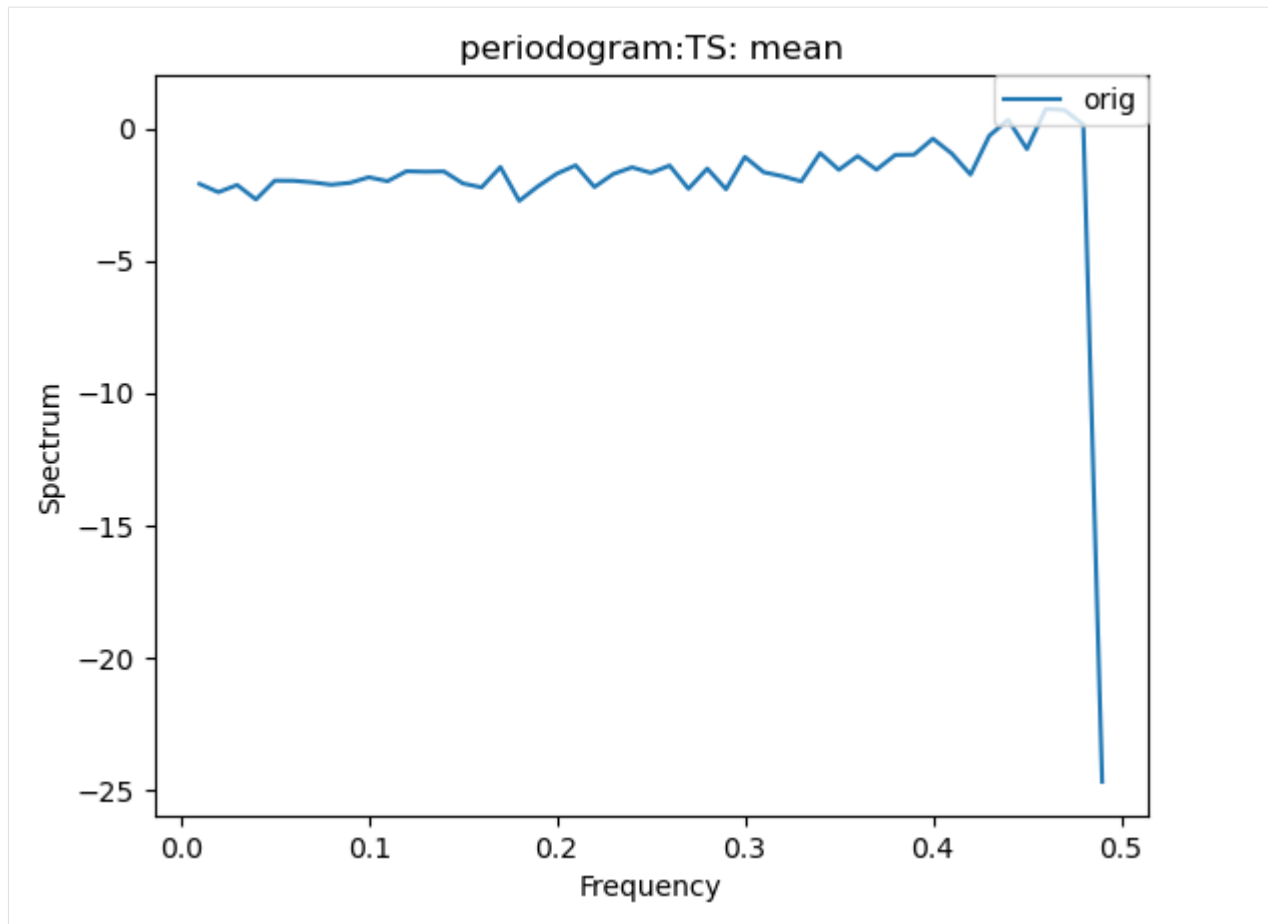
```
[27]: # Histogram of mean TS values in the col_ts (for 'orig' and 'zfpA1.0') dataset
ldcpy.plot(col_ts, "TS", sets=["orig", "zfpA1.0"], calc="mean", plot_type="histogram")
```



### Other Time-Series Plots

We can create a periodogram based on a dataset as well, by specifying a `plot_type` of “periodogram”.

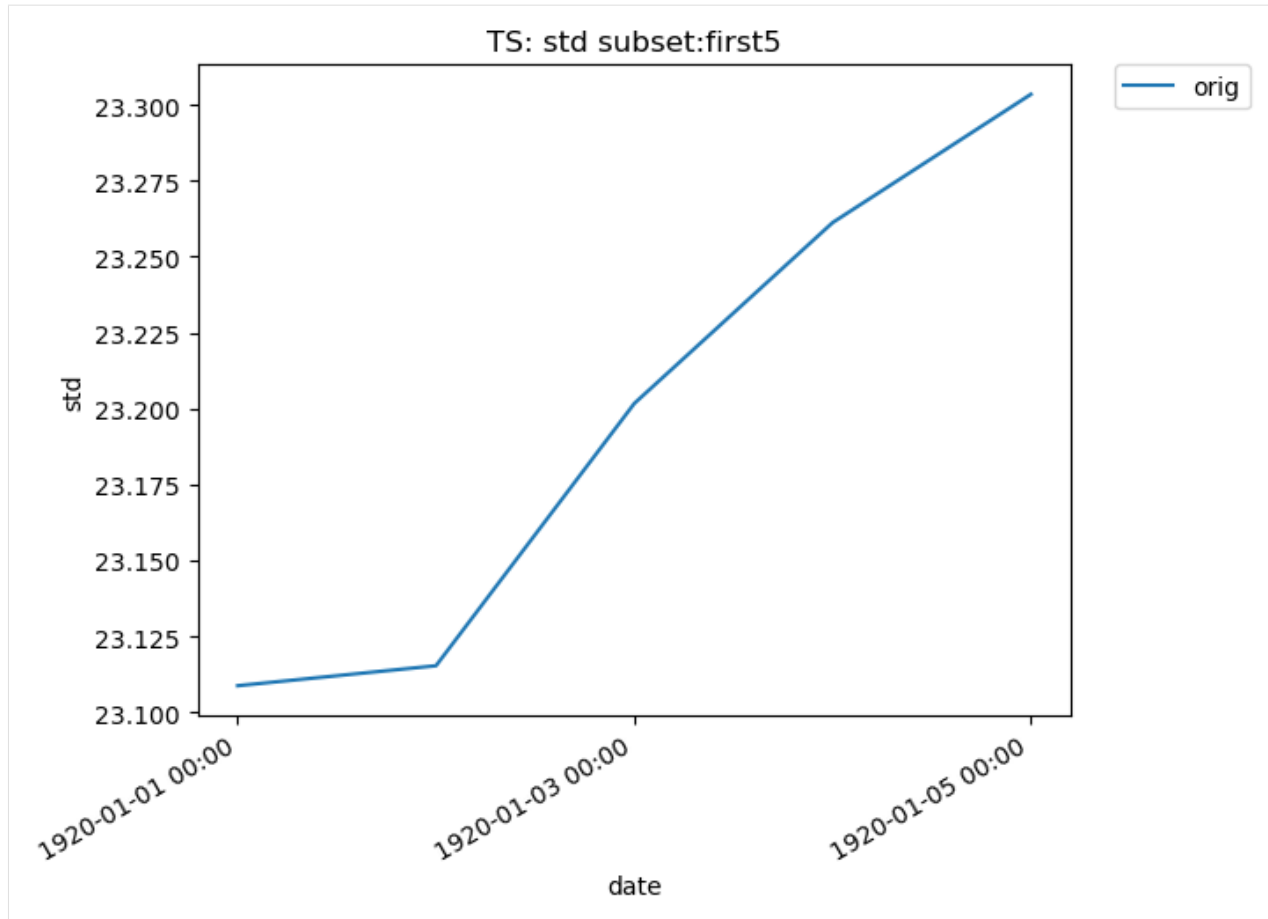
```
[28]: ldcpy.plot(col_ts, "TS", sets=["orig"], calc="mean", plot_type="periodogram")
```



## Subsetting

Subsetting is also possible on time-series data. The following plot makes use of the subset argument, which is used to plot metrics on only a portion of the data. A full list of available subsets is available [here](#). The following plot uses the 'first5' subset, which only shows the metric values for the first five time slices (in this case, days) of data:

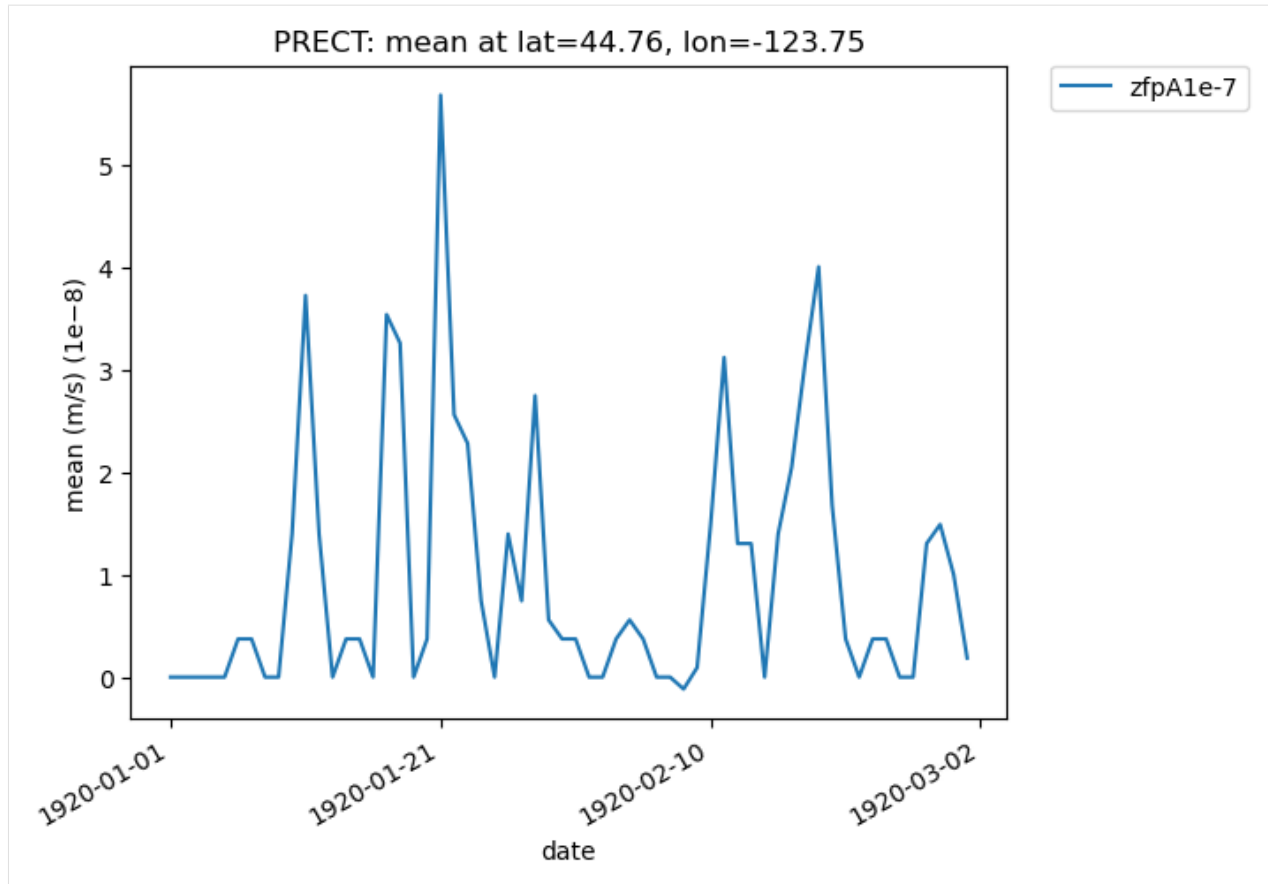
```
[29]: # Time-series plot of first five TS standard deviations in col_ts "orig" dataset
ldcpy.plot(
    col_ts,
    "TS",
    sets=["orig"],
    calc="std",
    plot_type="time_series",
    subset="first5",
)
```



Additionally, we can specify “lat” and “lon” keywords for time-series plots that give us a subset of the data at a single point, rather than averaging over all latitudes and longitudes. The nearest latitude and longitude point to the one specified is plotted (and the actual coordinates of the point can be found in the plot title). This plot, for example, shows the difference in mean rainfall between the compressed and original data at the location (44.76, -123.75):

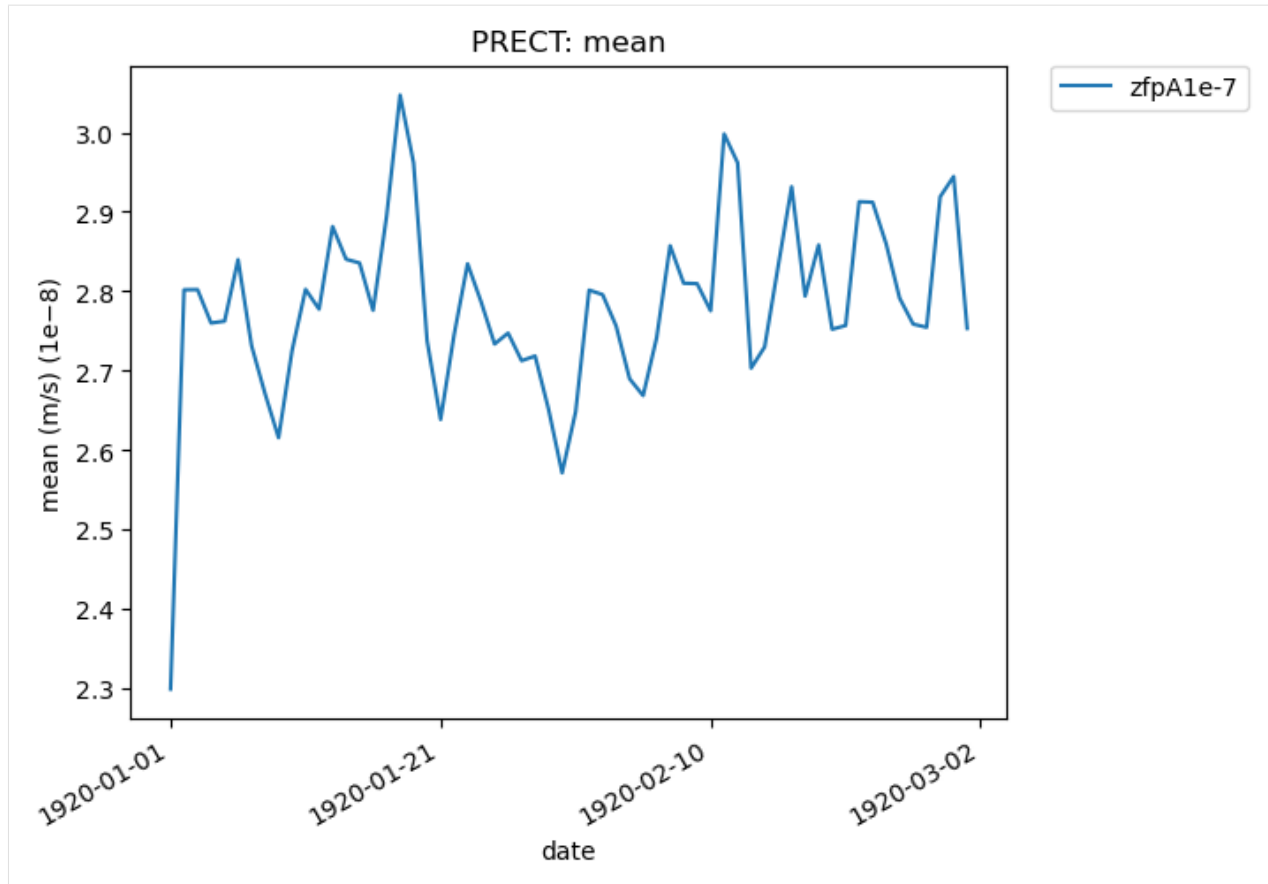
```
[30]: # Time series plot of PRECT mean data for col_prect "zfpA1e-7" dataset at the
      ↪ location (44.76, -123.75)

ldcpy.plot(
    col_prect,
    "PRECT",
    sets=["zfpA1e-7"],
    calc="mean",
    plot_type="time_series",
    lat=44.76,
    lon=-123.75,
)
```



```
[31]: # Time series plot of PRECt mean data for col_prec "zfpAle-7" dataset at the
      ↪ location (44.76, -123.75)

ldcpy.plot(col_prec, "PRECt", sets=["zfpAle-7"], calc="mean", plot_type="time_series
      ↪")
```

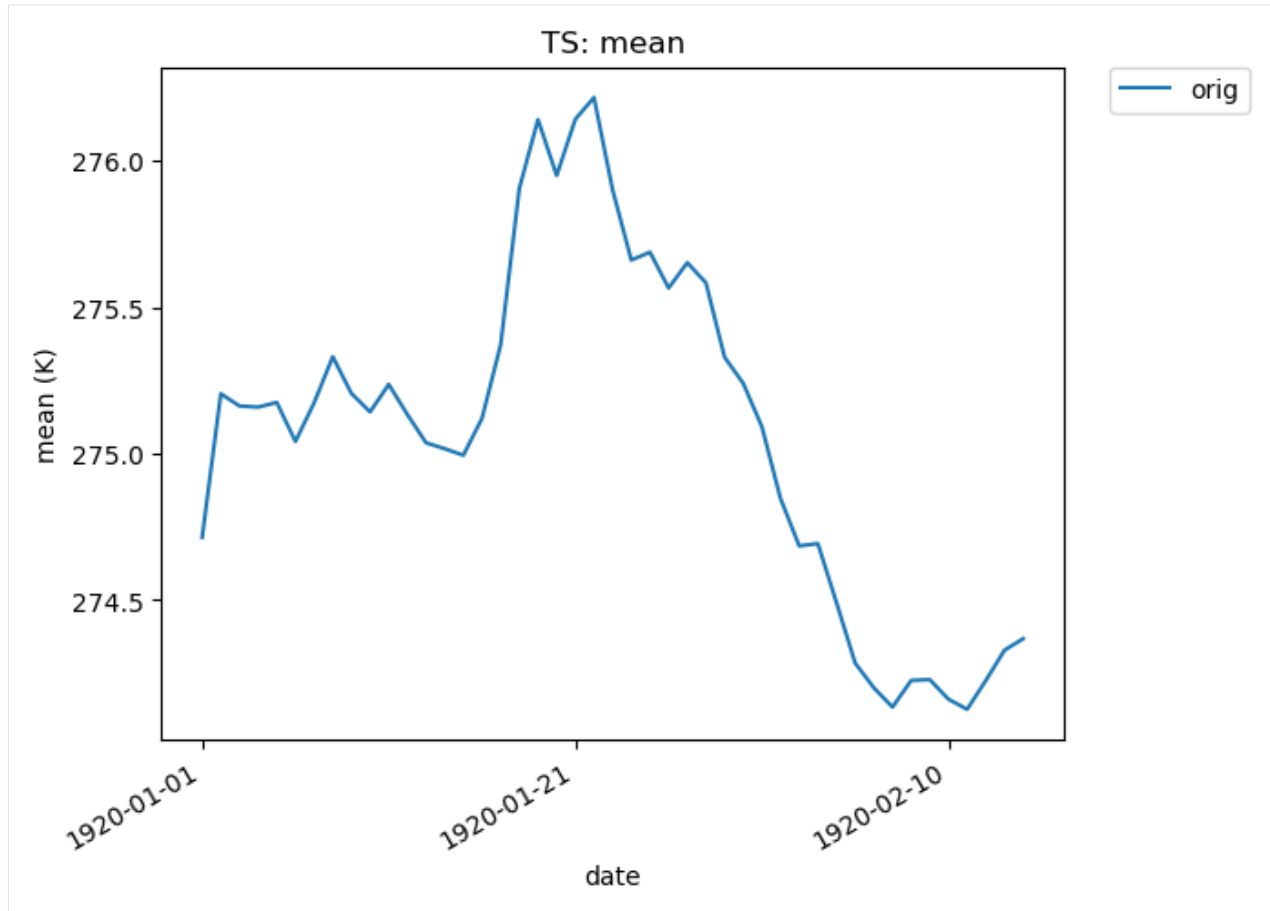


It is also possible to plot metrics for a subset of the data, by specifying the start and end indices of the data we are interested in. This command creates a time-series plot of the mean TS values for the first 45 days of data:

```
[32]: # Time series plot of first 45 TS mean data points for col_ts "orig" dataset

ldcpy.plot(
    col_ts,
    "TS",
    sets=["orig"],
    calc="mean",
    start=0,
    end=44,
    plot_type="time_series",
)
```





```
[ ]:
```

### 5.4.2 Using data from AWS

A significant amount of Earth System Model (ESM) data is publicly available online, including data from the CESM Large Ensemble, CMIP5, and CMIP6 datasets. For accessing a single file, we can specify the file (typically netcdf or zarr format) and its location and then use `fsspec` (the “Filesystem Spec+ python package) and `xarray` to create a `array.dataset`. For several files, the `intake_esm` python module (<https://github.com/intake/intake-esm>) is particularly nice for obtaining the data and put it into an `xarray.dataset`.

This notebook assumes familiarity with the Tutorial Notebook. It additionally shows how to gather data from an ESM collection, put it into a dataset, and then create simple plots using the data with `ldcpy`.

## Example Data

The example data we use is from the CESM Large Ensemble, member 31. This ensemble data has been lossily compressed and reconstructed as part of a blind evaluation study of lossy data compression in LENS (e.g., <http://www.cesm.ucar.edu/projects/community-projects/LENS/projects/lossy-data-compression.html> or <https://gmd.copernicus.org/articles/9/4381/2016/>).

Most of the data from the CESM Large Ensemble Project has been made available on Amazon Web Services (Amazon S3), see [http://ncar-aws-www.s3-website-us-west-2.amazonaws.com/CESM\\_LENS\\_on\\_AWS.htm](http://ncar-aws-www.s3-website-us-west-2.amazonaws.com/CESM_LENS_on_AWS.htm).

For comparison purposes, the original (non-compressed) data for Ensemble 31 has recently been made available on Amazon Web Services (Amazon S3) in the “ncar-cesm-lens-baker-lossy-compression-test” bucket.

```
[1]: # Add ldcpy root to system path
import sys

sys.path.insert(0, '../..../')

# Import ldcpy package
# Autoreloads package everytime the package is called, so changes to code will be_
↪reflected in the notebook if the above sys.path.insert(...) line is uncommented.
%load_ext autoreload
%autoreload 2
import fsspec
import intake
import xarray as xr

import ldcpy

# display the plots in this notebook
%matplotlib inline

# silence warnings
import warnings

warnings.filterwarnings("ignore")
```

## Method 1: using fsspec and xr.open\_zarr

First, specify the filesystem and location of the data. Here we are accessing the original data from CESM-LENS ensemble 31, which is available on Amazon S3 in the store named “ncar-cesm-lens-baker-lossy-compression-test” bucket.

First we listing all available files (which are timeseries files containing a single variable) for that dataset. Note that unlike in the TutorialNotebook (which used NetCDF files), these files are all zarr format. Both monthly and daily data is available.

```
[2]: fs = fsspec.filesystem("s3", anon=True)
stores = fs.ls("ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/") [1:]
stores[:]

[2]: ['ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-daily-
↪FLNS.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-daily-
↪FLNSC.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-daily-
↪FLUT.zarr',
```

(continues on next page)

(continued from previous page)

```

'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-daily-
↪FSNS.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-daily-
↪FSNSC.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-daily-
↪FSNTOA.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-daily-
↪ICEFRAC.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-daily-
↪LHFLX.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-daily-
↪PRECL.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-daily-
↪PRECSC.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-daily-
↪PRECSL.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-daily-
↪PRECT.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-daily-
↪PRECTMX.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-daily-
↪PSL.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-daily-
↪Q850.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-daily-
↪SHFLX.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-daily-
↪TMQ.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-daily-
↪TREFHT.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-daily-
↪TREFHTMN.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-daily-
↪TREFHTMX.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-daily-
↪TS.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-daily-
↪UBOT.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-daily-
↪WSPDSRFV.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-daily-
↪Z500.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-monthly-
↪FLNS.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-monthly-
↪FLNSC.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-monthly-
↪FLUT.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-monthly-
↪FSNS.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-monthly-
↪FSNSC.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-monthly-
↪FSNTOA.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-monthly-
↪ICEFRAC.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-monthly-
↪LHFLX.zarr',

```

(continues on next page)

(continued from previous page)

```

'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-monthly-
↪PRECC.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-monthly-
↪PRECL.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-monthly-
↪PRECSC.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-monthly-
↪PRECSL.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-monthly-
↪PSL.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-monthly-
↪Q.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-monthly-
↪SHFLX.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-monthly-
↪T.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-monthly-
↪TMQ.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-monthly-
↪TREFHT.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-monthly-
↪TREFHTMN.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-monthly-
↪TREFHTMX.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-monthly-
↪TS.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-monthly-
↪U.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-monthly-
↪V.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-20C-monthly-
↪Z3.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-daily-
↪FLNS.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-daily-
↪FLNSC.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-daily-
↪FLUT.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-daily-
↪FSNS.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-daily-
↪FSNSC.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-daily-
↪FSNTOA.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-daily-
↪ICEFRAC.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-daily-
↪LHFLX.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-daily-
↪PRECL.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-daily-
↪PRECSC.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-daily-
↪PRECSL.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-daily-
↪PRECT.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-daily-
↪PRECTMX.zarr',

```

(continues on next page)

(continued from previous page)

```

'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-daily-
↪PSL.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-daily-
↪Q850.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-daily-
↪SHFLX.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-daily-
↪TMQ.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-daily-
↪TREFHT.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-daily-
↪TREFHTMN.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-daily-
↪TREFHTMX.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-daily-
↪TS.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-daily-
↪UBOT.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-daily-
↪WSPDSRFAV.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-daily-
↪Z500.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-
↪monthly-FLNS.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-
↪monthly-FLNSC.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-
↪monthly-FLUT.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-
↪monthly-FSNS.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-
↪monthly-FSNSC.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-
↪monthly-FSNTOA.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-
↪monthly-ICEFRAC.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-
↪monthly-LHFLX.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-
↪monthly-PRECC.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-
↪monthly-PRECL.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-
↪monthly-PRECSC.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-
↪monthly-PRECSL.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-
↪monthly-PSL.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-
↪monthly-Q.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-
↪monthly-SHFLX.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-
↪monthly-T.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-
↪monthly-TMQ.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-
↪monthly-TREFHT.zarr',

```

(continues on next page)

(continued from previous page)

```
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-
↪monthly-TREFHTMN.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-
↪monthly-TREFHTMX.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-
↪monthly-TS.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-
↪monthly-U.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-
↪monthly-V.zarr',
'ncar-cesm-lens-baker-lossy-compression-test/lens-ens31/cesmle-atm-ens31-RCP85-
↪monthly-Z3.zarr']
```

Then we select the file from the store that we want and open it as an `xarray.Dataset` using `xr.open_zarr()`. Here we grab data for the first 2D daily variable, FLNS (net longwave flux at surface, in  $W/m^2$ ), in the list (accessed by its location – `stores[0]`).

```
[3]: store = fs.get_mapper(stores[0])
ds_flns = xr.open_zarr(store, consolidated=True)
ds_flns

[3]: <xarray.Dataset>
Dimensions:      (lat: 192, lon: 288, nbnd: 2, time: 31390)
Coordinates:
  * lat          (lat) float64 -90.0 -89.06 -88.12 -87.17 ... 88.12 89.06 90.0
  * lon          (lon) float64 0.0 1.25 2.5 3.75 5.0 ... 355.0 356.2 357.5 358.8
  * time         (time) object 1920-01-01 12:00:00 ... 2005-12-31 12:00:00
    time_bnds    (time, nbnd) object dask.array<chunksize=(15695, 2), meta=np.ndarray>
Dimensions without coordinates: nbnd
Data variables:
  FLNS          (time, lat, lon) float32 dask.array<chunksize=(576, 192, 288), meta=np.
↪ndarray>
Attributes:
  Conventions:      CF-1.0
  NCO:              netCDF Operators version 4.7.9 (Homepage = http://nco.s...
  Version:          $Name$
  case:             b.e11.B20TRC5CNBDRD.f09_g16.031
  host:             ys0219
  initial_file:     b.e11.B20TRC5CNBDRD.f09_g16.001.cam.i.1920-01-01-00000.nc
  logname:          mickelso
  revision_Id:      $Id$
  source:           CAM
  title:            UNSET
  topography_file:  /glade/p/cesmdata/cseg/inputdata/atm/cam/topo/USGS-gtop...
```

The above returned an `xarray.Dataset`.

Now let's grab the TMQ (Total vertically integrated precipitable water) and the TS (surface temperature data) and PRECT (precipitation rate) data from AWS.

```
[4]: # TMQ data
store2 = fs.get_mapper(stores[16])
ds_tmq = xr.open_zarr(store2, consolidated=True)
# TS data
store3 = fs.get_mapper(stores[20])
ds_ts = xr.open_zarr(store3, consolidated=True)
# PRECT data
```

(continues on next page)

(continued from previous page)

```
store4 = fs.get_mapper(stores[11])
ds_prect = xr.open_zarr(store4, consolidated=True)
```

Now we have the original data for FLNS and TMQ and TS and PRECT. Next we want to get the lossy compressed variants to compare with these.

## Method 2: Using intake\_esm

Now we will demonstrate using the intake\_esm module to get the lossy variants of the files retrieved above. We can use the intake\_esm module to search for and open several files as xarray.Dataset objects. The code below is modified from the intake\_esm documentation, available here: <https://intake-esm.readthedocs.io/en/latest/?badge=latest#overview>.

We want to use ensemble 31 data from the CESM-LENS collection on AWS, which (as explained above) has been subjected to lossy compression. Many catalogs for publicly available datasets are accessible via intake-esm can be found at <https://github.com/NCAR/intake-esm-datastore/tree/master/catalogs>, including for CESM-LENS. We can open that collection as follows (see here: <https://github.com/NCAR/esm-collection-spec/blob/master/collection-spec/collection-spec.md#attribute-object>):

```
[5]: aws_loc = (
      "https://raw.githubusercontent.com/NCAR/cesm-lens-aws/master/intake-catalogs/aws-
      ↪cesml-le.json"
    )
aws_col = intake.open_esm_datastore(aws_loc)
aws_col

<IPython.core.display.HTML object>
```

Next, we search for the subset of the collection (dataset and variables) that we are interested in. Let's grab FLNS, TMQ, and TS daily data from the atm component for our comparison (available data in this collection is listed here: [http://ncar-aws-www.s3-website-us-west-2.amazonaws.com/CESM\\_LENS\\_on\\_AWS.htm](http://ncar-aws-www.s3-website-us-west-2.amazonaws.com/CESM_LENS_on_AWS.htm)).

```
[6]: # we want daily data for FLNS, TMQ, and TS and PRECT
aws_col_subset = aws_col.search(
    component="atm",
    frequency="daily",
    experiment="20C",
    variable=["FLNS", "TS", "TMQ", "PRECT"],
)
# display header info to verify that we got the right variables
aws_col_subset.df.head()
```

```
[6]:  component frequency experiment variable \
0      atm      daily      20C      FLNS
1      atm      daily      20C     PRECT
2      atm      daily      20C      TMQ
3      atm      daily      20C       TS

                                     path \
0  s3://ncar-cesm-lens/atm/daily/cesmLE-20C-FLNS...
1  s3://ncar-cesm-lens/atm/daily/cesmLE-20C-PRECT...
2  s3://ncar-cesm-lens/atm/daily/cesmLE-20C-TMQ.zarr
3  s3://ncar-cesm-lens/atm/daily/cesmLE-20C-TS.zarr

                                variable_long_name  dim_per_tstep \
0                                net longwave flux at surface         2.0
1  total (convective and large-scale) precipitati...         2.0
```

(continues on next page)

(continued from previous page)

2	total (vertically integrated) precipitable water	2.0
3	surface temperature (radiative)	2.0

	start	end
0	1920-01-01 12:00:00	2005-12-31 12:00:00
1	1920-01-01 12:00:00	2005-12-31 12:00:00
2	1920-01-01 12:00:00	2005-12-31 12:00:00
3	1920-01-01 12:00:00	2005-12-31 12:00:00

Then we load matching catalog entries into xarray datasets ([https://intake-esm.readthedocs.io/en/latest/api.html#intake\\_esm.core.esm\\_datastore.to\\_dataset\\_dict](https://intake-esm.readthedocs.io/en/latest/api.html#intake_esm.core.esm_datastore.to_dataset_dict)). We create a dictionary of datasets:

```
[7]: dset_dict = aws_col_subset.to_dataset_dict(
      zarr_kwargs={"consolidated": True, "decode_times": True},
      storage_options={"anon": True},
      cdf_kwargs={"chunks": {}, "decode_times": False},
    )
dset_dict
```

--> The keys in the returned dictionary of datasets are constructed as follows:  
'component.experiment.frequency'

<IPython.core.display.HTML object>

```
[7]: {'atm.20C.daily': <xarray.Dataset>
  Dimensions:      (lat: 192, lon: 288, member_id: 40, nbnd: 2, time: 31390)
  Coordinates:
    * lat          (lat) float64 -90.0 -89.06 -88.12 -87.17 ... 88.12 89.06 90.0
    * lon          (lon) float64 0.0 1.25 2.5 3.75 5.0 ... 355.0 356.2 357.5 358.8
    * member_id    (member_id) int64 1 2 3 4 5 6 7 8 ... 34 35 101 102 103 104 105
    * time         (time) object 1920-01-01 12:00:00 ... 2005-12-31 12:00:00
    time_bnds      (time, nbnd) object dask.array<chunksize=(15695, 2), meta=np.ndarray>
  Dimensions without coordinates: nbnd
  Data variables:
    FLNS           (member_id, time, lat, lon) float32 dask.array<chunksize=(1, 576, 192,
    ↪ 288), meta=np.ndarray>
    PRECT          (member_id, time, lat, lon) float32 dask.array<chunksize=(1, 576, 192,
    ↪ 288), meta=np.ndarray>
    TMQ            (member_id, time, lat, lon) float32 dask.array<chunksize=(1, 576, 192,
    ↪ 288), meta=np.ndarray>
    TS             (member_id, time, lat, lon) float32 dask.array<chunksize=(1, 576, 192,
    ↪ 288), meta=np.ndarray>
  Attributes:
    nco_openmp_thread_number: 1
    revision_Id:              $Id$
    logname:                  mudryk
    initial_file:             b.e11.B20TRC5CNBDRD.f09_g16.001.cam.i.1920-01-...
    title:                    UNSET
    NCO:                      4.4.2
    Version:                  $Name$
    source:                   CAM
    Conventions:              CF-1.0
    important_note:           This data is part of the project 'Blind Evalua...
    topography_file:          /scratch/p/pjk/mudryk/cesm1_1_2_LENS/inputdata...
    intake_esm_dataset_key:   atm.20C.daily}
```

Check the dataset keys to ensure that what we want is present. Here we only have one entry in the dictionary as we



requested the same time period and output frequency for all variables:

```
[8]: dset_dict.keys()
[8]: dict_keys(['atm.20C.daily'])
```

Finally, put the dataset that we are interested from the dictionary into its own dataset variable. (We want the 20th century daily data – which is our only option.)

Also note from above that there are 40 ensemble members - and we just want ensemble 31 (member\_id = 30 as can be seen in the coordinates above).

```
[9]: aws_ds = dset_dict["atm.20C.daily"]
aws_ds = aws_ds.isel(member_id=30)
aws_ds

[9]: <xarray.Dataset>
Dimensions:    (lat: 192, lon: 288, nbnd: 2, time: 31390)
Coordinates:
  * lat        (lat) float64 -90.0 -89.06 -88.12 -87.17 ... 88.12 89.06 90.0
  * lon        (lon) float64 0.0 1.25 2.5 3.75 5.0 ... 355.0 356.2 357.5 358.8
  member_id    int64 31
  * time       (time) object 1920-01-01 12:00:00 ... 2005-12-31 12:00:00
  time_bnds    (time, nbnd) object dask.array<chunksize=(15695, 2), meta=np.ndarray>
Dimensions without coordinates: nbnd
Data variables:
  FLNS         (time, lat, lon) float32 dask.array<chunksize=(576, 192, 288), meta=np.
↪ndarray>
  PRECT        (time, lat, lon) float32 dask.array<chunksize=(576, 192, 288), meta=np.
↪ndarray>
  TMQ          (time, lat, lon) float32 dask.array<chunksize=(576, 192, 288), meta=np.
↪ndarray>
  TS           (time, lat, lon) float32 dask.array<chunksize=(576, 192, 288), meta=np.
↪ndarray>
Attributes:
  nco_openmp_thread_number: 1
  revision_Id:              $Id$
  logname:                  mudryk
  initial_file:             b.e11.B20TRC5CNBDRD.f09_g16.001.cam.i.1920-01-...
  title:                    UNSET
  NCO:                      4.4.2
  Version:                  $Name$
  source:                   CAM
  Conventions:              CF-1.0
  important_note:           This data is part of the project 'Blind Evalua...
  topography_file:          /scratch/p/pjk/mudryk/cesm1_1_2_LENS/inputdata...
  intake_esm_dataset_key:   atm.20C.daily
```

Now we have datasets for the original and the lossy compressed data for FLNS, TMQ, PRECT, and TS, which we can extract into a dataset for each variable:

```
[10]: # extract the three variables from aws_ds as datasets
aws_flns = aws_ds["FLNS"].to_dataset()
aws_tmq = aws_ds["TMQ"].to_dataset()
aws_ts = aws_ds["TS"].to_dataset()
aws_prect = aws_ds["PRECT"].to_dataset()
```

### 5.4.3 Use *ldcpy* to compare the original data to the lossy compressed data

To use *ldcpy*, we need to group the data that we want to compare (like variables) into dataset collections. In the Tutorial Notebook, we used `ldcpy.open_datasets()` to do this as we needed to get the data from the NetCDF files. Here we already loaded the data from AWS into datasets, so we just need to use `ldcpy.collect_datasets()` to form collections of the datasets that we want to compare.

`ldcpy.collect_datasets()` requires the following three arguments:

- *varnames* : the variable(s) of interest to combine across files (typically the timeseries file variable name)
- *list\_of\_ds* : a list of the xarray datasets
- *labels* : a corresponding list of names (or labels) for each dataset in the collection

Note: This function is a wrapper for `xarray.concat()`, and any additional key/value pairs passed in as a dictionary are used as arguments to `xarray.concat()`.

We will create 4 collections for *ldcpy* (one each for FLNS, TMQ, PRECT, and TS) and assign labels “original” and “lossy” to the respective datasets.

```
[11]: # FLNS collection
col_flns = ldcpy.collect_datasets(["FLNS"], [ds_flns, aws_flns], ["original", "lossy"])
col_flns
```

dataset size in GB 27.63

```
[11]: <xarray.Dataset>
Dimensions:      (collection: 2, lat: 382, lon: 288, time: 31390)
Coordinates:
  member_id      int64 31
  * lat           (lat) float64 -90.0 -89.06 -89.06 -88.12 ... 89.06 89.06 90.0
  * lon           (lon) float64 0.0 1.25 2.5 3.75 5.0 ... 355.0 356.2 357.5 358.8
  * time          (time) object 1920-01-01 12:00:00 ... 2005-12-31 12:00:00
  * collection    (collection) <U8 'original' 'lossy'
Data variables:
  FLNS            (collection, time, lat, lon) float32 dask.array<chunksize=(1, 576, 382, 288), meta=np.ndarray>
Attributes:
  Conventions:      CF-1.0
  NCO:              netCDF Operators version 4.7.9 (Homepage = http://nco.s...
  Version:          $Name$
  case:             b.e11.B20TRC5CNBDRD.f09_g16.031
  host:             ys0219
  initial_file:     b.e11.B20TRC5CNBDRD.f09_g16.001.cam.i.1920-01-01-00000.nc
  logname:          mickelso
  revision_Id:      $Id$
  source:           CAM
  title:            UNSET
  topography_file:  /glade/p/cesmdata/cseg/inputdata/atm/cam/topo/USGS-gtop...
```

```
[12]: # TMQ collection
col_tmq = ldcpy.collect_datasets(["TMQ"], [ds_tmq, aws_tmq], ["original", "lossy"])
col_tmq
```

dataset size in GB 27.63

```
[13]: # Ts collection
col_ts = ldcpy.collect_datasets(["TS"], [ds_ts, aws_ts], ["original", "lossy"])
col_ts
```

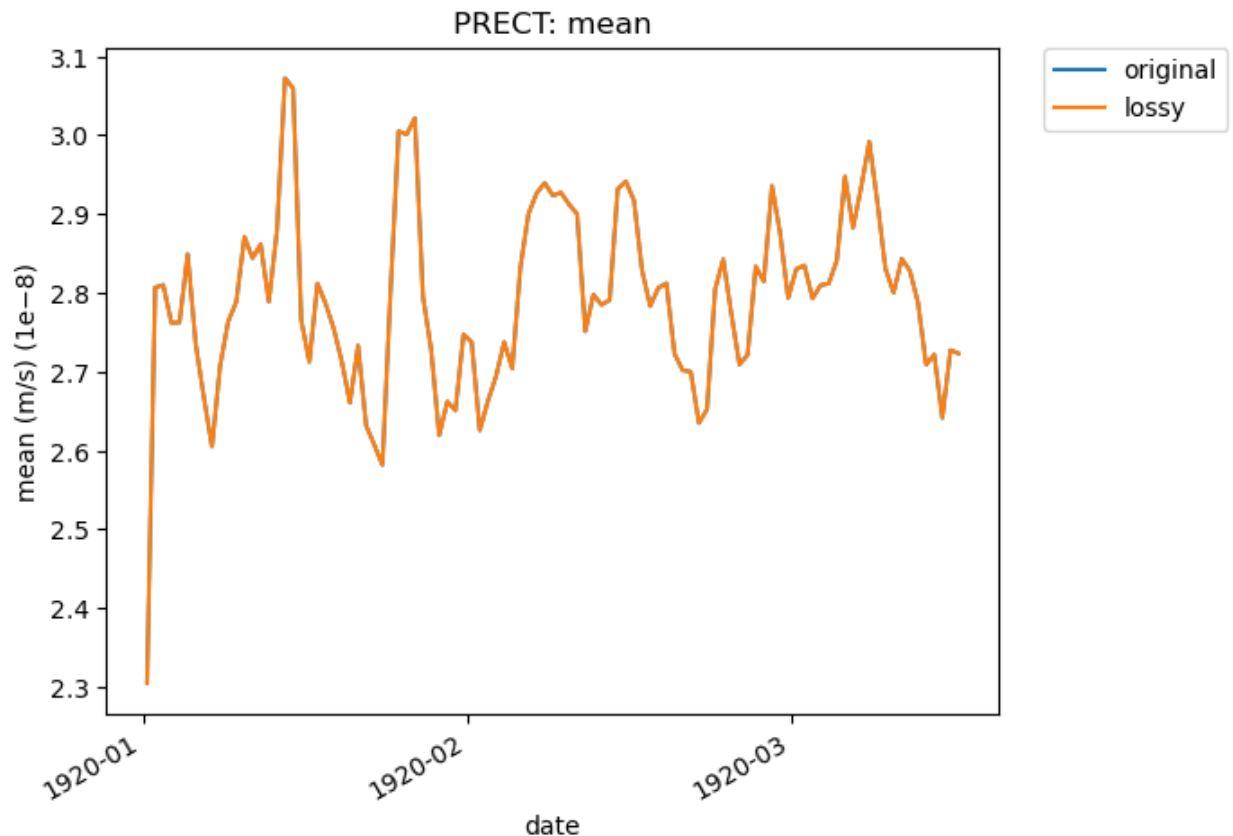
```
dataset size in GB 27.63
```

```
[14]: # PRECT collection
col_prect = ldcpy.collect_datasets(["PRECT"], [ds_prect, aws_prect], ["original",
↪ "lossy"])
```

```
dataset size in GB 27.63
```

Now that we have our collections, we can do some comparisons. Note that these are large files, so make sure you have sufficient compute/memory.

```
[16]: # Time-series plot of PRECT mean in col_ds 'original' dataset - first 100 days
ldcpy.plot(
    col_prect,
    "PRECT",
    sets=["original", "lossy"],
    calc="mean",
    plot_type="time_series",
    start=0,
    end=100,
)
```



```
[18]: # print statistics about 'original', 'lossy', and diff between the two datasets for_
↪ TMQ at time slice 365
ldcpy.compare_stats(col_tmq.isel(time=365), "TMQ", "original", "lossy")
```

```

mean original          : 16.586
mean lossy             : 16.541
mean diff              : 0.0047053

variance original      : 239.68
variance lossy         : 238.4

standard deviation original : 15.482
standard deviation lossy   : 15.44

max value original     : 71.03
max value lossy        : 71
min value original     : 0.28534
min value lossy        : 0.28516

max abs diff           : 0.0073329
min abs diff           : 0.0020776
mean abs diff          : 0.0047053
mean squared diff       : 2.2139e-05
root mean squared diff : 0.0053893
normalized root mean squared diff : 7.6178e-05
normalized max pointwise error : 0.00010365
pearson correlation coefficient : 0.99768
ks p-value             : 0.57304
spatial relative error(% > 0.0001) : 98.116
max spatial relative error : 0.0077153
ssim                   : 0.23137
ssim_fp                : nan

```

The original data for FLNS and TMQ and TS and PRECT (from Amazon S3 in the “ncar-cesm-lens-baker-lossy-compression-test” bucket) was loaded above using method 1. An alternative would be to create our own catalog for this data for use with intake-esm. To illustrate this, we created a `test_catalog.csv` and `test_collection.json` file for this particular simple example.

We first open our collection.

```

[19]: my_col_loc = "./collections/test_collection.json"
      my_col = intake.open_esm_datastore(my_col_loc)
      my_col

<IPython.core.display.HTML object>

```

```

[20]: # printing the head() gives us the file names
      my_col.df.head()

```

```

[20]: component frequency experiment variable \
0      atm      daily      20C      TS
1      atm      daily      20C      PRECT
2      atm      daily      20C      FLNS
3      atm      daily      20C      TMQ

                                path
0  s3://ncar-cesm-lens-baker-lossy-compression-te...
1  s3://ncar-cesm-lens-baker-lossy-compression-te...
2  s3://ncar-cesm-lens-baker-lossy-compression-te...
3  s3://ncar-cesm-lens-baker-lossy-compression-te...

```

Let’s load all of these into our dictionary! (So we don’t need to do the search to make a subset of variables as above in

## Method 2.)

```
[21]: my_dset_dict = my_col.to_dataset_dict(
      zarr_kwargs={"consolidated": True, "decode_times": True},
      storage_options={"anon": True},
    )
my_dset_dict

--> The keys in the returned dictionary of datasets are constructed as follows:
      'component.experiment.frequency'

<IPython.core.display.HTML object>
```

```
[21]: {'atm.20C.daily': <xarray.Dataset>
      Dimensions:    (lat: 192, lon: 288, nbnd: 2, time: 31390)
      Coordinates:
        * lat        (lat) float64 -90.0 -89.06 -88.12 -87.17 ... 88.12 89.06 90.0
        * lon        (lon) float64 0.0 1.25 2.5 3.75 5.0 ... 355.0 356.2 357.5 358.8
        * time       (time) object 1920-01-01 12:00:00 ... 2005-12-31 12:00:00
          time_bnds  (time, nbnd) object dask.array<chunksize=(15695, 2), meta=np.ndarray>
      Dimensions without coordinates: nbnd
      Data variables:
        TS           (time, lat, lon) float32 dask.array<chunksize=(576, 192, 288),
        ↪meta=np.ndarray>
        PRECT        (time, lat, lon) float32 dask.array<chunksize=(576, 192, 288),
        ↪meta=np.ndarray>
        FLNS         (time, lat, lon) float32 dask.array<chunksize=(576, 192, 288),
        ↪meta=np.ndarray>
        TMQ          (time, lat, lon) float32 dask.array<chunksize=(576, 192, 288),
        ↪meta=np.ndarray>
      Attributes:
        revision_Id:    $Id$
        logname:        mickelso
        initial_file:    b.e11.B20TRC5CNBDRD.f09_g16.001.cam.i.1920-01-01...
        case:           b.e11.B20TRC5CNBDRD.f09_g16.031
        host:           ys0219
        title:          UNSET
        NCO:            netCDF Operators version 4.7.9 (Homepage = http:...
        Version:        $Name$
        source:         CAM
        Conventions:    CF-1.0
        topography_file: /glade/p/cesmdata/cseg/inputdata/atm/cam/topo/US...
        intake_esm_dataset_key: atm.20C.daily}
```

```
[22]: # we again just want the 20th century daily data
my_ds = my_dset_dict["atm.20C.daily"]
my_ds
```

```
[22]: <xarray.Dataset>
      Dimensions:    (lat: 192, lon: 288, nbnd: 2, time: 31390)
      Coordinates:
        * lat        (lat) float64 -90.0 -89.06 -88.12 -87.17 ... 88.12 89.06 90.0
        * lon        (lon) float64 0.0 1.25 2.5 3.75 5.0 ... 355.0 356.2 357.5 358.8
        * time       (time) object 1920-01-01 12:00:00 ... 2005-12-31 12:00:00
          time_bnds  (time, nbnd) object dask.array<chunksize=(15695, 2), meta=np.ndarray>
      Dimensions without coordinates: nbnd
      Data variables:
        TS           (time, lat, lon) float32 dask.array<chunksize=(576, 192, 288), meta=np.
        ↪ndarray>
```

(continues on next page)

(continued from previous page)

```

    PRECT      (time, lat, lon) float32 dask.array<chunks=(576, 192, 288), meta=np.
↪ndarray>
    FLNS      (time, lat, lon) float32 dask.array<chunks=(576, 192, 288), meta=np.
↪ndarray>
    TMQ       (time, lat, lon) float32 dask.array<chunks=(576, 192, 288), meta=np.
↪ndarray>
Attributes:
  revision_Id:      $Id$
  logname:          mickelso
  initial_file:     b.e11.B20TRC5CNBDRD.f09_g16.001.cam.i.1920-01-01...
  case:             b.e11.B20TRC5CNBDRD.f09_g16.031
  host:             ys0219
  title:            UNSET
  NCO:              netCDF Operators version 4.7.9 (Homepage = http:...
  Version:          $Name$
  source:           CAM
  Conventions:      CF-1.0
  topography_file:  /glade/p/cesmdata/cseg/inputdata/atm/cam/topo/US...
  intake_esm_dataset_key: atm.20C.daily

```

Now we can make a dataset for each variable as before.

```

[23]: my_ts = my_ds["TS"].to_dataset()
      my_tmq = my_ds["TMQ"].to_dataset()
      my_prect = my_ds["PRECT"].to_dataset()
      my_flns = my_ds["FLNS"].to_dataset()

```

And now we can form new collections as before and do comparisons...

```
[ ]:
```

#### 5.4.4 NCAR JupyterHub Large Data Example Notebook

*Note: If you do not have access to the NCAR machine, please look at the AWS-LENS example notebook instead.*

This notebook demonstrates how to compare large datasets on glade with ldcpy. In particular, we will look at data from CESM-LENS1 project (<http://www.cesm.ucar.edu/projects/community-projects/LENS/data-sets.html>). In doing so, we will start a DASK client from Jupyter. This notebook is meant to be run on NCAR's JupyterHub (<https://jupyterhub.ucar.edu>). We will use a subset of the CESM-LENS1 data on glade is located in /glade/p/cisl/asap/ldcpy\_sample\_data/lens.

We assume that you have a copy of the ldcpy code on NCAR's glade filesystem, obtained via: `git clone https://github.com/NCAR/ldcpy.git`

When you launch your NCAR JupyterHub session, you will need to indicate a machine (Cheyenne or Casper) and then you will need your charge account. You can then launch the session and navigate to this notebook.

NCAR's JupyterHub documentation: <https://www2.cisl.ucar.edu/resources/jupyterhub-ncar>

Here's another good resource for using NCAR's JupyterHub: <https://ncar-hackathons.github.io/jupyterlab-tutorial/jhub.html>

**You need to run your notebook with the “cmip6-201910” kernel (choose from the dropdown in the upper left.)**

Note that the compressed data that we are using was generated for this paper:

Allison H. Baker, Dorit M. Hammerling, Sheri A. Mickelson, Haiying Xu, Martin B. Stolpe, Phillipe Naveau, Ben Sanderson, Imme Ebert-Uphoff, Savini Samarasinghe, Francesco De Simone, Francesco Carbone, Christian

N. Gencarelli, John M. Dennis, Jennifer E. Kay, and Peter Lindstrom, “Evaluating Lossy Data Compression on Climate Simulation Data within a Large Ensemble.” *Geoscientific Model Development*, 9, pp. 4381-4403, 2016 (<https://gmd.copernicus.org/articles/9/4381/2016/>)

## Setup

Let’s set up our environment. First, make sure that you are using the `cmip6-2019.10` kernel. Then you will need to modify the path below to indicate where you have cloned `ldcpy`. (*Note: soon we will install `ldcpy` on Cheyenne/Casper in the `cmip6-2019.10` kernel .*)

If you want to use the `dask` dashboard, then the `dask.config` link must be set below (modify for your path in your browser).

```
[1]: # Make sure you are using the cmip6-2019.10 kernel

# Add ldcpy root to system path (MODIFY FOR YOUR LDCPY CODE LOCATION)
import sys

sys.path.insert(0, '/glade/u/home/abaker/repos/ldcpy')
import ldcpy

# Display output of plots directly in Notebook
%matplotlib inline
# Automatically reload module if it is edited
%reload_ext autoreload
%autoreload 2

# silence warnings
import warnings

warnings.filterwarnings("ignore")

# if you want to use the DASK dashboard on Casper, then modify the below and run
# import dask
# dask.config.set({'distributed.dashboard.link' : 'https://jupyterhub.ucar.edu/dav/
↪user/abaker/proxy/{port}/status'})

# if you want to use the DASK dashboard on Cheyenne, then modify the below and run
import dask

dask.config.set(
    {'distributed.dashboard.link': 'https://jupyterhub.ucar.edu/ch/user/abaker/proxy/
↪{port}/status'}
)

[1]: <dask.config.set at 0x2af0285f9f50>
```

**Connect to DASK distributed cluster (Cheyenne uses PBS / Casper uses slurm):**

(The cluster object is for a single compute node: <https://jobqueue.dask.org/en/latest/howitworks.html>)

```
[2]: # start the dask scheduler

# for Casper
# from dask_jobqueue import SLURMCluster
# cluster = SLURMCluster(memory="40GB", cores=4, processes=1, walltime="02:00:00",
# ↪project="NIOW0001")

# for Cheyenne
from dask_jobqueue import PBSCluster

cluster = PBSCluster(
    queue="regular",
    walltime="02:00:00",
    project="NIOW0001",
    memory="109GB",
    resource_spec="select=1:ncpus=9:mem=109GB",
    cores=36,
    processes=9,
)

# scale as needed
cluster.adapt(minimum_jobs=1, maximum_jobs=30)
cluster

VBox(children=(HTML(value='<h2>PBSCluster</h2>'), HBox(children=(HTML(value='\n<div>\n
↪n <style scoped>\n      .d...
```

The scheduler creates a normal-looking job script that it can submit multiple times to the queue:

```
[3]: # Look at the job script (optional)
print(cluster.job_script())

#!/usr/bin/env bash

#PBS -N dask-worker
#PBS -q regular
#PBS -A NIOW0001
#PBS -l select=1:ncpus=9:mem=109GB
#PBS -l walltime=02:00:00

/ncar/usr/jupyterhub/envs/cmip6-201910/bin/python -m distributed.cli.dask_worker tcp://
↪/10.148.11.122:37213 --nthreads 4 --nprocs 9 --memory-limit 12.11GB --name name --
↪nanny --death-timeout 60 --interface ib0
```

```
[4]: from dask.distributed import Client

# Connect client to the remote dask workers
client = Client(cluster)
client

[4]: <Client: 'tcp://10.148.11.122:37213' processes=0 threads=0, memory=0 B>
```

*Note: click on the dashboard link above to see your DASK tasks*



## The sample data on the glade filesystem

In /glade/p/cisl/asap/ldcpy\_sample\_data/lens on glade, we have TS (surface temperature), PRECT (precipitation rate), and PS (surface pressure) data from CESM-LENS1. These all are 2D variables. TS and PRECT have daily output, and PS has monthly output. We have the compressed and original versions of all these variables that we would like to compare with ldcpy.

First we list what is in this directory (two subdirectories):

```
[5]: # list directory contents
import os

os.listdir("/glade/p/cisl/asap/ldcpy_sample_data/lens")

[5]: ['README.txt', 'lossy', 'orig']
```

Now we look at the contents of each subdirectory. We have 14 files in each, consisting of 2 different timeseries files for each variable (1920-2005 and 2006-2080).

```
[6]: # list lossy directory contents (files that have been lossy compressed and
↳reconstructed)
lossy_files = os.listdir("/glade/p/cisl/asap/ldcpy_sample_data/lens/lossy")
lossy_files
```

```
[6]: ['c.CCN3.monthly.192001-200512.nc',
      'c.FLNS.monthly.192001-200512.nc',
      'c.LHFLX.daily.20060101-20801231.nc',
      'c.TS.daily.19200101-20051231.nc',
      'c.U.monthly.200601-208012.nc',
      'c.TREFHT.monthly.192001-200512.nc',
      'c.LHFLX.daily.19200101-20051231.nc',
      'c.TMQ.monthly.192001-200512.nc',
      'c.PRECT.daily.20060101-20801231.nc',
      'c.PS.monthly.200601-208012.nc',
      'c.FLNS.monthly.200601-208012.nc',
      'c.Q.monthly.192001-200512.nc',
      'c.U.monthly.192001-200512.nc',
      'c.CLOUD.monthly.192001-200512.nc',
      'c.PRECT.daily.19200101-20051231.nc',
      'c.CLOUD.monthly.200601-208012.nc',
      'c.TS.daily.20060101-20801231.nc',
      'c.Q.monthly.200601-208012.nc',
      'c.CCN3.monthly.200601-208012.nc',
      'c.TREFHT.monthly.200601-208012.nc',
      'c.TMQ.monthly.200601-208012.nc',
      'c.PS.monthly.192001-200512.nc']
```

```
[7]: # list orig (i.e., uncompressed) directory contents
orig_files = os.listdir("/glade/p/cisl/asap/ldcpy_sample_data/lens/orig")
orig_files
```

```
[7]: ['PRECT.daily.20060101-20801231.nc',
      'TMQ.monthly.192001-200512.nc',
      'PS.monthly.192001-200512.nc',
      'PRECT.daily.19200101-20051231.nc',
      'TS.daily.20060101-20801231.nc',
      'PS.monthly.200601-208012.nc',
      'TS.daily.19200101-20051231.nc',
      'TMQ.monthly.200601-208012.nc',
```

(continues on next page)

(continued from previous page)

```
'FLNS.monthly.192001-200512.nc',
'FLNS.monthly.200601-200812.nc']
```

We can look at how big these files are...

```
[8]: print("Original files")
for f in orig_files:
    print(
        f,
        " ",
        os.stat("/glade/p/cisl/asap/ldcpy_sample_data/lens/orig/" + f).st_size / 1e9,
        "GB",
    )
```

```
Original files
PRECT.daily.20060101-20081231.nc    4.909326733 GB
TMQ.monthly.192001-200512.nc      0.160075734 GB
PS.monthly.192001-200512.nc       0.12766304 GB
PRECT.daily.19200101-20051231.nc   5.629442994 GB
TS.daily.20060101-20081231.nc     3.435295036 GB
PS.monthly.200601-200812.nc       0.111186435 GB
TS.daily.19200101-20051231.nc     3.962086636 GB
TMQ.monthly.200601-200812.nc      0.139301278 GB
FLNS.monthly.192001-200512.nc     0.170014618 GB
FLNS.monthly.200601-200812.nc     0.148417532 GB
```

## Open datasets

First, let's look at the original and reconstructed files for the monthly surface Pressure (PS) data for 1920-2006. We begin by using `ldcpy.open_dataset()` to open the files of interest into our dataset collection. Usually we want chunks to be 100-150MB, but this is machine and app dependent.

```
[9]: # load the first 86 years of montly surface pressure into a collection

col_PS = ldcpy.open_datasets(
    ["PS"],
    [
        "/glade/p/cisl/asap/ldcpy_sample_data/lens/orig/PS.monthly.192001-200512.nc",
        "/glade/p/cisl/asap/ldcpy_sample_data/lens/lossy/c.PS.monthly.192001-200512.nc",
    ],
    ["orig", "lossy"],
    chunks={"time": 500},
)
col_PS
```

dataset size in GB 0.46

```
[9]: <xarray.Dataset>
Dimensions:      (collection: 2, lat: 192, lon: 288, time: 1032)
Coordinates:
  * lat          (lat) float64 -90.0 -89.06 -88.12 -87.17 ... 88.12 89.06 90.0
  * lon          (lon) float64 0.0 1.25 2.5 3.75 5.0 ... 355.0 356.2 357.5 358.8
  * time         (time) object 1920-02-01 00:00:00 ... 2006-01-01 00:00:00
  * collection   (collection) <U5 'orig' 'lossy'
```

(continues on next page)

(continued from previous page)

```

Data variables:
  PS          (collection, time, lat, lon) float32 dask.array<chunks=(1, 500, 192, 288), meta=np.ndarray>
Attributes:
  Conventions: CF-1.0
  source:      CAM
  case:        b.e11.B20TRC5CNBDRD.f09_g16.031
  title:       UNSET
  logname:     mickelso
  host:        ys0219
  Version:     $Name$
  revision_Id: $Id$
  initial_file: b.e11.B20TRC5CNBDRD.f09_g16.001.cam.i.1920-01-01-00000.nc
  topography_file: /glade/p/cesmdata/cseg/inputdata/atm/cam/topo/USGS-gtop...
  history:     Tue Nov 3 13:51:10 2020: ncks -L 5 PS.monthly.192001-2...
  NCO:         netCDF Operators version 4.7.9 (Homepage = http://nco.s...

```

## Data comparison

Now we use the ldcpy package features to compare the data.

## Surface Pressure

Let's look at the comparison statistics at the first timeslice for PS.

```
[10]: ps0 = col_PS.isel(time=0)
ldcpy.compare_stats(ps0, "PS", "orig", "lossy")
```

```

mean orig          : 96750
mean lossy         : 96734
mean diff          : 15.806

variance orig      : 8.4254e+07
variance lossy     : 8.4249e+07

standard deviation orig : 9179.1
standard deviation lossy : 9178.8

max value orig     : 1.0299e+05
max value lossy    : 1.0298e+05
min value orig     : 51967
min value lossy    : 51952

max abs diff       : 31.992
min abs diff       : 0
mean abs diff      : 15.806
mean squared diff  : 249.83
root mean squared diff : 18.303
normalized root mean squared diff : 0.0003587
normalized max pointwise error : 0.00062698
pearson correlation coefficient : 1
ks p-value         : 1.6583e-05
spatial relative error(% > 0.0001) : 69.085

```

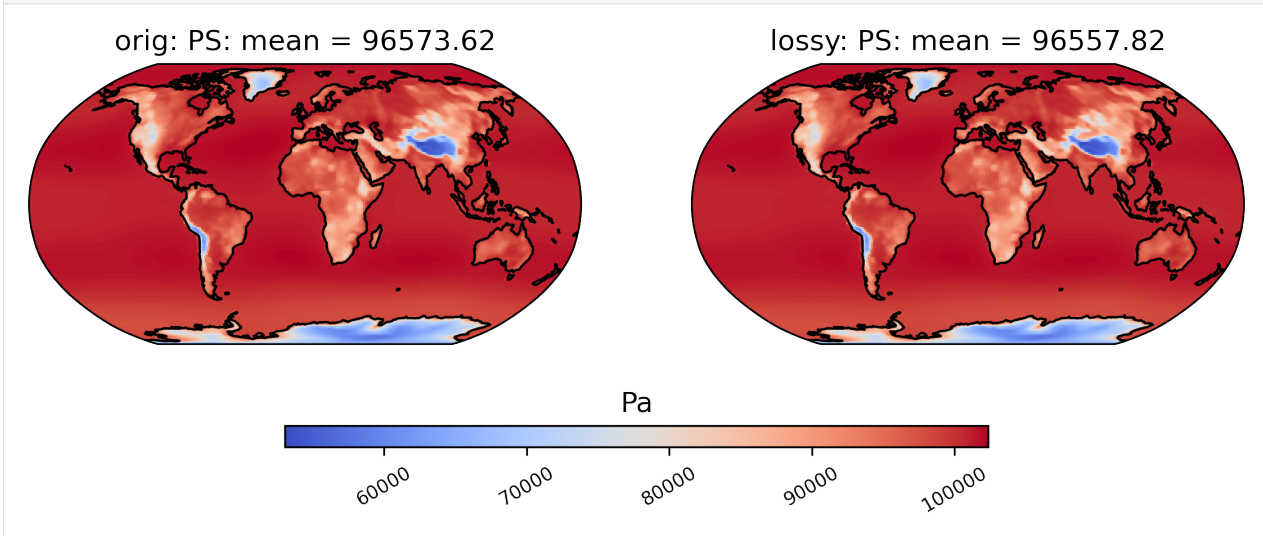
(continues on next page)

(continued from previous page)

```
max spatial relative error      : 0.00048247
ssim                           : 0.99809
ssim_fp                         : 0.98467
```

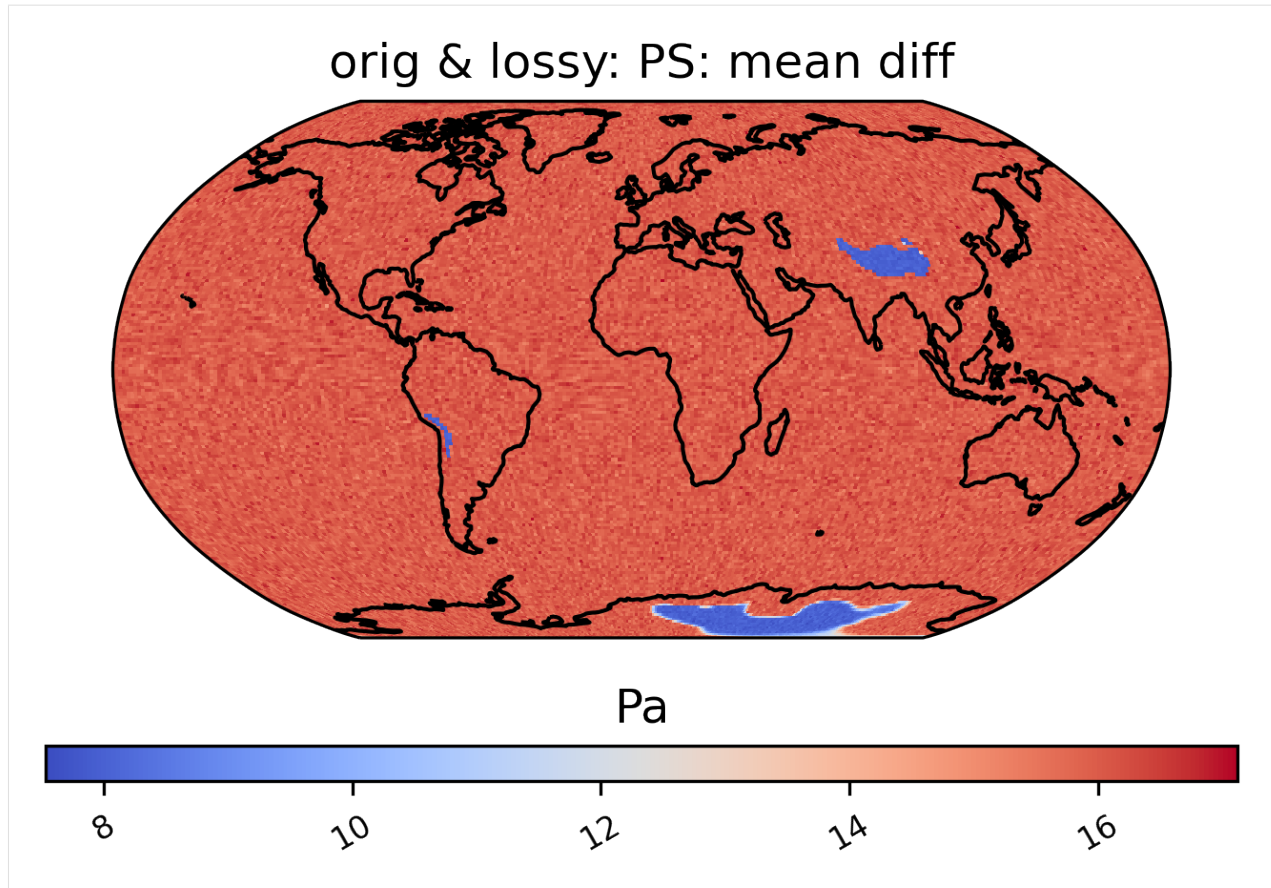
Now we make a plot to compare the mean PS values across time in the orig and lossy datasets.

```
[11]: # comparison between mean PS values (over the 86 years) in col_PS orig and lossy_
      ↪ datasets
ldcpy.plot(col_PS, "PS", sets=["orig", "lossy"], calc="mean")
```



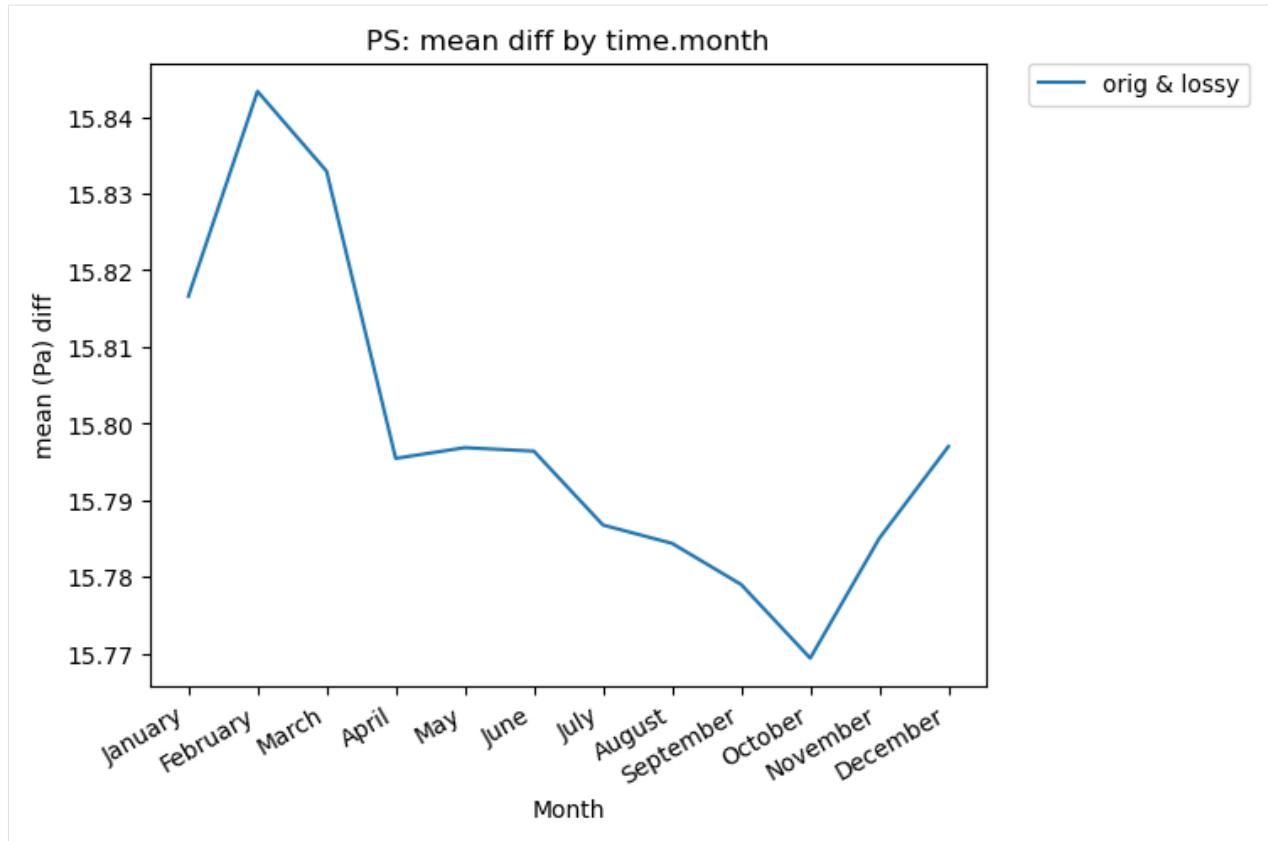
Now we instead show the difference plot for the above plots.

```
[12]: # diff between mean PS values in col_PS orig and lossy datasets
ldcpy.plot(col_PS, "PS", sets=["orig", "lossy"], calc="mean", calc_type="diff")
```



We can also look at mean differences over time. Here we are looking at the spatial averages and then grouping by day of the year. If doing a timeseries plot for this much data, using “group\_by” is often a good idea.

```
[13]: # Time-series plot of mean PS differences between col_PS orig and col_PS lossy,
      ↪ datasets grouped by month of year
ldcpy.plot(
    col_PS,
    "PS",
    sets=["orig", "lossy"],
    calc="mean",
    plot_type="time_series",
    group_by="time.month",
    calc_type="diff",
)
```



```
[14]: # Time-series plot of PS mean (grouped by month) in the original and lossy datasets
ldcpy.plot(
    col_PS,
    "PS",
    sets=["orig", "lossy"],
    calc="mean",
    plot_type="time_series",
    group_by="time.month",
)
```



```
[15]: del col_PS
```

## Surface Temperature

Now let's open the daily surface temperature (TS) data for 1920-2005 into a collection. These are larger files than the monthly PS data.

```
[16]: # load the first 86 years of daily surface temperature (TS) data
col_TS = ldcpy.open_datasets(
    ["TS"],
    [
        "/glade/p/cisl/asap/ldcpy_sample_data/lens/orig/TS.daily.19200101-20051231.nc
    ↪ ",
        "/glade/p/cisl/asap/ldcpy_sample_data/lens/lossy/c.TS.daily.19200101-20051231.
    ↪ nc",
    ],
    ["orig", "lossy"],
    chunks={"time": 500},
)
col_TS
```

dataset size in GB 13.89

```
[16]: <xarray.Dataset>
Dimensions:      (collection: 2, lat: 192, lon: 288, time: 31390)
Coordinates:
  * lat           (lat) float64 -90.0 -89.06 -88.12 -87.17 ... 88.12 89.06 90.0
  * lon           (lon) float64 0.0 1.25 2.5 3.75 5.0 ... 355.0 356.2 357.5 358.8
  * time          (time) object 1920-01-01 00:00:00 ... 2005-12-31 00:00:00
  * collection    (collection) <U5 'orig' 'lossy'
Data variables:
  TS              (collection, time, lat, lon) float32 disk.array<chunksizes=(1, 500, 192, 288), meta=np.ndarray>
Attributes:
  Conventions:    CF-1.0
  source:         CAM
  case:           b.e11.B20TRC5CNBDRD.f09_g16.031
  title:          UNSET
  logname:        mickelso
  host:           ys0219
  Version:        $Name$
  revision_Id:    $Id$
  initial_file:   b.e11.B20TRC5CNBDRD.f09_g16.001.cam.i.1920-01-01-00000.nc
  topography_file: /glade/p/cesmdata/cseg/inputdata/atm/cam/topo/USGS-gtop...
  history:        Tue Nov 3 13:56:03 2020: ncks -L 5 TS.daily.19200101-2...
  NCO:            netCDF Operators version 4.7.9 (Homepage = http://nco.s...
```

Look at the first time slice (time = 0) statistics:

```
[17]: ldcpy.compare_stats(col_TS.isel(time=0), "TS", "orig", "lossy")
```

```
mean orig           : 274.71
mean lossy          : 274.66
mean diff           : 0.054906

variance orig       : 533.98
variance lossy      : 533.43

standard deviation orig : 23.108
standard deviation lossy : 23.096

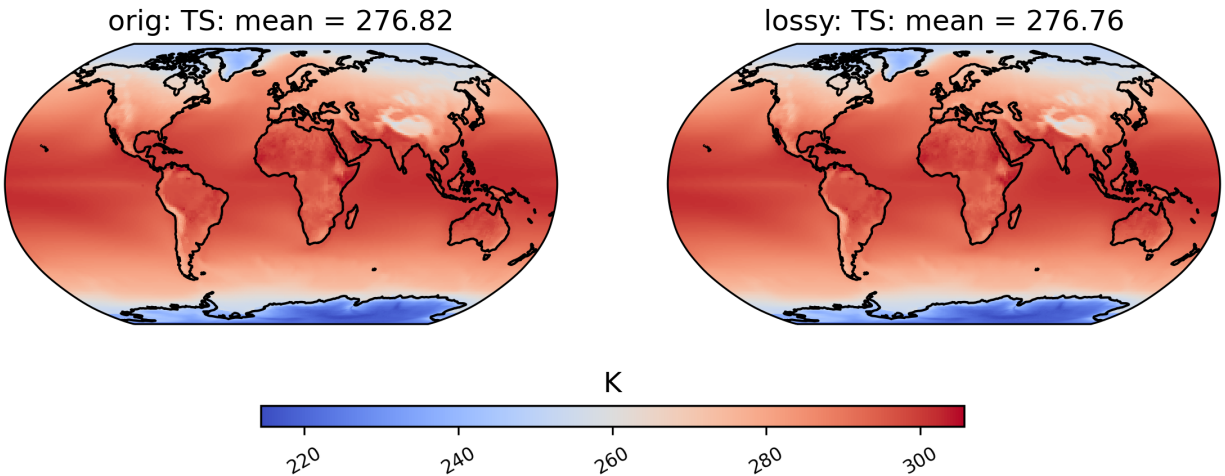
max value orig      : 315.58
max value lossy     : 315.5
min value orig      : 216.73
min value lossy     : 216.69

max abs diff        : 0.12497
min abs diff        : 0
mean abs diff       : 0.054906
mean squared diff   : 0.0030146
root mean squared diff : 0.06527
normalized root mean squared diff : 0.0006603
normalized max pointwise error : 0.0012642
pearson correlation coefficient : 1
ks p-value          : 0.36817
spatial relative error(% > 0.0001) : 73.293
max spatial relative error : 0.00048733
ssim                : 0.9985
ssim_fp             : 0.99794
```

Now we compare mean TS over time in a plot:

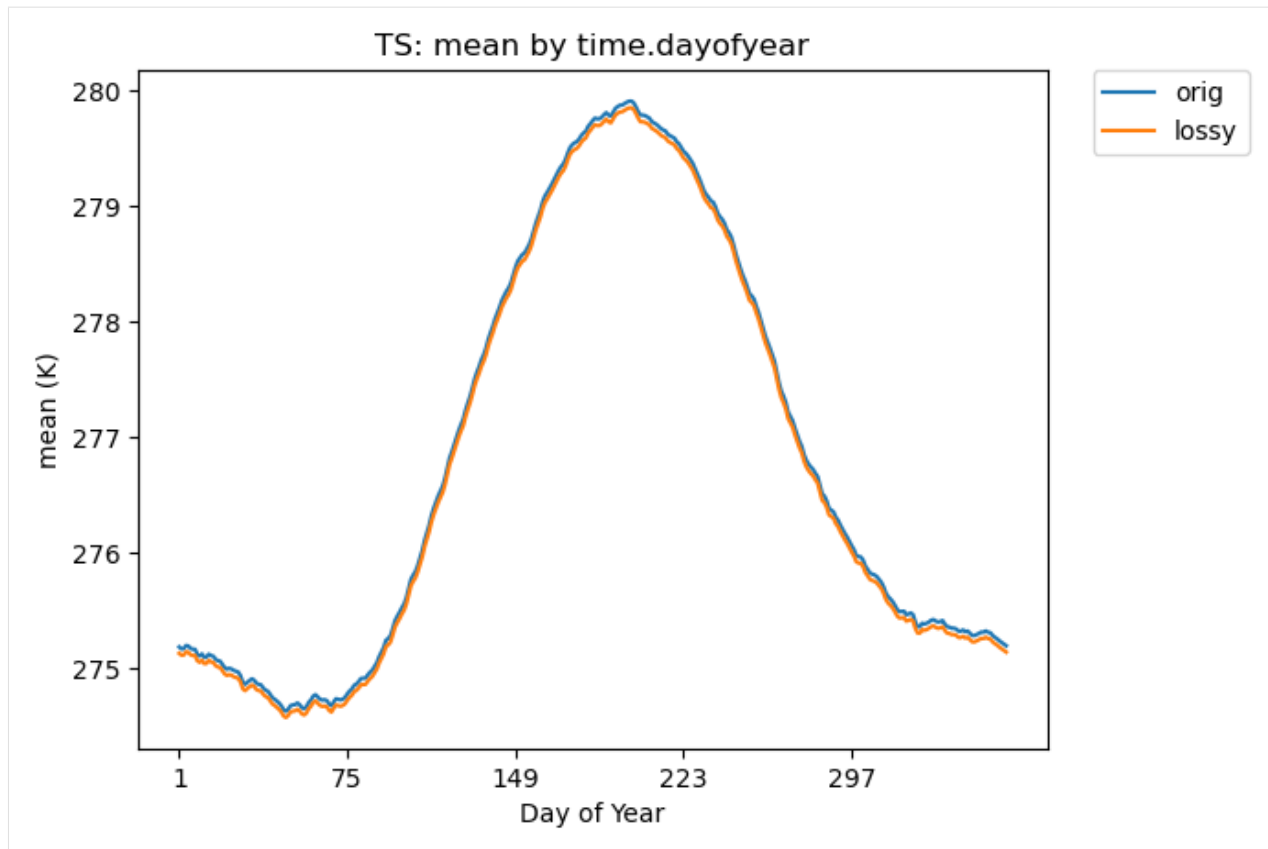


```
[18]: # comparison between mean TS values in col_TS orig and lossy datasets
ldcpy.plot(col_TS, "TS", sets=["orig", "lossy"], calc="mean")
```



Below we do a time series plot and group by day of the year. (Note that the `group_by` functionality is not fast.)

```
[19]: # Time-series plot of TS means (grouped by days) in the original and lossy datasets
ldcpy.plot(
    col_TS,
    "TS",
    sets=["orig", "lossy"],
    calc="mean",
    plot_type="time_series",
    group_by="time.dayofyear",
)
```



Let's delete the PS and TS data to free up memory.

```
[20]: del col_TS
```

## Precipitation Rate

Now let's open the daily precipitation rate (PRECT) data for 2006-2080 into a collection.

```
[21]: # load the last 75 years of PRECT data
col_PRECT = ldcpy.open_datasets(
    ["PRECT"],
    [
        "/glade/p/cisl/asap/ldcpy_sample_data/lens/orig/PRECT.daily.20060101-20801231.
↪nc",
        "/glade/p/cisl/asap/ldcpy_sample_data/lens/lossy/c.PRECT.daily.20060101-
↪20801231.nc",
    ],
    ["orig", "lossy"],
    chunks={"time": 500},
)
col_PRECT
```

dataset size in GB 12.11

```
[21]: <xarray.Dataset>
Dimensions:      (collection: 2, lat: 192, lon: 288, time: 27375)
```

(continues on next page)

(continued from previous page)

```

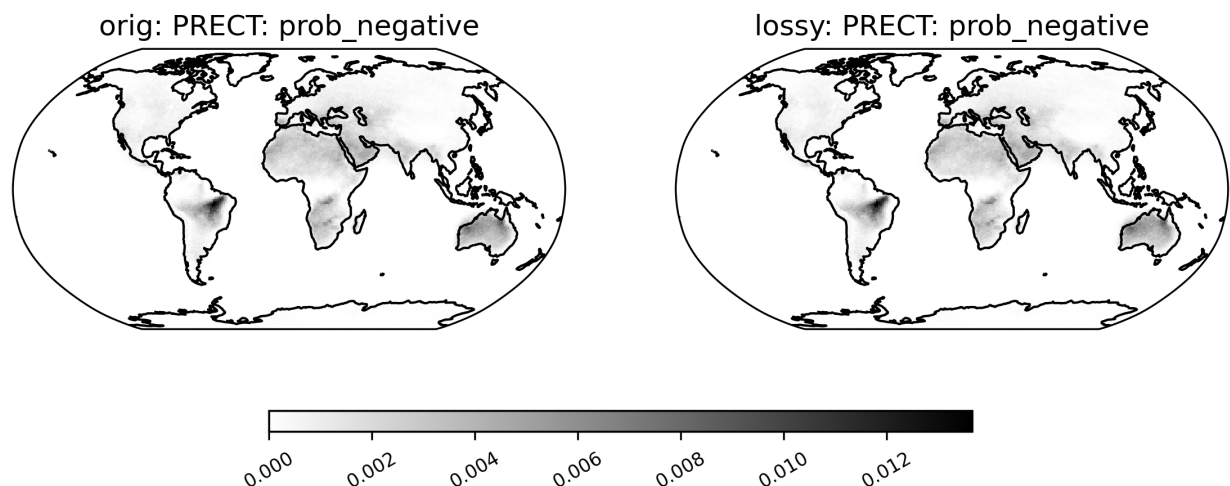
Coordinates:
  * lat          (lat) float64 -90.0 -89.06 -88.12 -87.17 ... 88.12 89.06 90.0
  * lon          (lon) float64 0.0 1.25 2.5 3.75 5.0 ... 355.0 356.2 357.5 358.8
  * time         (time) object 2006-01-01 00:00:00 ... 2080-12-31 00:00:00
  * collection   (collection) <U5 'orig' 'lossy'
Data variables:
  PRECT          (collection, time, lat, lon) float32 dask.array<chunks=(1, 500, 192, 288), meta=np.ndarray>
Attributes:
  Conventions:   CF-1.0
  source:        CAM
  case:          b.e11.BRCP85C5CNBDRD.f09_g16.031
  title:         UNSET
  logname:       mickelso
  host:          ys1023
  Version:       $Name$
  revision_Id:   $Id$
  initial_file:  b.e11.B20TRC5CNBDRD.f09_g16.031.cam.i.2006-01-01-00000.nc
  topography_file: /glade/p/cesmdata/cseg/inputdata/atm/cam/topo/USGS-gtop...
  history:       Tue Nov 3 14:13:51 2020: ncks -L 5 PRECT.daily.2006010...
  NCO:          netCDF Operators version 4.7.9 (Homepage = http://nco.s...

```

```
[22]: # compare probability of negative rainfall (and get ssim)
```

```
ldcpy.plot(
    col_PRECT,
    "PRECT",
    sets=["orig", "lossy"],
    calc="prob_negative",
    color="binary",
    calc_ssim=True,
)
```

```
SSIM 1 & 2 = 1.00000
```

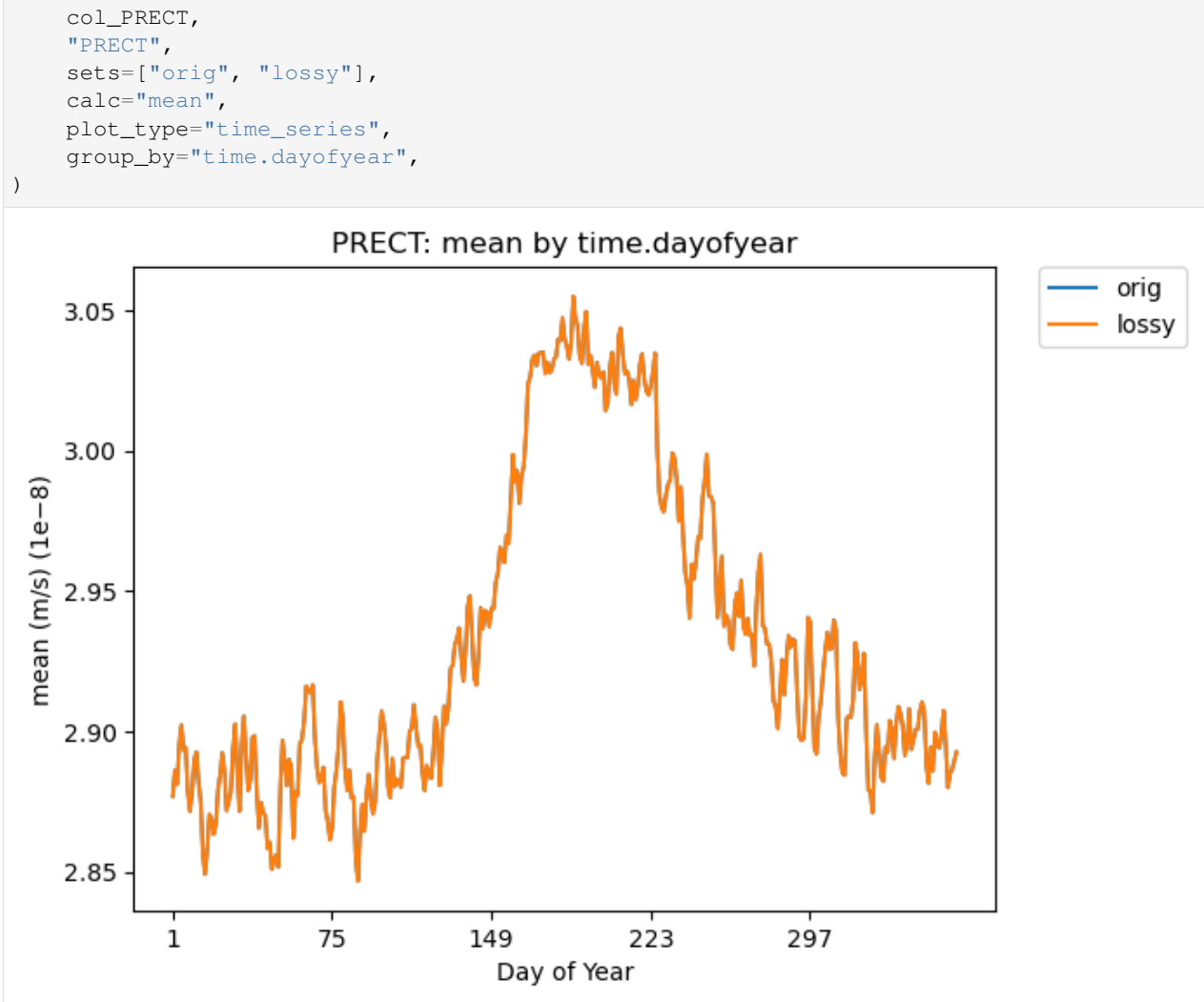


Mean PRECT over time...

```
[23]: # Time-series plot of PRECT mean in 'orig' dataset
ldcpy.plot(
```

(continues on next page)

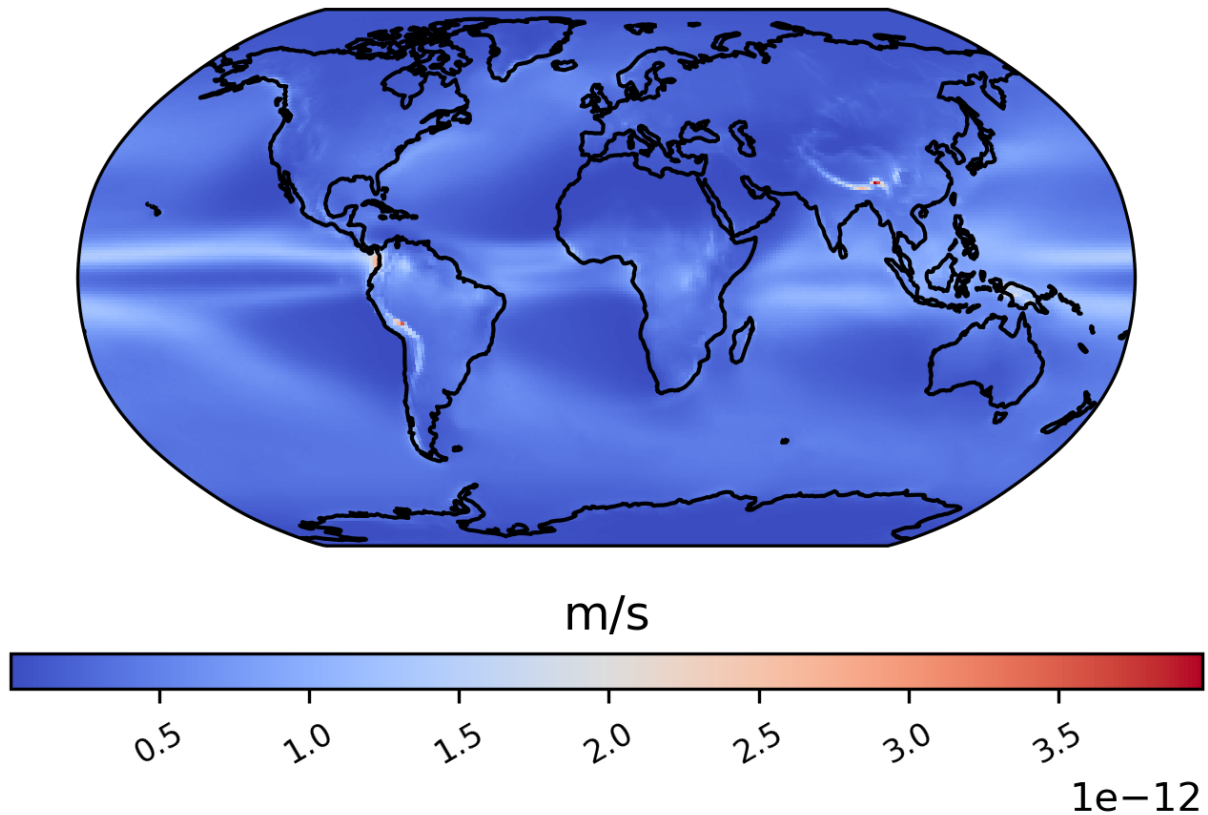
(continued from previous page)



Now look at mean over time spatial plot:

```
[24]: # diff between mean PRECT values across the entire timeseries
ldcpy.plot(
    col_PRECT,
    "PRECT",
    sets=["orig", "lossy"],
    calc="mean",
    calc_type="diff",
)
```

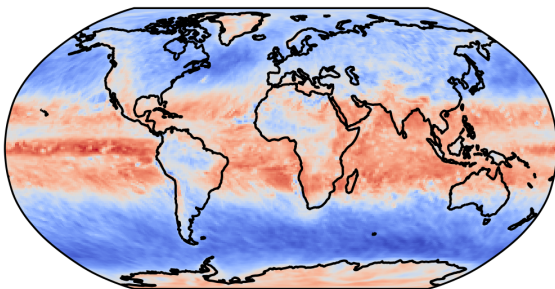
orig & lossy: PRECT: mean diff



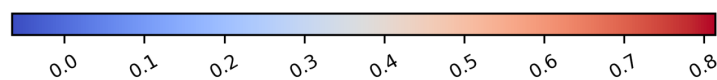
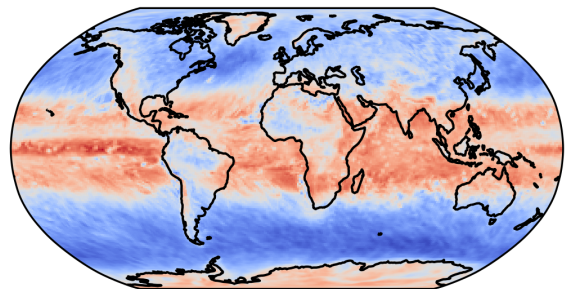
Calculating the correlation of the lag-1 values ... for the first 10 years

```
[25]: # plot of lag-1 correlation of PRECT values
ldcpy.plot(col_PRECT, "PRECT", sets=["orig", "lossy"], calc="lag1", start=0, end=3650)
```

orig: PRECT: lag1



lossy: PRECT: lag1



## CAGEO Plots

Comparing the sz and zfp compressors with a number of metrics (specified with “calc” and “calc\_type” in the plot routine).

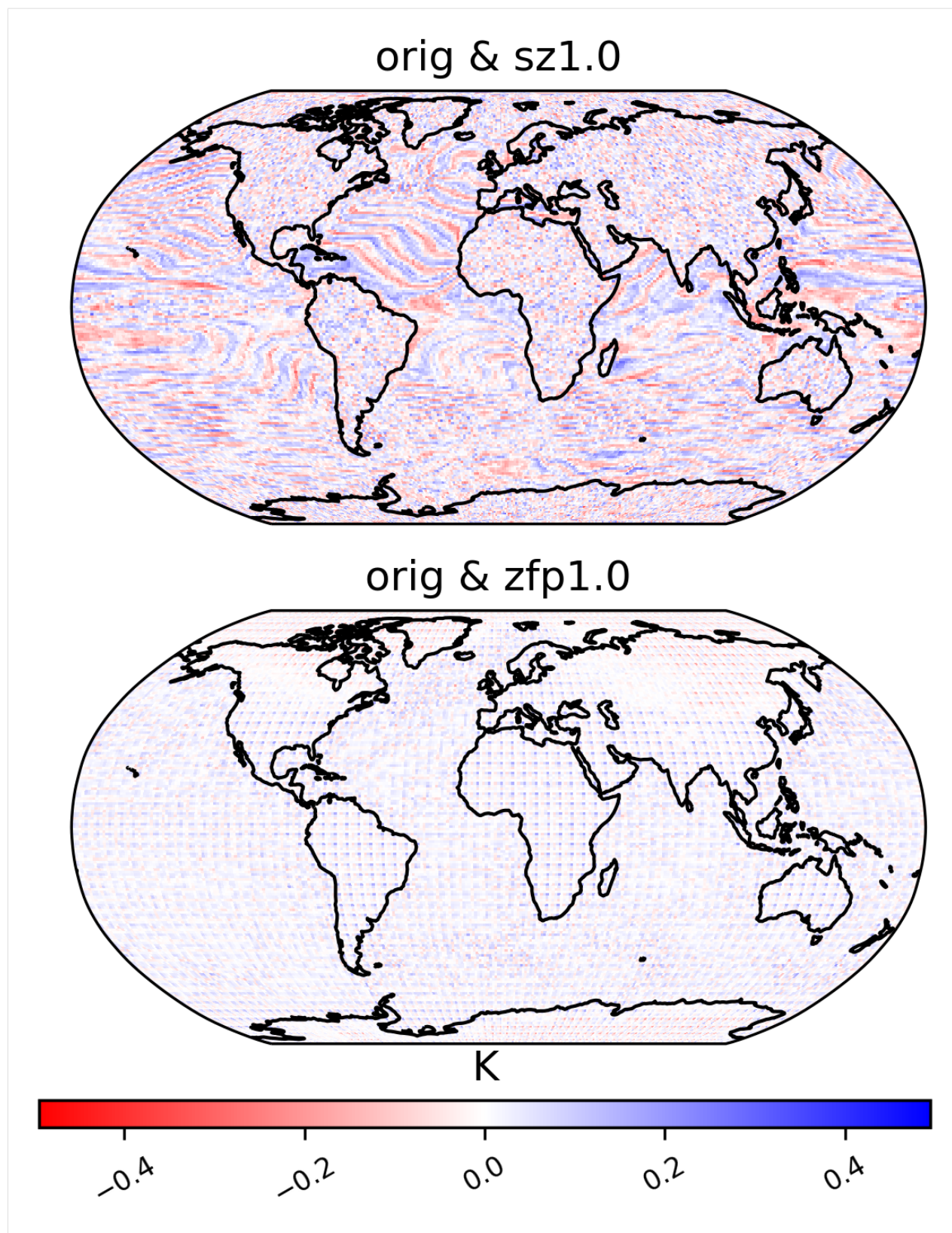
```
[26]: col_TS = ldcpy.open_datasets(  
    ["TS"],  
    [  
        "/glade/p/cisl/asap/ldcpy_sample_data/cageo/orig/b.e11.B20TRC5CNBDRD.f09_g16.  
→030.cam.h1.TS.19200101-20051231.nc",  
        "/glade/p/cisl/asap/ldcpy_sample_data/cageo/lossy/sz1.0.TS.nc",  
        "/glade/p/cisl/asap/ldcpy_sample_data/cageo/lossy/zfp1.0.TS.nc",  
    ],  
    ["orig", "sz1.0", "zfp1.0"],  
    chunks={"time": 500},  
)  
col_PRECT = ldcpy.open_datasets(  
    ["PRECT"],  
    [  
        "/glade/p/cisl/asap/ldcpy_sample_data/cageo/orig/b.e11.B20TRC5CNBDRD.f09_g16.  
→030.cam.h1.PRECT.19200101-20051231.nc",  
        "/glade/p/cisl/asap/ldcpy_sample_data/cageo/lossy/sz1e-8.PRECT.nc",  
        "/glade/p/cisl/asap/ldcpy_sample_data/cageo/lossy/zfp1e-8.PRECT.nc",  
    ],  
    ["orig", "sz1e-8", "zfp1e-8"],  
    chunks={"time": 500},  
)
```

dataset size in GB 20.83

dataset size in GB 20.83

difference in mean

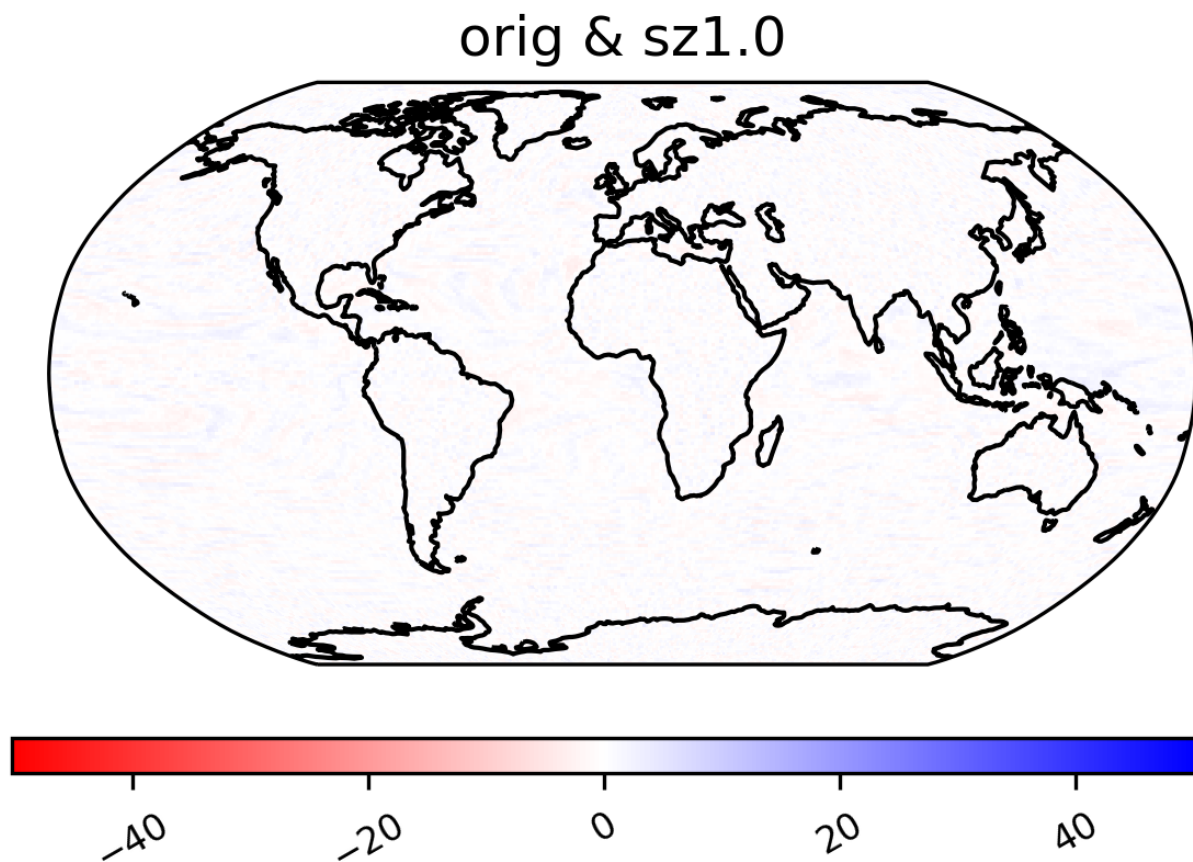
```
[27]: ldcpy.plot(  
    col_TS,  
    "TS",  
    sets=["orig", "sz1.0", "zfp1.0"],  
    calc="mean",  
    calc_type="diff",  
    color="bwr_r",  
    start=0,  
    end=50,  
    tex_format=False,  
    vert_plot=True,  
    short_title=True,  
    axes_symmetric=True,  
)
```



Zscore



```
[28]: ldcpy.plot(
    col_TS,
    "TS",
    sets=["orig", "sz1.0"],
    calc="zscore",
    calc_type="metric_of_diff",
    color="bwr_r",
    start=0,
    end=1000,
    tex_format=False,
    vert_plot=True,
    short_title=True,
    axes_symmetric=True,
)
```



Pooled variance ratio

```
[29]: ldcpy.plot(
    col_TS,
    "TS",
    sets=["orig", "sz1.0"],
    calc="pooled_var_ratio",
    color="bwr_r",
    calc_type="metric_of_diff",
    transform="log",
```

(continues on next page)



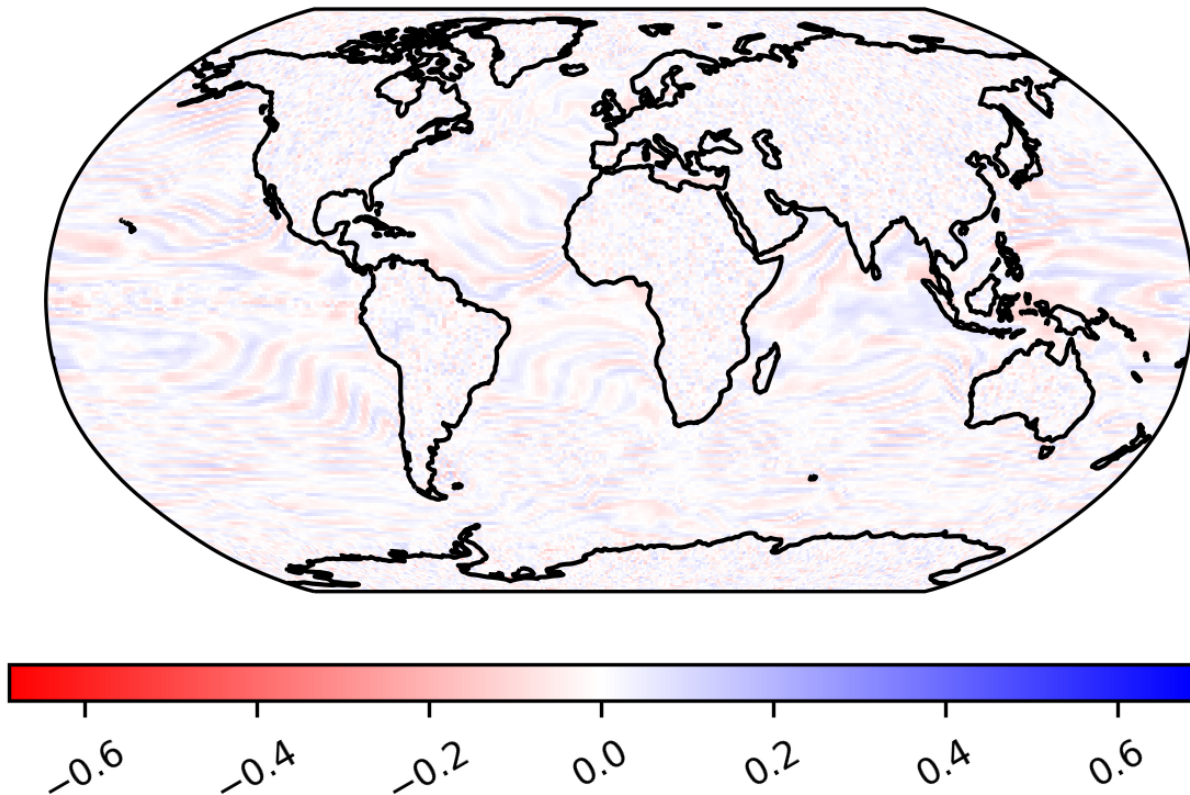
(continued from previous page)

```

start=0,
end=100,
tex_format=False,
vert_plot=True,
short_title=True,
axes_symmetric=True,
)

```

orig & sz1.0

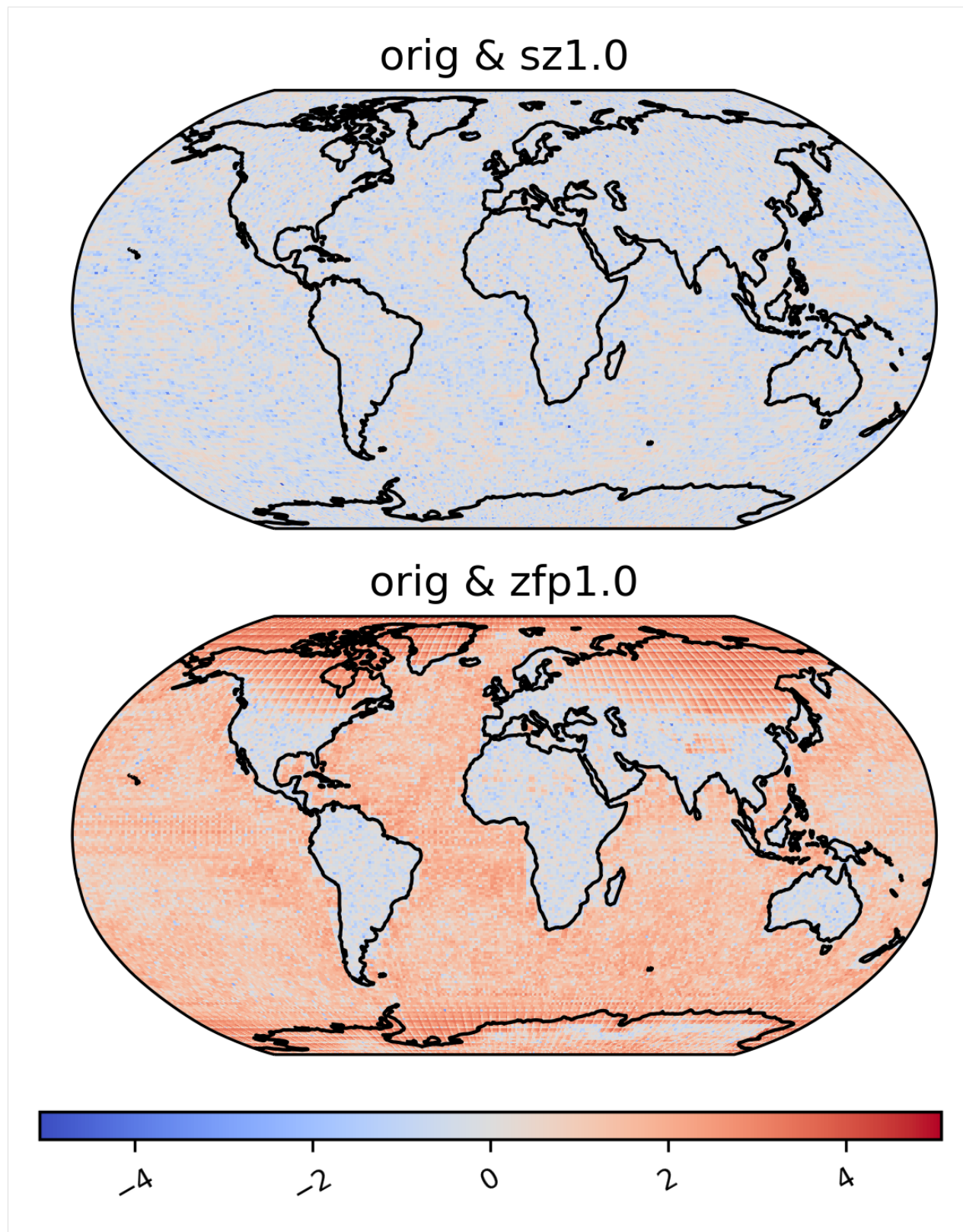


annual harmonic relative ratio

```

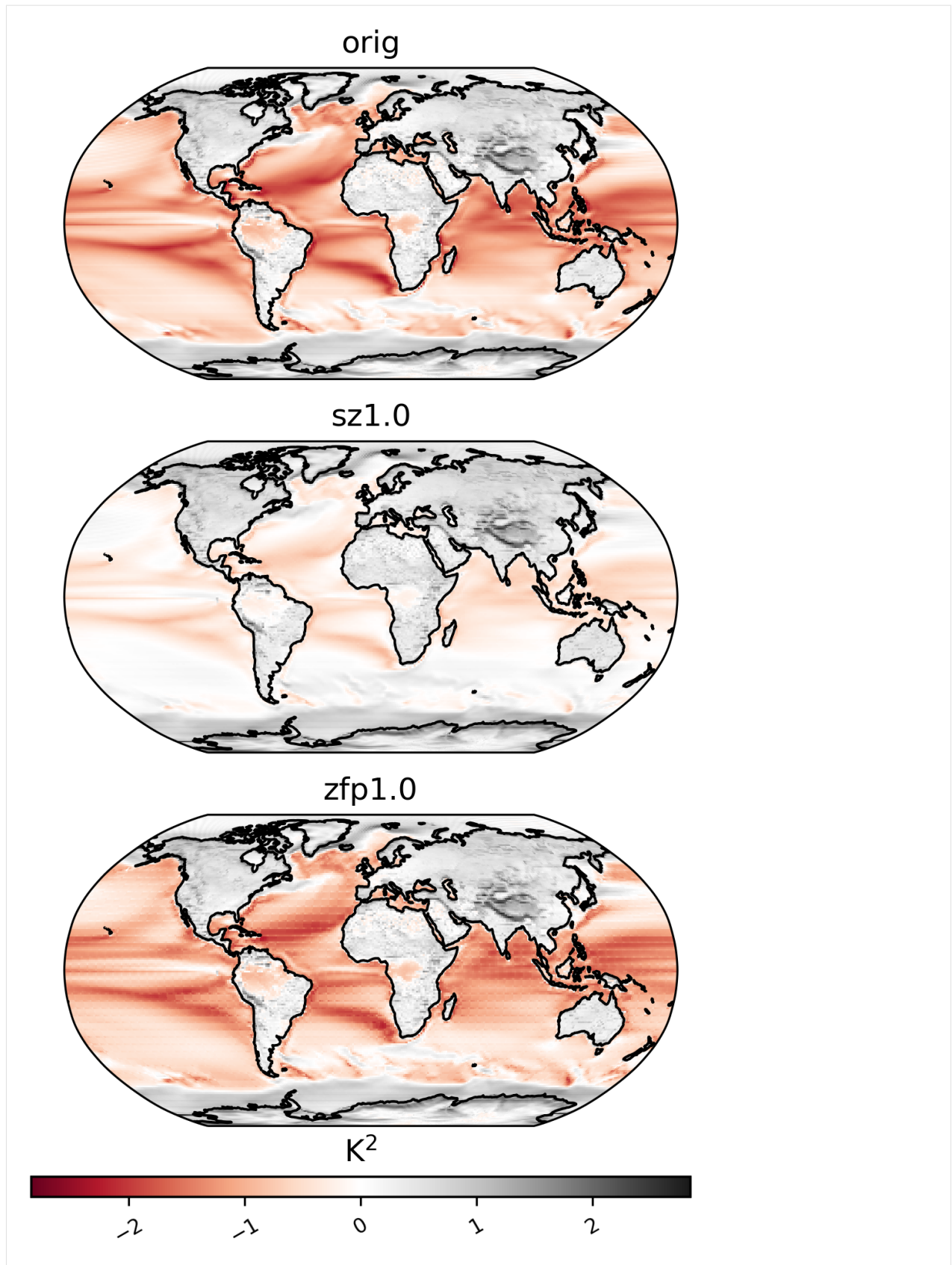
[30]: col_TS["TS"] = col_TS.TS.chunk({"time": -1, "lat": 10, "lon": 10})
ldcpy.plot(
    col_TS,
    "TS",
    sets=["orig", "sz1.0", "zfp1.0"],
    calc="ann_harmonic_ratio",
    transform="log",
    calc_type="metric_of_diff",
    tex_format=False,
    axes_symmetric=True,
    vert_plot=True,
    short_title=True,
)
col_TS["TS"] = col_TS.TS.chunk({"time": 500, "lat": 192, "lon": 288})

```



NS contrast variance

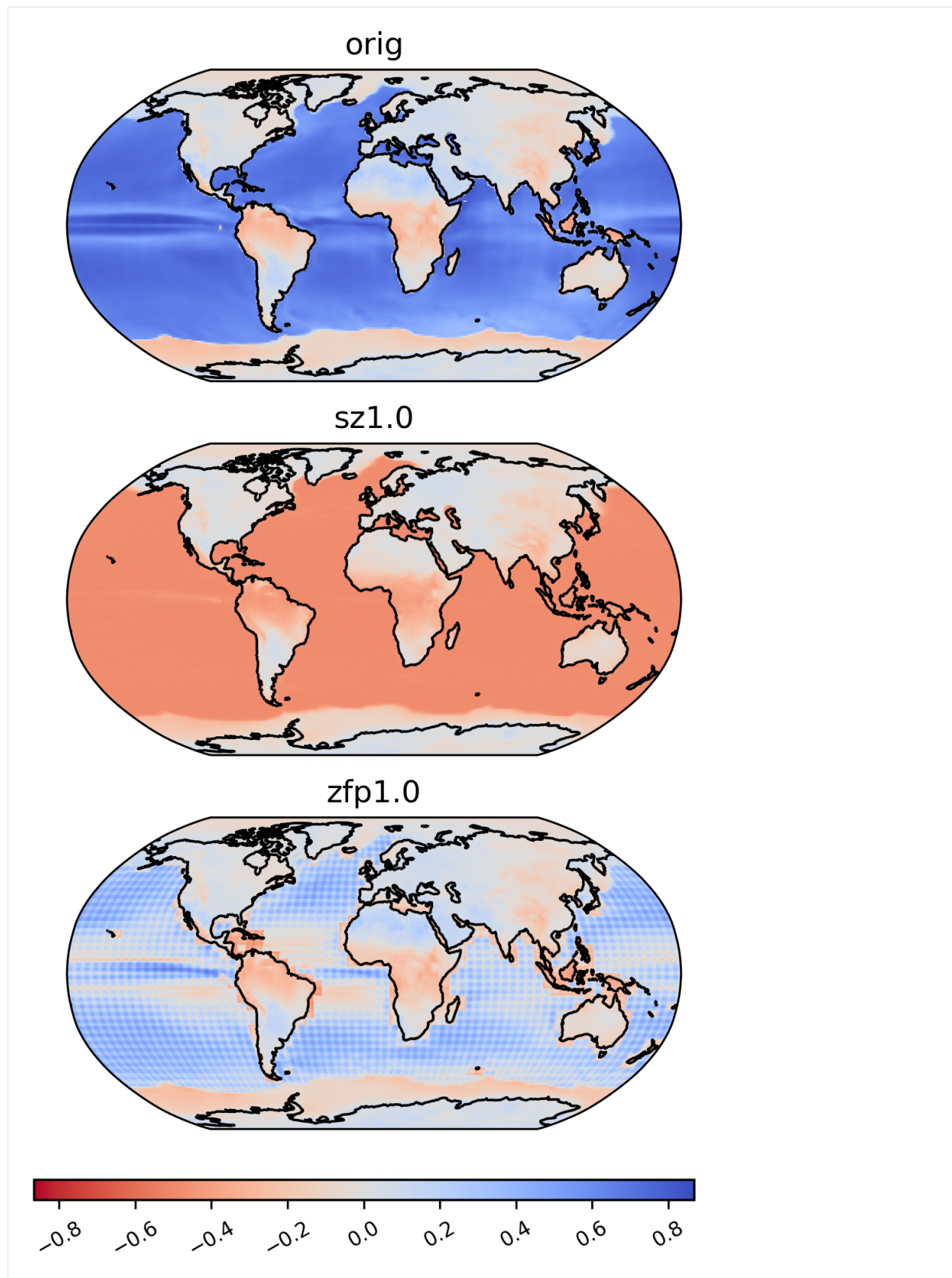
```
[31]: ldcpy.plot(  
    col_TS,  
    "TS",  
    sets=["orig", "sz1.0", "zfp1.0"],  
    calc="ns_con_var",  
    color="RdGy",  
    calc_type="raw",  
    transform="log",  
    axes_symmetric=True,  
    tex_format=False,  
    vert_plot=True,  
    short_title=True,  
)
```



deseasonalized lag-1 autocorrelations of first differences

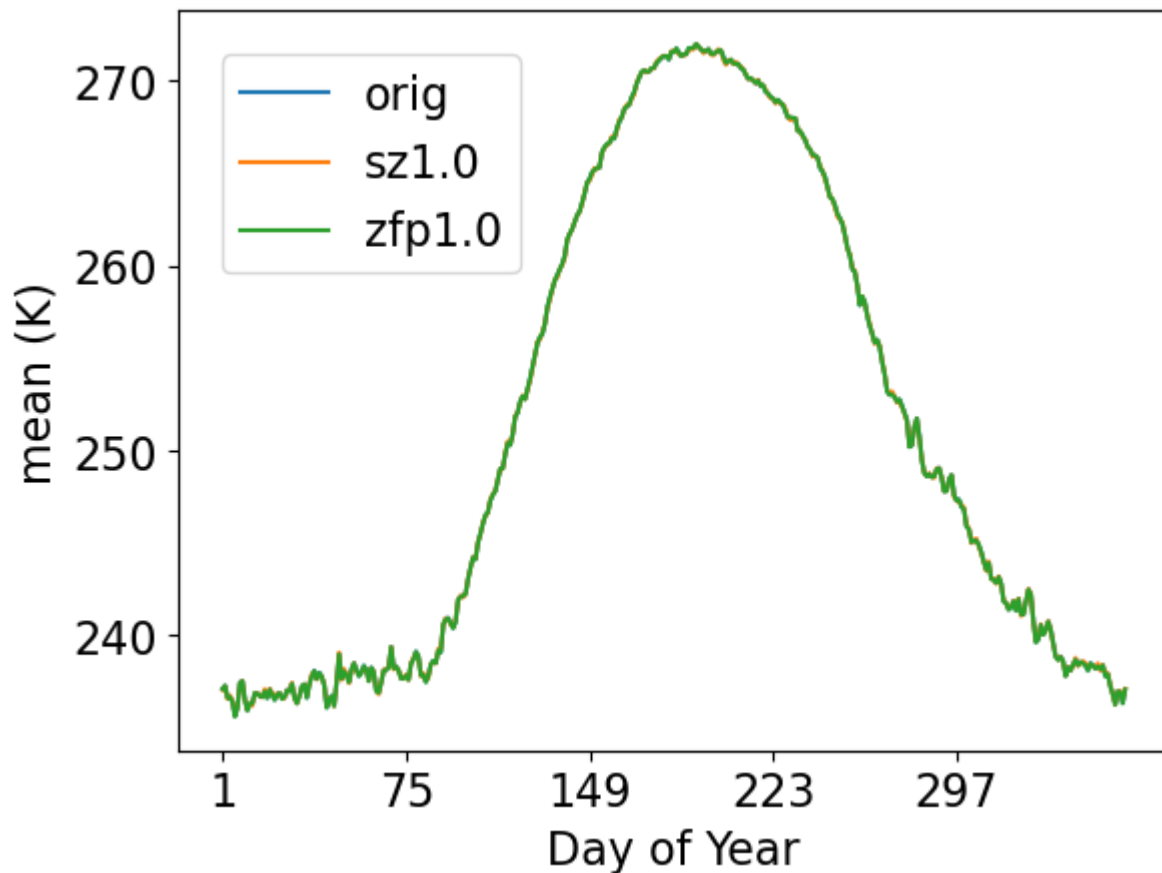
```
[32]: ldcpy.plot(  
    col_TS,  
    "TS",  
    sets=["orig", "sz1.0", "zfp1.0"],  
    calc="lag1_first_difference",  
    color="coolwarm_r",  
    calc_type="raw",  
    tex_format=False,  
    vert_plot=True,  
    axes_symmetric=True,  
    short_title=True,  
)
```





mean by day of year

```
[33]: ldcpy.plot(
    col_TS,
    "TS",
    sets=["orig", "sz1.0", "zfp1.0"],
    calc="mean",
    plot_type="time_series",
    group_by="time.dayofyear",
    legend_loc="upper left",
    lat=90,
    lon=0,
    vert_plot=True,
    tex_format=False,
    short_title=True,
)
```

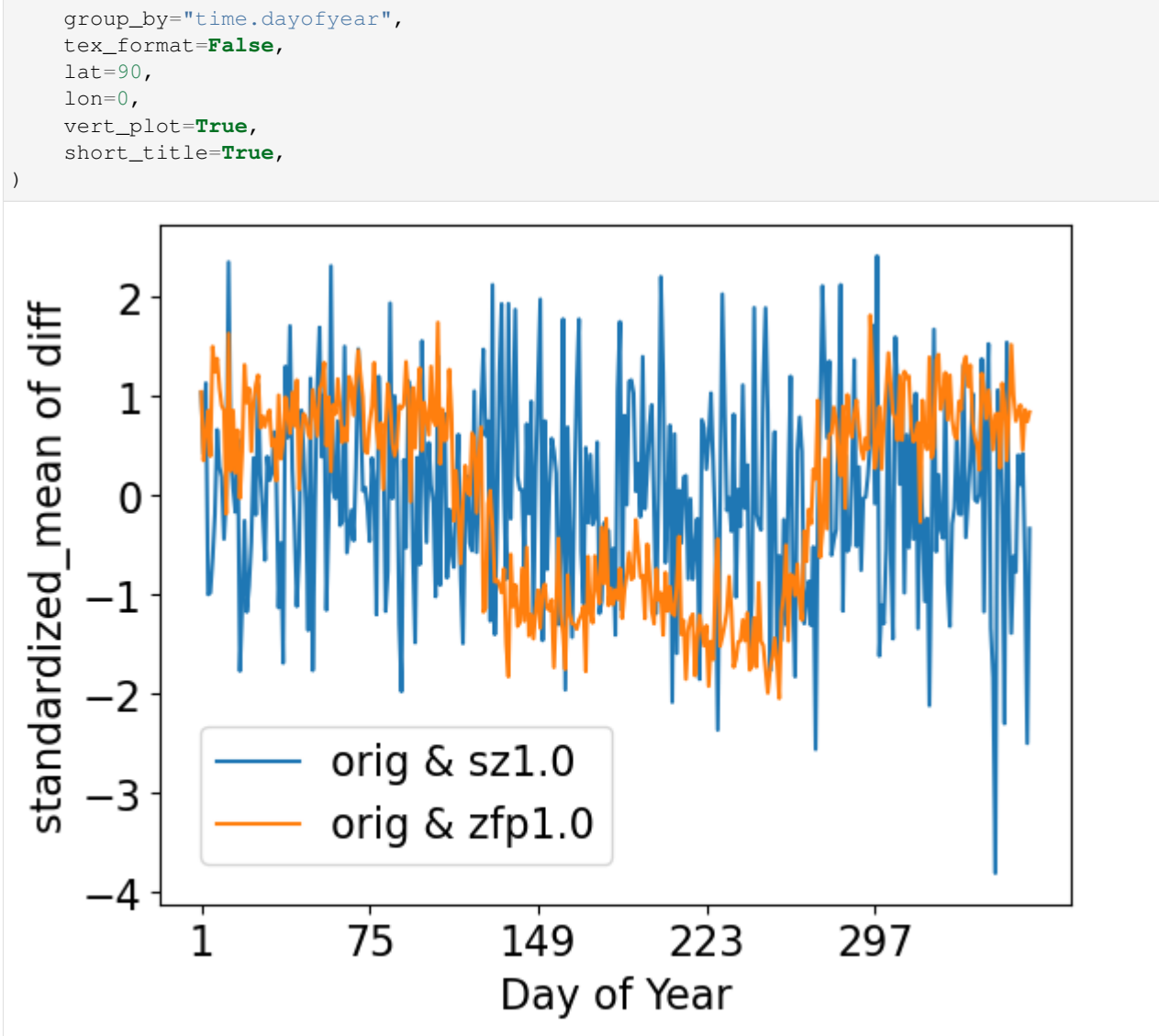


standardized mean errors by day of year

```
[34]: ldcpy.plot(
    col_TS,
    "TS",
    sets=["orig", "sz1.0", "zfp1.0"],
    calc="standardized_mean",
    legend_loc="lower left",
    calc_type="metric_of_diff",
    plot_type="time_series",
)
```

(continues on next page)

(continued from previous page)



```
[35]: del col_TS
```

Do any other comparisons you wish ... and then clean up!

```
[36]: cluster.close()
      client.close()
```

```
[ ]:
```

This Notebook demonstrates comparing different levels of compression on several test variables in the CESM-LENS1 dataset (<http://www.cesm.ucar.edu/projects/community-projects/LENS/data-sets.html>). In doing so, we will start a DASK client from Jupyter.



### 5.4.5 Setup

Before using this notebook, we recommend familiarity with the TutorialNotebook and the LargeDataGladeNotebook. This notebook is meant to be run on NCAR's JupyterHub (<https://jupyterhub.ucar.edu>). The subset of CESM-LENS1 data on glade is located in `/glade/p/cisl/asap/abaker/compression_samples/cam-lens`.

When you launch your NCAR JupyterHub session, you will need to indicate a machine (Cheyenne or Casper) and then you will need your charge account. You can then launch the session and navigate to this notebook.

NCAR's JupyterHub documentation: <https://www2.cisl.ucar.edu/resources/jupyterhub-ncar>

**You need to run your notebook with the “cmip6-201910” kernel (choose from the dropdown in the upper left.)**

When you launch your NCAR JupyterHub session, you will need to indicate a machine (Cheyenne or Casper) and then you will need your charge account. You can then launch the session and navigate to this notebook.

NCAR's JupyterHub documentation: <https://www2.cisl.ucar.edu/resources/jupyterhub-ncar>

```
[1]: # Make sure you are using the cmpi6-2019.10 kernel

# Add ldcpy root to system path (MODIFY FOR YOUR LDCPY CODE LOCATION)
import sys

sys.path.insert(0, '/glade/u/home/abaker/repos/ldcpy')
import ldcpy

# Display output of plots directly in Notebook
%matplotlib inline
# Automatically reload module if it is edited
%reload_ext autoreload
%autoreload 2

# silence warnings
import warnings

warnings.filterwarnings("ignore")
```

Start DASK...

```
[2]: # start the dask scheduler

# for Casper
# from dask_jobqueue import SLURMCluster
# cluster = SLURMCluster(memory="40GB", cores=4, processes=1, walltime="02:00:00",
↳project="NIOW0001")

# for Cheyenne
from dask_jobqueue import PBSCluster

cluster = PBSCluster(
    queue="regular",
    walltime="02:00:00",
    project="NIOW0001",
    memory="109GB",
    resource_spec="select=1:ncpus=9:mem=109GB",
    cores=36,
    processes=9,
)
```

(continues on next page)

(continued from previous page)

```
# scale as needed
cluster.adapt(minimum_jobs=1, maximum_jobs=30)
cluster

VBox(children=(HTML(value='<h2>PBSCluster</h2>'), HBox(children=(HTML(value='\n<div>\n
↳n <style scoped>\n      .d...
```

Connect to the client...

```
[3]: from dask.distributed import Client

# Connect client to the remote dask workers
client = Client(cluster)
client

[3]: <Client: 'tcp://10.148.11.122:44459' processes=0 threads=0, memory=0 B>
```

## 5.4.6 Sample Data

In the glade directory mentionned above (/glade/p/cisl/asap/abaker/compression\_samples/cam-lens), we have an “orig” directory with the original (uncompressed) version of each sample variable: Q, FLNS, TREFHT, TMQ, PS, U, PRECT, LHFLX, CCN3, TS, and CLOUD. Then there is a directory for each sample variable with a variety of compressed variants.

Note that each test variables has a two timeseries slices: 1920-2005 and 2006-2080. The range of data is indicated in the filename along with whether the data is “daily” or “monthly” timeslices.

Variables U, CLOUD, and Qare 3D, and the rest are 2D.

```
[4]: # list directory contents
import os

os.listdir("/glade/p/cisl/asap/abaker/compression_samples/cam-lens")

[4]: ['README.txt',
      'Q',
      'FLNS',
      'TREFHT',
      'TMQ',
      'PS',
      'U',
      'PRECT',
      'LHFLX',
      'CCN3',
      'orig',
      'CLOUD',
      'TS']

[5]: # List the original data files
os.listdir("/glade/p/cisl/asap/abaker/compression_samples/cam-lens/orig")

[5]: ['CLOUD.monthly.200601-208012.nc',
      'LHFLX.daily.20060101-20801231.nc',
      'TREFHT.monthly.192001-200512.nc',
      'U.monthly.192001-200512.nc',
      'PRECT.daily.20060101-20801231.nc',
```

(continues on next page)

(continued from previous page)

```
'CLOUD.monthly.192001-200512.nc',
'TMQ.monthly.192001-200512.nc',
'CCN3.monthly.200601-208012.nc',
'PS.monthly.192001-200512.nc',
'CCN3.monthly.192001-200512.nc',
'PRECT.daily.19200101-20051231.nc',
'TS.daily.20060101-20801231.nc',
'PS.monthly.200601-208012.nc',
'TS.daily.19200101-20051231.nc',
'Q.monthly.192001-200512.nc',
'Q.monthly.200601-208012.nc',
'TMQ.monthly.200601-208012.nc',
'U.monthly.200601-208012.nc',
'TREFHT.monthly.200601-208012.nc',
'FLNS.monthly.192001-200512.nc',
'FLNS.monthly.200601-208012.nc',
'LHFLX.daily.19200101-20051231.nc']
```

```
[6]: # List the compressed data files for TMQ
os.listdir("/glade/p/cisl/asap/abaker/compression_samples/cam-lens/TMQ")
```

```
[6]: ['zfp.p22.TMQ.monthly.192001-200512.nc',
'zfp.p18.TMQ.monthly.200601-208012.nc',
'zfp.p22.TMQ.monthly.200601-208012.nc',
'zfp.p20.TMQ.monthly.192001-200512.nc',
'zfp.p12.TMQ.monthly.200601-208012.nc',
'zfp.p8.TMQ.monthly.192001-200512.nc',
'zfp.p18.TMQ.monthly.192001-200512.nc',
'zfp.p12.TMQ.monthly.192001-200512.nc',
'zfp.p8.TMQ.monthly.200601-208012.nc',
'zfp.p10.TMQ.monthly.192001-200512.nc',
'fpzip16.TMQ.monthly.200601-208012.nc',
'zfp.p16.TMQ.monthly.192001-200512.nc',
'zfp.p20.TMQ.monthly.200601-208012.nc',
'zfp.p10.TMQ.monthly.200601-208012.nc',
'zfp.p14.TMQ.monthly.200601-208012.nc',
'zfp.p16.TMQ.monthly.200601-208012.nc',
'zfp.p14.TMQ.monthly.192001-200512.nc',
'fpzip16.TMQ.monthly.192001-200512.nc']
```

Make a collection with some of the TMQ data to compare. We'll look at the 2006-2080 data. Make sure to use “useful” labels. You can't mix the 2006-2080 and 1920-2005 data in the same collection as they have different numbers of time slices.

The “fpzip” version is from the blind evaluation. For the “zfp” versions, the p indicates the precision parameter (the higher the number after p, the more accurate).

```
[7]: col_tmq = ldcpy.open_datasets(
    ["TMQ"],
    [
        "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/orig/TMQ.monthly.
↪200601-208012.nc",
        "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/TMQ/fpzip16.TMQ.
↪monthly.200601-208012.nc",
        "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/TMQ/zfp.p8.TMQ.
↪monthly.200601-208012.nc",
        "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/TMQ/zfp.p10.TMQ.
↪monthly.200601-208012.nc",
```

(continues on next page)

(continued from previous page)

```

        "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/TMQ/zfp.p12.TMQ.
↪monthly.200601-208012.nc",
        "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/TMQ/zfp.p14.TMQ.
↪monthly.200601-208012.nc",
        "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/TMQ/zfp.p16.TMQ.
↪monthly.200601-208012.nc",
        "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/TMQ/zfp.p18.TMQ.
↪monthly.200601-208012.nc",
        "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/TMQ/zfp.p20.TMQ.
↪monthly.200601-208012.nc",
        "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/TMQ/zfp.p22.TMQ.
↪monthly.200601-208012.nc",
    ],
    [
        "orig",
        "fpzip16",
        "zfp-p8",
        "zfp-p10",
        "zfp-p12",
        "zfp-p14",
        "zfp-p16",
        "zfp-p18",
        "zfp-p20",
        "zfp-p22",
    ],
    chunks={"time": 700},
)
col_tmq

```

dataset size in GB 1.99

```

[7]: <xarray.Dataset>
Dimensions:      (collection: 10, lat: 192, lon: 288, time: 900)
Coordinates:
  * lat          (lat) float64 -90.0 -89.06 -88.12 -87.17 ... 88.12 89.06 90.0
  * lon          (lon) float64 0.0 1.25 2.5 3.75 5.0 ... 355.0 356.2 357.5 358.8
  * time         (time) object 2006-02-01 00:00:00 ... 2081-01-01 00:00:00
  * collection   (collection) <U7 'orig' 'fpzip16' ... 'zfp-p20' 'zfp-p22'
Data variables:
  TMQ            (collection, time, lat, lon) float32 dask.array<chunks=(1, 700,
↪192, 288), meta=np.ndarray>
Attributes:
  Conventions:   CF-1.0
  source:        CAM
  case:          b.e11.BRCP85C5CNBDRD.f09_g16.031
  title:         UNSET
  logname:       mickelso
  host:          ys1023
  Version:       $Name$
  revision_Id:   $Id$
  initial_file:  b.e11.B20TRC5CNBDRD.f09_g16.031.cam.i.2006-01-01-00000.nc
  topography_file: /glade/p/cesmdata/cseg/inputdata/atm/cam/topo/USGS-gtop...
  history:       Tue Nov 3 14:06:43 2020: ncks -L 5 TMQ.monthly.200601-...
  NCO:           netCDF Operators version 4.7.9 (Homepage = http://nco.s...

```

```

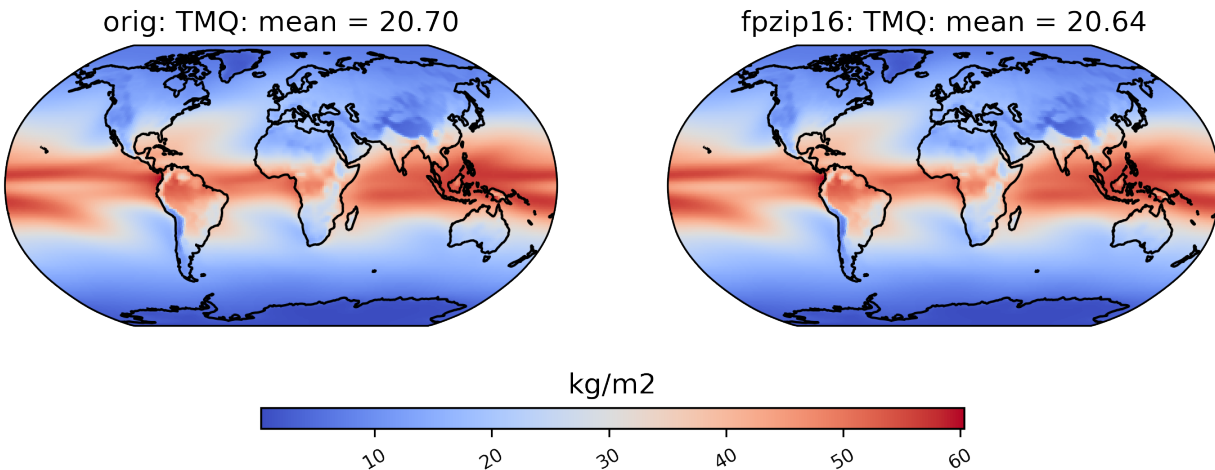
[8]: # first compare the mean TMQ values in col_tmq orig and fpzip datasets (which is in_
↪the released lens data)

```

(continues on next page)

(continued from previous page)

```
ldcpy.plot(col_tmq, "TMQ", sets=["orig", "fpzip16"], calc="mean")
```



Let's look at some statistics for the time slice= 100 of the data (the data has 900 time slices)

```
[10]: ldcpy.compare_stats(col_tmq.isel(time=100), "TMQ", "orig", "fpzip16")
```

```
mean orig           : 18.83
mean fpzip16        : 18.777
mean diff           : 0.052669

variance orig       : 220.84
variance fpzip16    : 219.64

standard deviation orig : 14.861
standard deviation fpzip16 : 14.82

max value orig      : 56.345
max value fpzip16   : 56.25
min value orig      : 0.16246
min value fpzip16   : 0.16211

max abs diff        : 0.24997
min abs diff        : 0
mean abs diff       : 0.052669
mean squared diff   : 0.002774
root mean squared diff : 0.077692
normalized root mean squared diff : 0.0013828
normalized max pointwise error : 0.0044493
pearson correlation coefficient : 1
ks p-value          : 0.0067025
spatial relative error(% > 0.0001) : 98.168
max spatial relative error : 0.0077457
ssim                : 0.99844
ssim_fp            : 0.99862
```

```
[12]: # The compression ratio (CR) for fpzip was ~3.4x better than lossless
# zfp - p10 has a similar CR - look at it's stats
ldcpy.compare_stats(col_tmq.isel(time=100), "TMQ", "orig", "zfp-p10")
```

```
mean orig          : 18.83
mean zfp-p10       : 18.848
mean diff          : -0.018232

variance orig      : 220.84
variance zfp-p10   : 221.26

standard deviation orig : 14.861
standard deviation zfp-p10 : 14.875

max value orig     : 56.345
max value zfp-p10  : 56.5
min value orig     : 0.16246
min value zfp-p10  : 0.16162

max abs diff       : 0.53075
min abs diff       : 2.5332e-07
mean abs diff      : 0.054665
mean squared diff   : 0.00033242
root mean squared diff : 0.08664
normalized root mean squared diff : 0.0015421
normalized max pointwise error : 0.0082268
pearson correlation coefficient : 0.99998
ks p-value         : 0.13637
spatial relative error(% > 0.0001) : 37.299
max spatial relative error : 0.030925
ssim               : 0.99751
ssim_fp            : 0.99584
```

```
[14]: # zfp - p12 has a CR of 2.5x, time=100)
ldcpy.compare_stats(col_tmq.isel(time=0), "TMQ", "orig", "zfp-p12")
```

```
mean orig          : 17.373
mean zfp-p12       : 17.377
mean diff          : -0.0040598

variance orig      : 241.72
variance zfp-p12   : 241.82

standard deviation orig : 15.547
standard deviation zfp-p12 : 15.551

max value orig     : 59.1
max value zfp-p12  : 59.109
min value orig     : 0.43742
min value zfp-p12  : 0.4375

max abs diff       : 0.14342
min abs diff       : 0
mean abs diff      : 0.013363
mean squared diff   : 1.6482e-05
root mean squared diff : 0.022262
normalized root mean squared diff : 0.0003795
normalized max pointwise error : 0.0020729
pearson correlation coefficient : 1
```

(continues on next page)

(continued from previous page)

```

ks p-value           : 0.39259
spatial relative error(% > 0.0001) : 35.408
max spatial relative error : 0.0053198
ssim                 : 0.99942
ssim_fp              : 0.99969

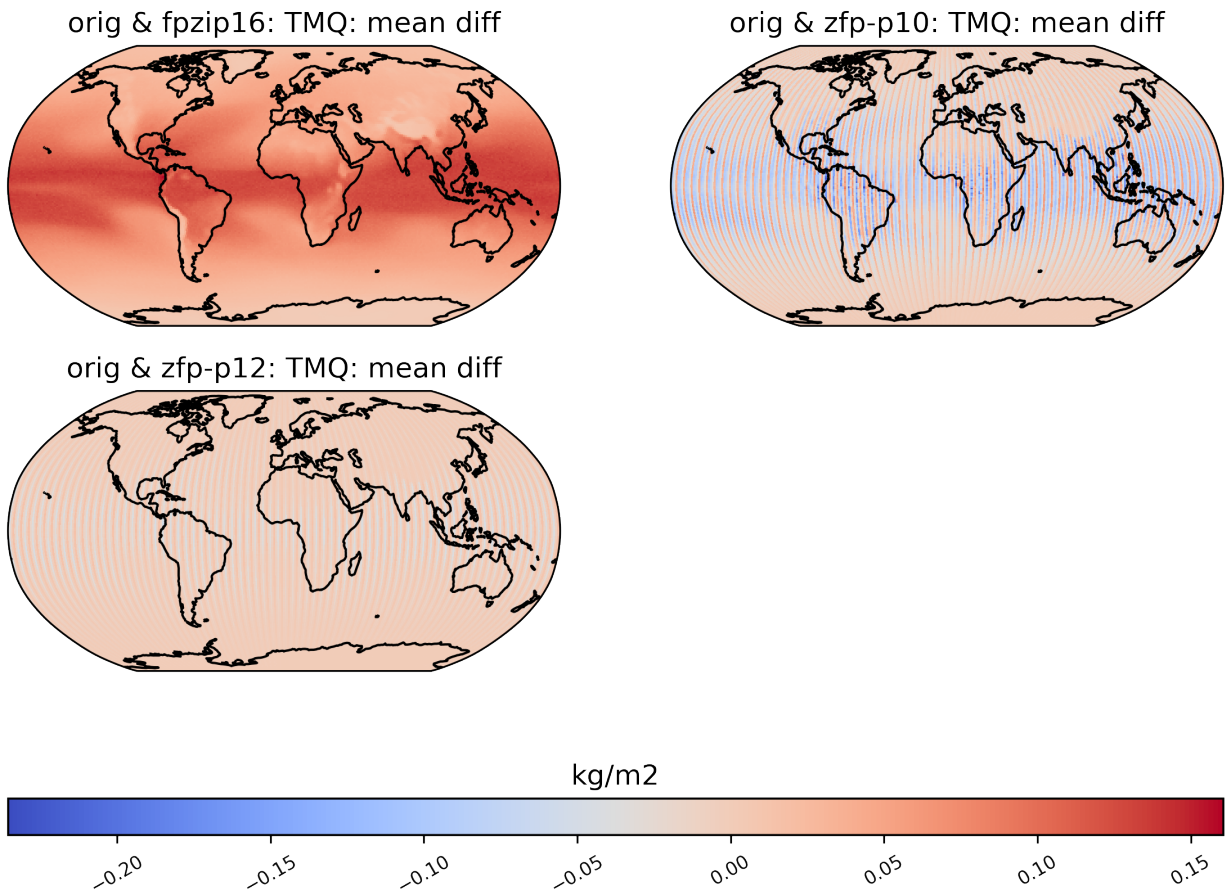
```

Let's look at something more interesting...

```

[15]: # diff between mean TS values in col_ds "orig" and "zfpA1.0" datasets
ldcpy.plot(
  col_tmq,
  "TMQ",
  sets=["orig", "fpzip16", "zfp-p10", "zfp-p12"],
  calc="mean",
  calc_type="diff",
)

```

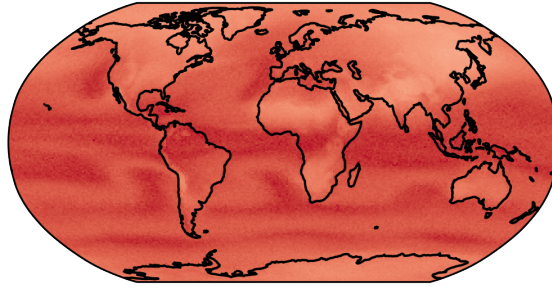


```

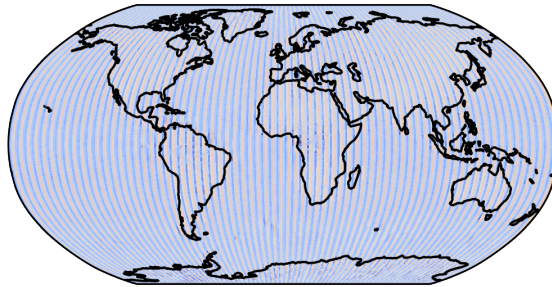
[16]: # plot of z-score under null hypothesis that orig value= compressed value
ldcpy.plot(
  col_tmq,
  "TMQ",
  sets=["orig", "fpzip16", "zfp-p10", "zfp-p12"],
  calc="zscore",
  calc_type="metric_of_diff",
)

```

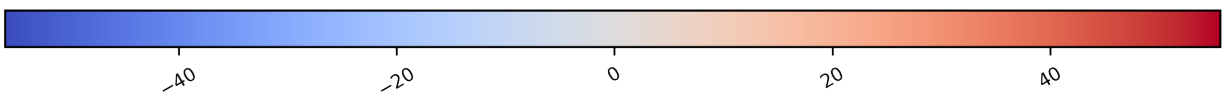
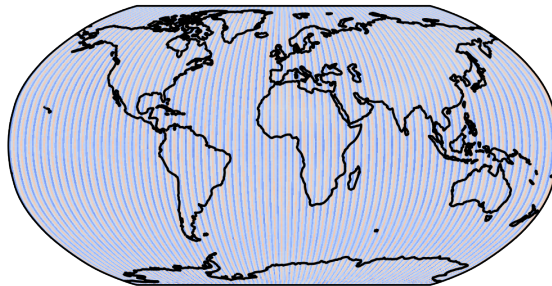
orig & fpzip16: TMQ: zscore: cutoff inf, % sig: 100.00 metric\_of\_diff



orig & zfp-p10: TMQ: zscore: cutoff inf, % sig: 79.52 metric\_of\_diff

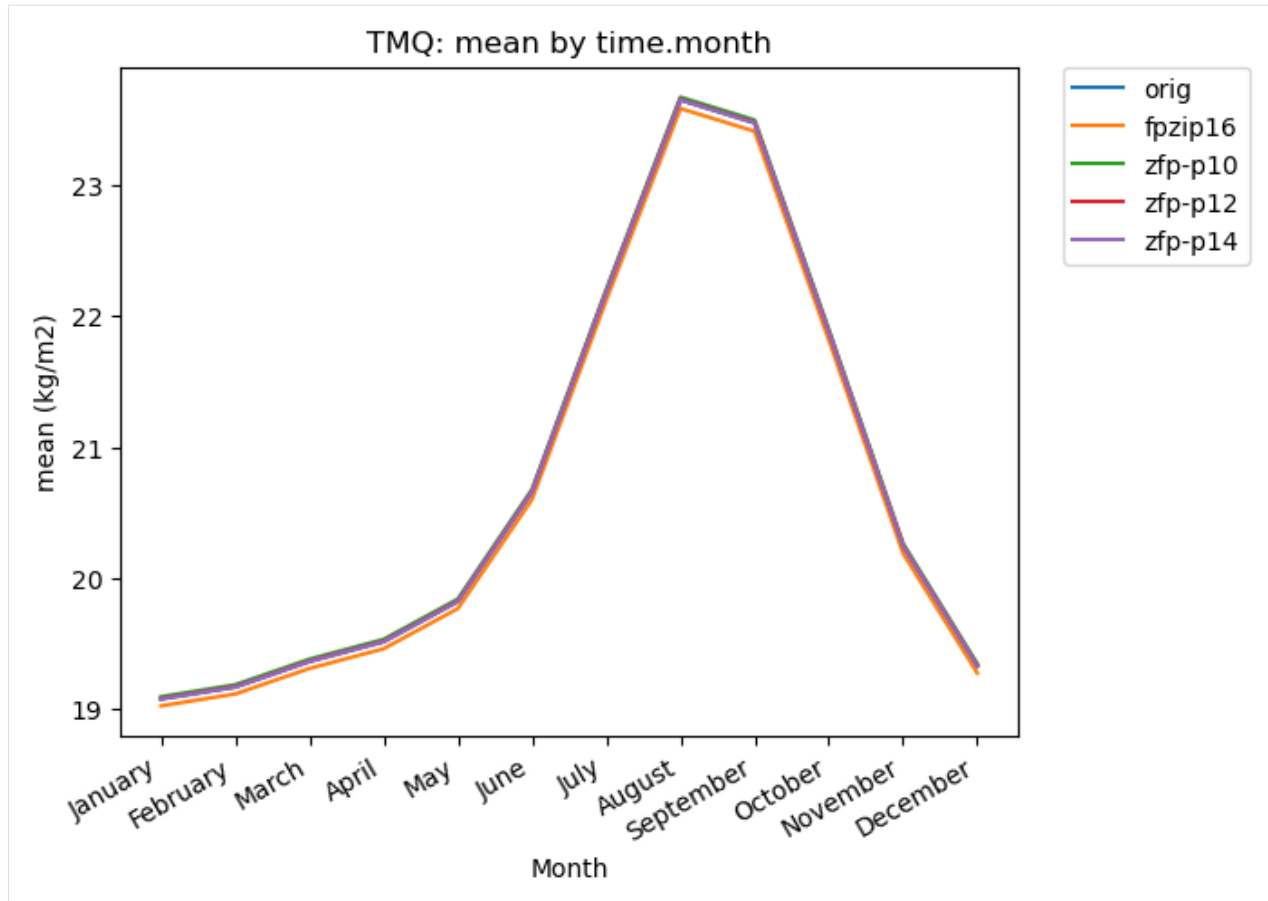


orig & zfp-p12: TMQ: zscore: cutoff inf, % sig: 73.63 metric\_of\_diff



```
[17]: # Time-series plot of TMQ mean in the original and lossy datasets
ldcpy.plot(
    col_tmq,
    "TMQ",
    sets=["orig", "fpzip16", "zfp-p10", "zfp-p12", "zfp-p14"],
    calc="mean",
    plot_type="time_series",
    group_by="time.month",
)
```





Now we look at other variables below

```
[18]: # List the compressed data files for TREFHT
os.listdir("/glade/p/cisl/asap/abaker/compression_samples/cam-lens/TREFHT")

[18]: ['zfp.p22.TREFHT.monthly.200601-208012.nc',
'zfp.p16.TREFHT.monthly.200601-208012.nc',
'zfp.p14.TREFHT.monthly.200601-208012.nc',
'zfp.p16.TREFHT.monthly.192001-200512.nc',
'zfp.p10.TREFHT.monthly.192001-200512.nc',
'zfp.p14.TREFHT.monthly.192001-200512.nc',
'zfp.p20.TREFHT.monthly.200601-208012.nc',
'zfp.p10.TREFHT.monthly.200601-208012.nc',
'zfp.p18.TREFHT.monthly.192001-200512.nc',
'zfp.p8.TREFHT.monthly.192001-200512.nc',
'zfp.p22.TREFHT.monthly.192001-200512.nc',
'fpzip20.TREFHT.monthly.200601-208012.nc',
'zfp.p12.TREFHT.monthly.200601-208012.nc',
'zfp.p20.TREFHT.monthly.192001-200512.nc',
'zfp.p12.TREFHT.monthly.192001-200512.nc',
'zfp.p18.TREFHT.monthly.200601-208012.nc',
'fpzip20.TREFHT.monthly.192001-200512.nc',
'zfp.p8.TREFHT.monthly.200601-208012.nc']
```

```
[19]: col_trefht = ldcpy.open_datasets(
    ["TREFHT"],
    [
        "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/orig/TREFHT.monthly.
↪200601-208012.nc",
        "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/TREFHT/fpzip20.TREFHT.
↪monthly.200601-208012.nc",
        "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/TREFHT/zfp.p8.TREFHT.
↪monthly.200601-208012.nc",
        "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/TREFHT/zfp.p10.TREFHT.
↪monthly.200601-208012.nc",
        "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/TREFHT/zfp.p12.TREFHT.
↪monthly.200601-208012.nc",
        "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/TREFHT/zfp.p14.TREFHT.
↪monthly.200601-208012.nc",
        "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/TREFHT/zfp.p16.TREFHT.
↪monthly.200601-208012.nc",
        "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/TREFHT/zfp.p18.TREFHT.
↪monthly.200601-208012.nc",
        "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/TREFHT/zfp.p20.TREFHT.
↪monthly.200601-208012.nc",
        "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/TREFHT/zfp.p22.TREFHT.
↪monthly.200601-208012.nc",
    ],
    [
        "orig",
        "fpzip20",
        "zfp-p8",
        "zfp-p10",
        "zfp-p12",
        "zfp-p14",
        "zfp-p16",
        "zfp-p18",
        "zfp-p20",
        "zfp-p22",
    ],
    chunks={"time": 700},
)
col_trefht
```

dataset size in GB 1.99

```
[19]: <xarray.Dataset>
Dimensions:      (collection: 10, lat: 192, lon: 288, time: 900)
Coordinates:
  * lat          (lat) float64 -90.0 -89.06 -88.12 -87.17 ... 88.12 89.06 90.0
  * lon          (lon) float64 0.0 1.25 2.5 3.75 5.0 ... 355.0 356.2 357.5 358.8
  * time         (time) object 2006-02-01 00:00:00 ... 2081-01-01 00:00:00
  * collection   (collection) <U7 'orig' 'fpzip20' ... 'zfp-p20' 'zfp-p22'
Data variables:
  TREFHT         (collection, time, lat, lon) float32 dask.array<chunksize=(1, 700,
↪192, 288), meta=np.ndarray>
Attributes:
  Conventions:   CF-1.0
  source:        CAM
  case:          b.e11.BRCP85C5CNBDRD.f09_g16.031
  title:         UNSET
```

(continues on next page)

(continued from previous page)

```

logname:      mickelso
host:         ys1023
Version:      $Name$
revision_Id:  $Id$
initial_file: b.e11.B20TRC5CNBDRD.f09_g16.031.cam.i.2006-01-01-00000.nc
topography_file: /glade/p/cesmdata/cseg/inputdata/atm/cam/topo/USGS-gtop...
history:      Tue Nov  3 14:11:24 2020: ncks -L 5 TREFHT.monthly.2006...
NCO:         netCDF Operators version 4.7.9 (Homepage = http://nco.s...

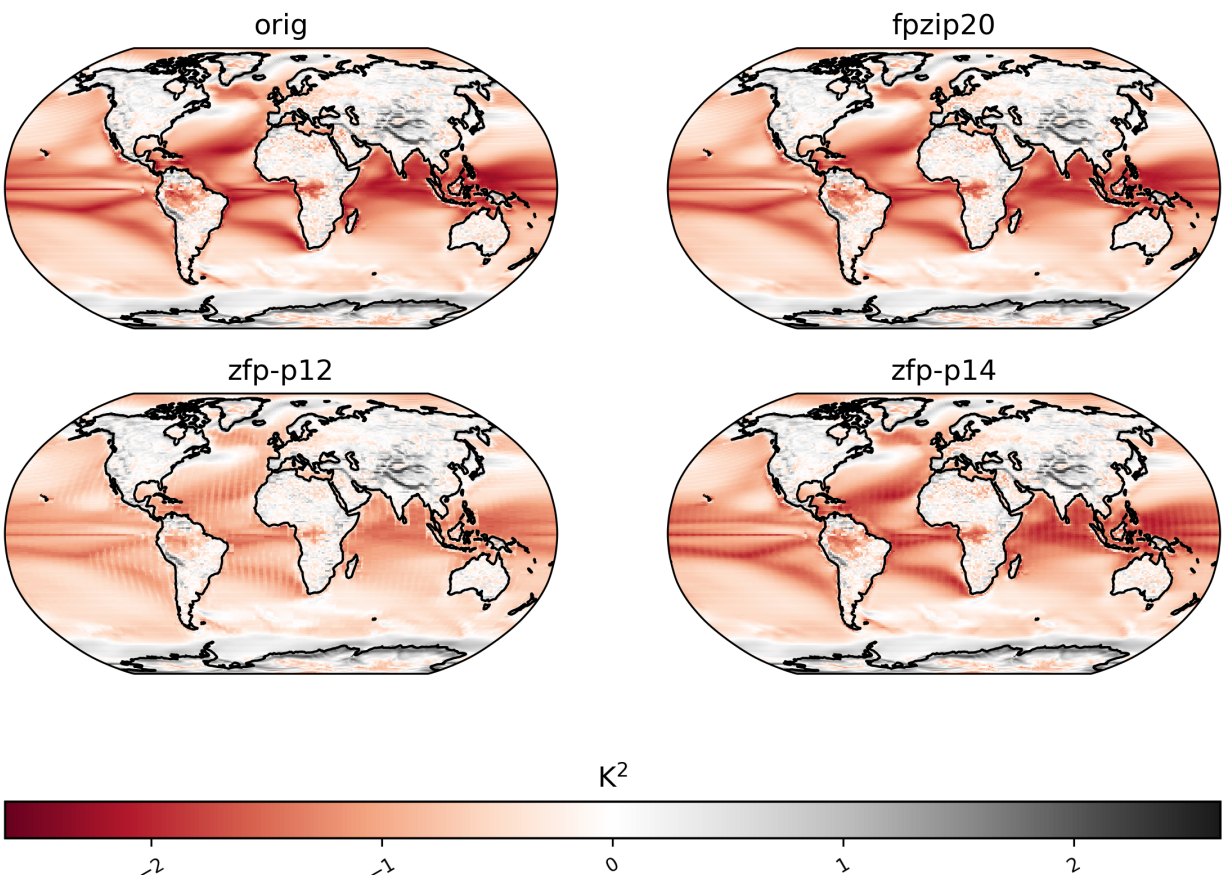
```

```
[20]: # Let's look and compare the north south contrast variances
```

```

ldcpy.plot(
    col_trefht,
    "TREFHT",
    sets=["orig", "fpzip20", "zfp-p12", "zfp-p14"],
    calc="ns_con_var",
    color="RdGy",
    calc_type="raw",
    transform="log",
    axes_symmetric=True,
    short_title=True,
)

```



```
[21]: # List the compressed data files for LHFLX
```

```
os.listdir("/glade/p/cisl/asap/abaker/compression_samples/cam-lens/LHFLX")
```

```
[21]: ['zfp.p14.LHFLX.daily.19200101-20051231.nc',
      'zfp.p12.LHFLX.daily.19200101-20051231.nc',
      'zfp.p18.LHFLX.daily.20060101-20801231.nc',
      'zfp.p12.LHFLX.daily.20060101-20801231.nc',
      'zfp.p16.LHFLX.daily.20060101-20801231.nc',
      'zfp.p22.LHFLX.daily.20060101-20801231.nc',
      'fpzip16.LHFLX.daily.19200101-20051231.nc',
      'zfp.p10.LHFLX.daily.19200101-20051231.nc',
      'fpzip16.LHFLX.daily.20060101-20801231.nc',
      'zfp.p8.LHFLX.daily.19200101-20051231.nc',
      'zfp.p8.LHFLX.daily.20060101-20801231.nc',
      'zfp.p20.LHFLX.daily.20060101-20801231.nc',
      'zfp.p14.LHFLX.daily.20060101-20801231.nc',
      'zfp.p10.LHFLX.daily.20060101-20801231.nc',
      'zfp.p16.LHFLX.daily.19200101-20051231.nc',
      'zfp.p22.LHFLX.daily.19200101-20051231.nc',
      'zfp.p18.LHFLX.daily.19200101-20051231.nc',
      'zfp.p20.LHFLX.daily.19200101-20051231.nc']
```

```
[22]: col_lhflx = ldcpy.open_datasets(
      ["LHFLX"],
      [
          "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/orig/LHFLX.daily.
↪20060101-20801231.nc",
          "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/LHFLX/fpzip16.LHFLX.
↪daily.20060101-20801231.nc",
          "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/LHFLX/zfp.p8.LHFLX.
↪daily.20060101-20801231.nc",
          "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/LHFLX/zfp.p10.LHFLX.
↪daily.20060101-20801231.nc",
          "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/LHFLX/zfp.p12.LHFLX.
↪daily.20060101-20801231.nc",
          "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/LHFLX/zfp.p14.LHFLX.
↪daily.20060101-20801231.nc",
          "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/LHFLX/zfp.p16.LHFLX.
↪daily.20060101-20801231.nc",
          "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/LHFLX/zfp.p18.LHFLX.
↪daily.20060101-20801231.nc",
          "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/LHFLX/zfp.p20.LHFLX.
↪daily.20060101-20801231.nc",
          "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/LHFLX/zfp.p22.LHFLX.
↪daily.20060101-20801231.nc",
      ],
      [
          "orig",
          "fpzip16",
          "zfp-p8",
          "zfp-p10",
          "zfp-p12",
          "zfp-p14",
          "zfp-p16",
          "zfp-p18",
          "zfp-p20",
          "zfp-p22",
      ],
      chunks={"time": 100},
  )
```

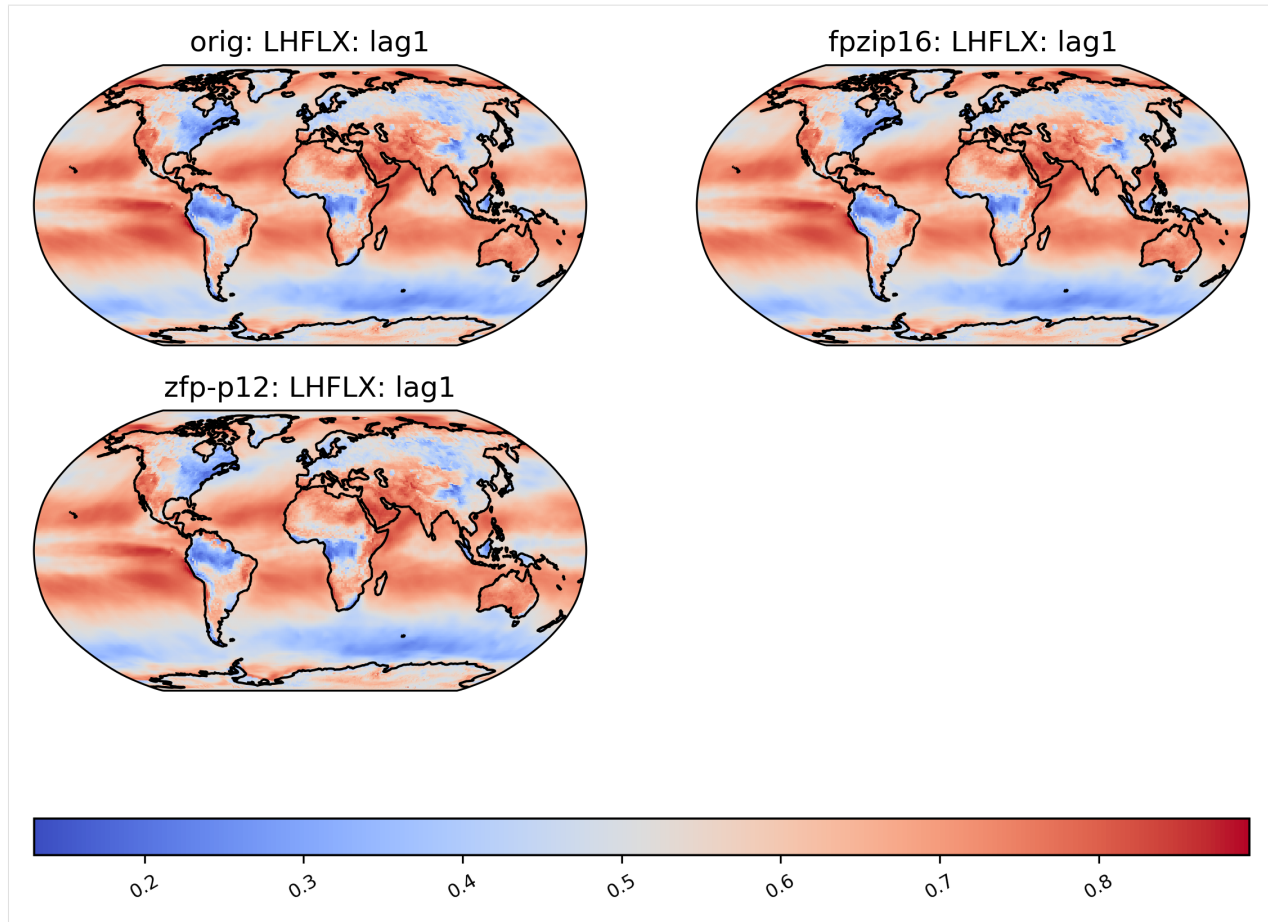
(continues on next page)

(continued from previous page)

```
col_lhflx
dataset size in GB 60.55

[22]: <xarray.Dataset>
Dimensions:      (collection: 10, lat: 192, lon: 288, time: 27375)
Coordinates:
  * lat          (lat) float64 -90.0 -89.06 -88.12 -87.17 ... 88.12 89.06 90.0
  * lon          (lon) float64 0.0 1.25 2.5 3.75 5.0 ... 355.0 356.2 357.5 358.8
  * time         (time) object 2006-01-01 00:00:00 ... 2080-12-31 00:00:00
  * collection   (collection) <U7 'orig' 'fpzip16' ... 'zfp-p20' 'zfp-p22'
Data variables:
  LHFLX          (collection, time, lat, lon) float32 dask.array<chunksize=(1, 100, ↵
↵192, 288), meta=np.ndarray>
Attributes:
  Conventions:    CF-1.0
  source:         CAM
  case:           b.e11.BRCP85C5CNBDRD.f09_g16.031
  title:          UNSET
  logname:        mickelso
  host:           ys1023
  Version:        $Name$
  revision_Id:    $Id$
  initial_file:   b.e11.B20TRC5CNBDRD.f09_g16.031.cam.i.2006-01-01-00000.nc
  topography_file: /glade/p/cesmdata/cseg/inputdata/atm/cam/topo/USGS-gtop...
  history:        Wed Nov 11 19:15:10 2020: ncks -L 5 LHFLX.daily.2006010...
  NCO:            netCDF Operators version 4.7.9 (Homepage = http://nco.s...
```

```
[26]: # plot of lag-1 correlation of LHFLX values for the first 10 years (NOTE: daily data ↵
↵takes longer)
ldcpy.plot(
    col_lhflx,
    "LHFLX",
    sets=["orig", "fpzip16", "zfp-p12"],
    calc="lag1",
    start=0,
    end=3650,
)
```



```
[27]: # List the compressed data files for Q
os.listdir("/glade/p/cisl/asap/abaker/compression_samples/cam-lens/Q")
```

```
[27]: ['zfp.p18.Q.monthly.192001-200512.nc',
      'zfp.p10.Q.monthly.200601-208012.nc',
      'zfp.p16.Q.monthly.200601-208012.nc',
      'zfp.p20.Q.monthly.200601-208012.nc',
      'zfp.p8.Q.monthly.200601-208012.nc',
      'zfp.p10.Q.monthly.192001-200512.nc',
      'zfp.p22.Q.monthly.200601-208012.nc',
      'fpzip20.Q.monthly.192001-200512.nc',
      'fpzip20.Q.monthly.200601-208012.nc',
      'zfp.p18.Q.monthly.200601-208012.nc',
      'zfp.p12.Q.monthly.200601-208012.nc',
      'zfp.p20.Q.monthly.192001-200512.nc',
      'zfp.p22.Q.monthly.192001-200512.nc',
      'zfp.p14.Q.monthly.200601-208012.nc',
      'zfp.p8.Q.monthly.192001-200512.nc',
      'zfp.p14.Q.monthly.192001-200512.nc',
      'zfp.p12.Q.monthly.192001-200512.nc',
      'zfp.p16.Q.monthly.192001-200512.nc']
```

```
[28]: col_q = ldcpy.open_datasets(
      ["Q"],
      [
```

(continues on next page)

(continued from previous page)

```

        "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/orig/Q.monthly.200601-
↪208012.nc",
        "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/Q/fpzip20.Q.monthly.
↪200601-208012.nc",
        "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/Q/zfp.p8.Q.monthly.
↪200601-208012.nc",
        "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/Q/zfp.p10.Q.monthly.
↪200601-208012.nc",
        "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/Q/zfp.p12.Q.monthly.
↪200601-208012.nc",
        "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/Q/zfp.p14.Q.monthly.
↪200601-208012.nc",
        "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/Q/zfp.p16.Q.monthly.
↪200601-208012.nc",
        "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/Q/zfp.p18.Q.monthly.
↪200601-208012.nc",
        "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/Q/zfp.p20.Q.monthly.
↪200601-208012.nc",
        "/glade/p/cisl/asap/abaker/compression_samples/cam-lens/Q/zfp.p22.Q.monthly.
↪200601-208012.nc",
    ],
    [
        "orig",
        "fpzip20",
        "zfp-p8",
        "zfp-p10",
        "zfp-p12",
        "zfp-p14",
        "zfp-p16",
        "zfp-p18",
        "zfp-p20",
        "zfp-p22",
    ],
    chunks={"time": 100},
)
col_q

```

dataset size in GB 59.72

```

[28]: <xarray.Dataset>
Dimensions:      (collection: 10, lat: 192, lev: 30, lon: 288, time: 900)
Coordinates:
  * lat          (lat) float64 -90.0 -89.06 -88.12 -87.17 ... 88.12 89.06 90.0
  * lev          (lev) float64 3.643 7.595 14.36 24.61 ... 957.5 976.3 992.6
  * lon          (lon) float64 0.0 1.25 2.5 3.75 5.0 ... 355.0 356.2 357.5 358.8
  * time         (time) object 2006-02-01 00:00:00 ... 2081-01-01 00:00:00
  * collection   (collection) <U7 'orig' 'fpzip20' ... 'zfp-p20' 'zfp-p22'
Data variables:
  Q              (collection, time, lev, lat, lon) float32 dask.array<chunksize=(1,
↪100, 30, 192, 288), meta=np.ndarray>
Attributes:
  Conventions:    CF-1.0
  source:         CAM
  case:           b.e11.BRCP85C5CNBDRD.f09_g16.031
  title:          UNSET
  logname:        mickelso
  host:           ys1023

```

(continues on next page)



(continued from previous page)

```

Version:      $Name$
revision_Id:  $Id$
initial_file: b.e11.B20TRC5CNBDRD.f09_g16.031.cam.i.2006-01-01-00000.nc
topography_file: /glade/p/cesmdata/cseg/inputdata/atm/cam/topo/USGS-gtop...
history:      Wed Nov 11 17:21:59 2020: ncks -L 5 Q.monthly.200601-20...
NCO:          netCDF Operators version 4.7.9 (Homepage = http://nco.s...

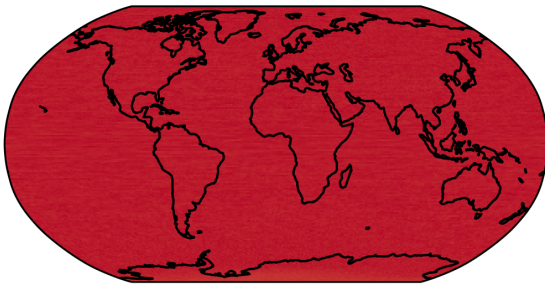
```

```

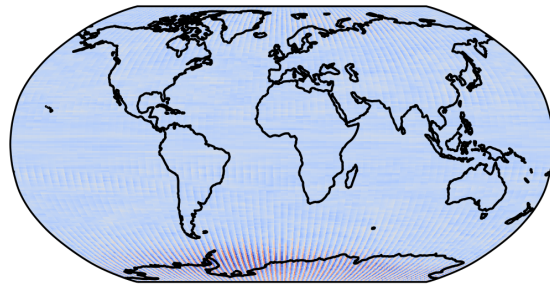
[29]: # diff between mean Q values across the entire timeseries
ldcpy.plot(
    col_q, "Q", sets=["orig", "fpzip20", "zfp-p12", "zfp-p14"], calc="mean", calc_
    ↪type="diff", lev=0
)

```

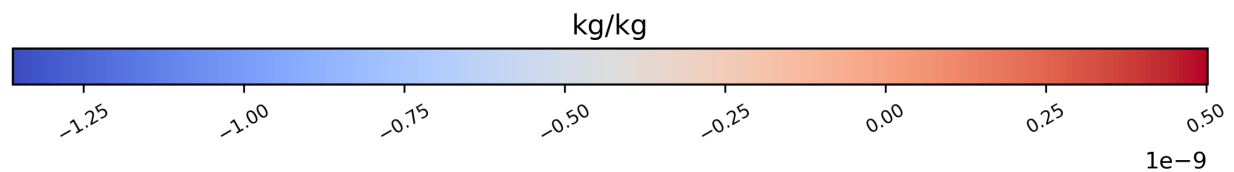
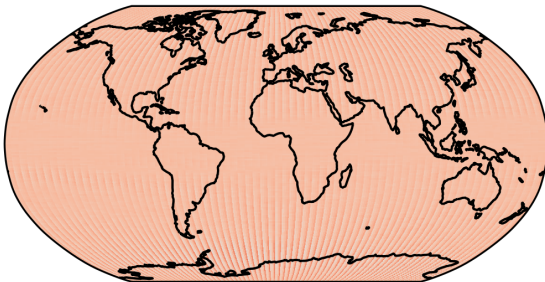
orig &amp; fpzip20: Q: mean diff



orig &amp; zfp-p12: Q: mean diff



orig &amp; zfp-p14: Q: mean diff

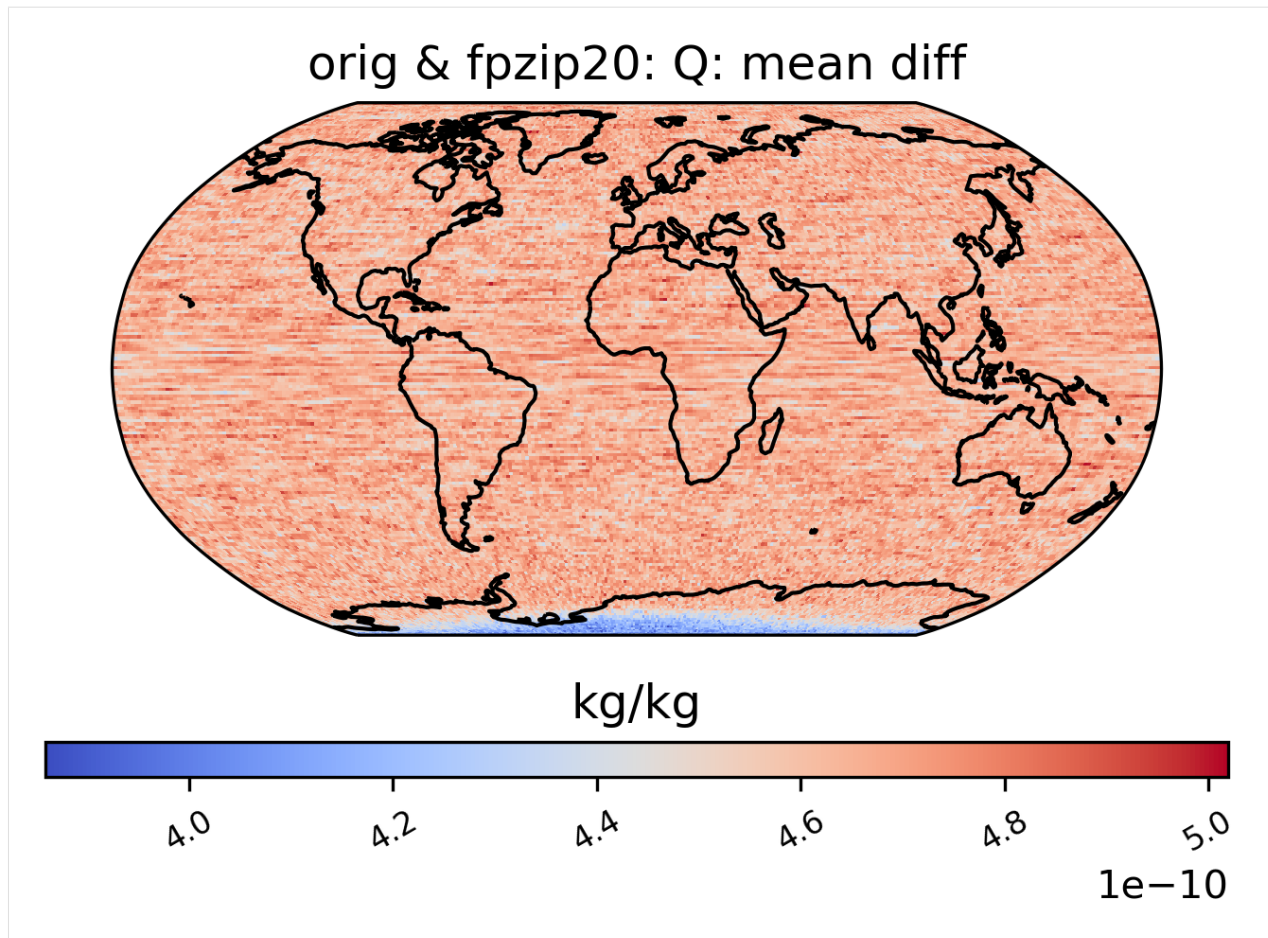


```

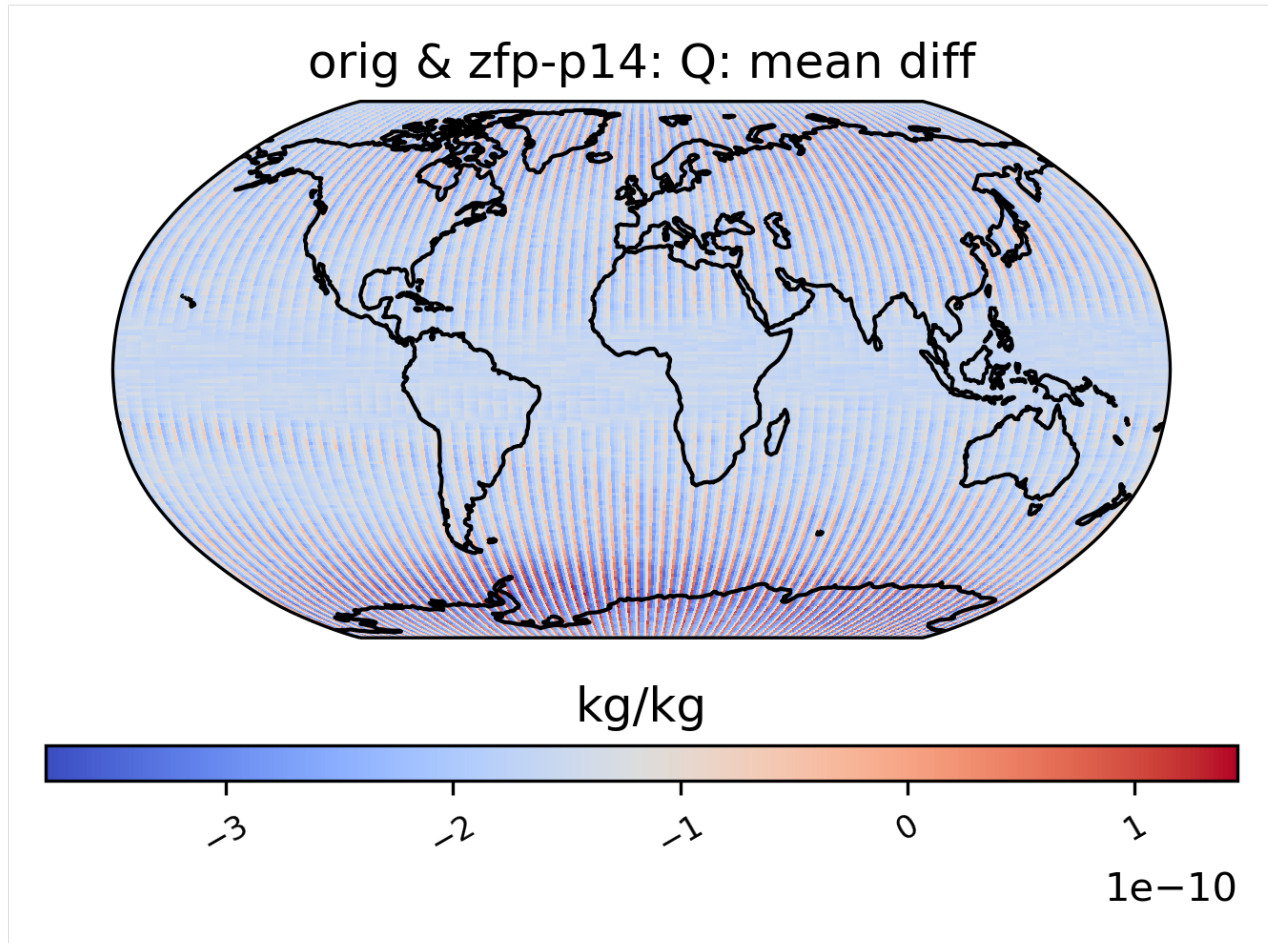
[30]: # diff between mean Q values across the entire timeseries - look at just one from_
    ↪above
ldcpy.plot(col_q, "Q", sets=["orig", "fpzip20"], calc="mean", calc_type="diff", lev=0)

```





```
[31]: # diff between mean Q values across the entire timeseries - look at just one from_
      ↪ above
ldcpy.plot(col_q, "Q", sets=["orig", "zfp-pl4"], calc="mean", calc_type="diff", lev=0)
```



```
[35]: # Note: since q is 3D, need to select a level and a time slice
ldcpy.compare_stats(col_q.isel(time=0, lev=0), "Q", "orig", "fpzip20")
```

```
mean orig           : 2.2156e-06
mean fpzip20        : 2.2151e-06
mean diff           : 4.6242e-10

variance orig       : 9.5578e-15
variance fpzip20    : 9.5584e-15

standard deviation orig : 9.7765e-08
standard deviation fpzip20 : 9.7768e-08

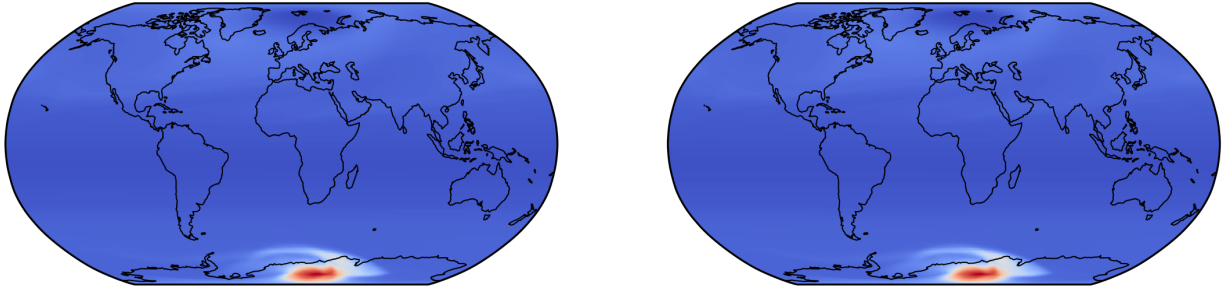
max value orig      : 3.2387e-06
max value fpzip20   : 3.2382e-06
min value orig      : 2.1375e-06
min value fpzip20   : 2.1374e-06

max abs diff        : 9.311e-10
min abs diff        : 0
mean abs diff       : 4.6242e-10
mean squared diff   : 2.1383e-19
root mean squared diff : 5.351e-10
normalized root mean squared diff : 0.0004859
```

(continues on next page)

(continued from previous page)

```
normalized max pointwise error      : 0.00084549
pearson correlation coefficient      : 1
ks p-value                          : 6.4126e-18
spatial relative error(% > 0.0001) : 75.96
max spatial relative error          : 0.00043437
ssim                                : 0.99865
ssim_fp                             : 0.96366
```



```
[36]: # Disconnect when finished
      cluster.close()
      client.close()
```

```
[ ]:
```



## INDICES AND TABLES

- `genindex`
- `modindex`
- `search`



## A

`annual_harmonic_relative_ratio()` (*ldcpy.metrics.DatasetMetrics* property), 19  
`annual_harmonic_relative_ratio_pct_sig()` (*ldcpy.metrics.DatasetMetrics* property), 19

## C

`check_metrics()` (in module *ldcpy.util*), 13  
`collect_datasets()` (in module *ldcpy.util*), 14  
`compare_stats()` (in module *ldcpy.util*), 14  
`covariance()` (*ldcpy.metrics.DiffMetrics* property), 20

## D

*DatasetMetrics* (class in *ldcpy.metrics*), 18  
*DiffMetrics* (class in *ldcpy.metrics*), 20

## E

`ew_con_var()` (*ldcpy.metrics.DatasetMetrics* property), 19

## G

`get_diff_metric()` (*ldcpy.metrics.DiffMetrics* method), 20  
`get_metric()` (*ldcpy.metrics.DatasetMetrics* method), 18  
`get_single_metric()` (*ldcpy.metrics.DatasetMetrics* method), 19

## K

`ks_p_value()` (*ldcpy.metrics.DiffMetrics* property), 20

## L

`lag1()` (*ldcpy.metrics.DatasetMetrics* property), 19  
`lag1_first_difference()` (*ldcpy.metrics.DatasetMetrics* property), 19  
*ldcpy.metrics*  
 module, 18  
*ldcpy.plot*  
 module, 16

*ldcpy.util*  
 module, 13

## M

`mae_day_max()` (*ldcpy.metrics.DatasetMetrics* property), 19  
`max_spatial_rel_error()` (*ldcpy.metrics.DiffMetrics* property), 20  
`mean()` (*ldcpy.metrics.DatasetMetrics* property), 19  
`mean_abs()` (*ldcpy.metrics.DatasetMetrics* property), 19  
`mean_squared()` (*ldcpy.metrics.DatasetMetrics* property), 19  
*MetricsPlot* (class in *ldcpy.plot*), 16  
 module  
   *ldcpy.metrics*, 18  
   *ldcpy.plot*, 16  
   *ldcpy.util*, 13

## N

`normalized_max_pointwise_error()` (*ldcpy.metrics.DiffMetrics* property), 20  
`normalized_root_mean_squared()` (*ldcpy.metrics.DiffMetrics* property), 20  
`ns_con_var()` (*ldcpy.metrics.DatasetMetrics* property), 19

## O

`odds_positive()` (*ldcpy.metrics.DatasetMetrics* property), 19  
`open_datasets()` (in module *ldcpy.util*), 15

## P

`pearson_correlation_coefficient()` (*ldcpy.metrics.DiffMetrics* property), 21  
`plot()` (in module *ldcpy.plot*), 16  
`pooled_variance()` (*ldcpy.metrics.DatasetMetrics* property), 19  
`pooled_variance_ratio()` (*ldcpy.metrics.DatasetMetrics* property), 20  
`prob_negative()` (*ldcpy.metrics.DatasetMetrics* property), 20

`prob_positive()` (*ldcpy.metrics.DatasetMetrics*  
*property*), 20

## R

`root_mean_squared()` (*ld-*  
*cpy.metrics.DatasetMetrics property*), 20

## S

`spatial_rel_error()` (*ldcpy.metrics.DiffMetrics*  
*property*), 21

`ssim_value()` (*ldcpy.metrics.DiffMetrics property*),  
21

`ssim_value_fp()` (*ldcpy.metrics.DiffMetrics prop-*  
*erty*), 21

`ssim_value_fp_old()` (*ldcpy.metrics.DiffMetrics*  
*property*), 21

`standardized_mean()` (*ld-*  
*cpy.metrics.DatasetMetrics property*), 20

`std()` (*ldcpy.metrics.DatasetMetrics property*), 20

`subset_data()` (*in module ldcpy.util*), 15

## T

`tex_escape()` (*in module ldcpy.plot*), 18

`time_series_plot()` (*ldcpy.plot.MetricsPlot*  
*method*), 16

## V

`variance()` (*ldcpy.metrics.DatasetMetrics property*),  
20

## Z

`zscore()` (*ldcpy.metrics.DatasetMetrics property*), 20

`zscore_cutoff()` (*ldcpy.metrics.DatasetMetrics*  
*property*), 20

`zscore_percent_significant()` (*ld-*  
*cpy.metrics.DatasetMetrics property*), 20