

CHESS Course: An Introduction to FAIR Data Management for Geoscientists, Syllabus

Day 1 - Session 4

Documentation of data

Torill Hamre, Nansen Environmental and Remote Sensing Center

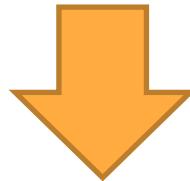


Session Topics

- Tools for documenting data
 - Rosetta (web application), NCO/CDO (command line), Python (netcdf4)
- Examples in Python
 - netcdf4, xarray
- Validation of metadata
 - What is actually validated?
- Validation tools for NetCDF/CF
 - NorDataNet validator, PUMA validator
- Rosetta in more detail
 - Profiles, time series, trajectory
 - Template concept and benefits
 - Converting a CTD profile from Seabird sensor
- Introduction to assignment

Learning Objectives

- Recognize the value of documenting your data
- Identify some tools for documenting data
- Understanding the difference between different categories of data
- Identify some tools for validating metadata
- Recognize the value of using standards



Assignment 1:
Documenting and formatting YOUR dataset
(alternative: provided data)

Preparing your data for publication

- Why is metadata important?
 - Documents the data
 - Content
 - Processing and quality control
 - Reference material
 - Contact for data provider
 - Data usage rights (license)
 - ...
- How ensure enough metadata?
 - Include in planning of data collection and processing



Preparing your data for publication

- Compile both metadata and other documentation
- Metadata
 - **Discovery metadata**, e.g. area covered, time period, parameters, ...
 - **Usage metadata**, e.g. processing history, units, data quality, ...
 - Avoid making new standards; utilize the existing ones
- Decide on suitable tools for formatting
 - Depends on formats used in processing and analysis tools
 - Select a **standard data format** and **metadata standard**
 - Data formats should be **self-describing**
 - Files should **include all needed metadata**

Preparing your data for publication

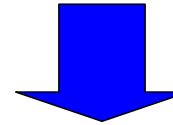
Why tools?

- To add documentation to dataset files

- **Discovery metadata**
- **Usage metadata**



```
WELL-18.CSV - Notepad
Data file for Datalogger.
COMPANY : <>company name>
COMP.STATUS : 
DATE : 03/06/2009
TIME : 09:54:45
FILENAME : C:\WELL-18.CSV
CREATED BY : SWS Diver-Office 2.0.1.2
[BEGINNING OF DATA]
[Logger settings]
Instrument type =Micro-Diver-15
Status =Started =0
Serial number =.00-C2999 215.
Instrument number =
Location =MW-18
Sample period =410
Sample method =T
Number of channels =2
[Channel 1]
Identification =WATER HEAD
Reference level =400.0 cm
Range =-1750.0 cm
Master level =0 m
Altitude =0 m
[Channel 2]
Identification =TEMPERATURE
Reference level =-20.00 °C
Range =100.00 °C
```

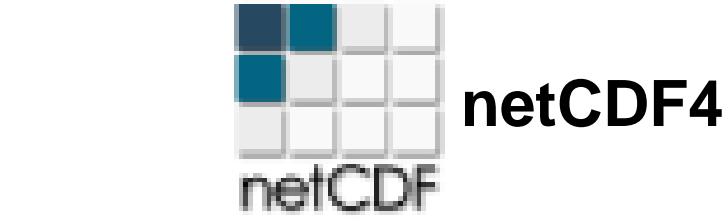
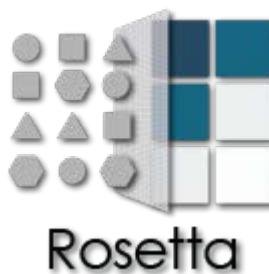


(data + metadata)

Tools for documenting metadata

What tools?

- GIS (commercial, free)
- Data processing and analysis tools
(commercial, free)
- Rosetta (free)

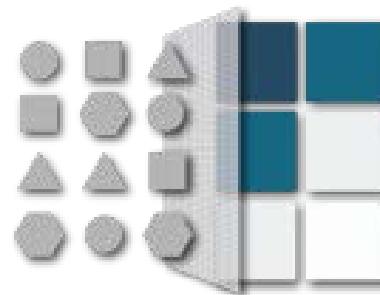


Tools for documenting metadata

Rosetta

Main features

- Read text files (CSV)
- Adds standard metadata
- Generates NetCDF
- Templates for conversion
- Web application



Tools for documenting metadata

NCO/CDO

Main features

- Reads, manipulates and stores NetCDF
- Add/modify metadata
- Remapping, subsetting
- Analysis (e.g. statistics)
- Supports NetCDF, DAP, HDF4, HDF5
- Command line



What is NCO?

The netCDF Operators (NCO) comprise about a dozen standalone, command-line programs that take [netCDF](#), [HDF](#), and/or [DAP](#) files as input, then operate (e.g., derive new fields, compute statistics, print, hyperslab, manipulate metadata, regrid) and output the results to screen or files in text, binary, or netCDF formats. NCO aids analysis of gridded and unstructured scientific data. The shell-command style of NCO allows users to manipulate and analyze files interactively, or with expressive scripts that avoid some overhead of higher-level programming environments.

Traditional geoscience data analysis requires users to work with numerous flat (data in one level or namespace) files. In that paradigm instruments or models produce, and then repositories archive and distribute, and then researchers request and analyze, collections of flat files. NCO works well with that paradigm, yet it also embodies the necessary algorithms to transition geoscience data analysis from relying solely on traditional (or “flat”) datasets to allowing newer hierarchical (or “nested”) datasets.

The next logical step is to support and enable combining all datastreams that meet user-specified criteria into a single or small number of files that hold *all* the science-relevant data organized in hierarchical structures. NCO (and no other software to our knowledge) can do this now. We call the resulting data storage, distribution, and analysis paradigm Group-Oriented Data Analysis and Distribution ([GODAD](#)). GODAD lets the scientific question organize the data, not the *ad hoc* granularity of all relevant datasets. The [User Guide](#) illustrates [GODAD](#) techniques for climate data analysis:

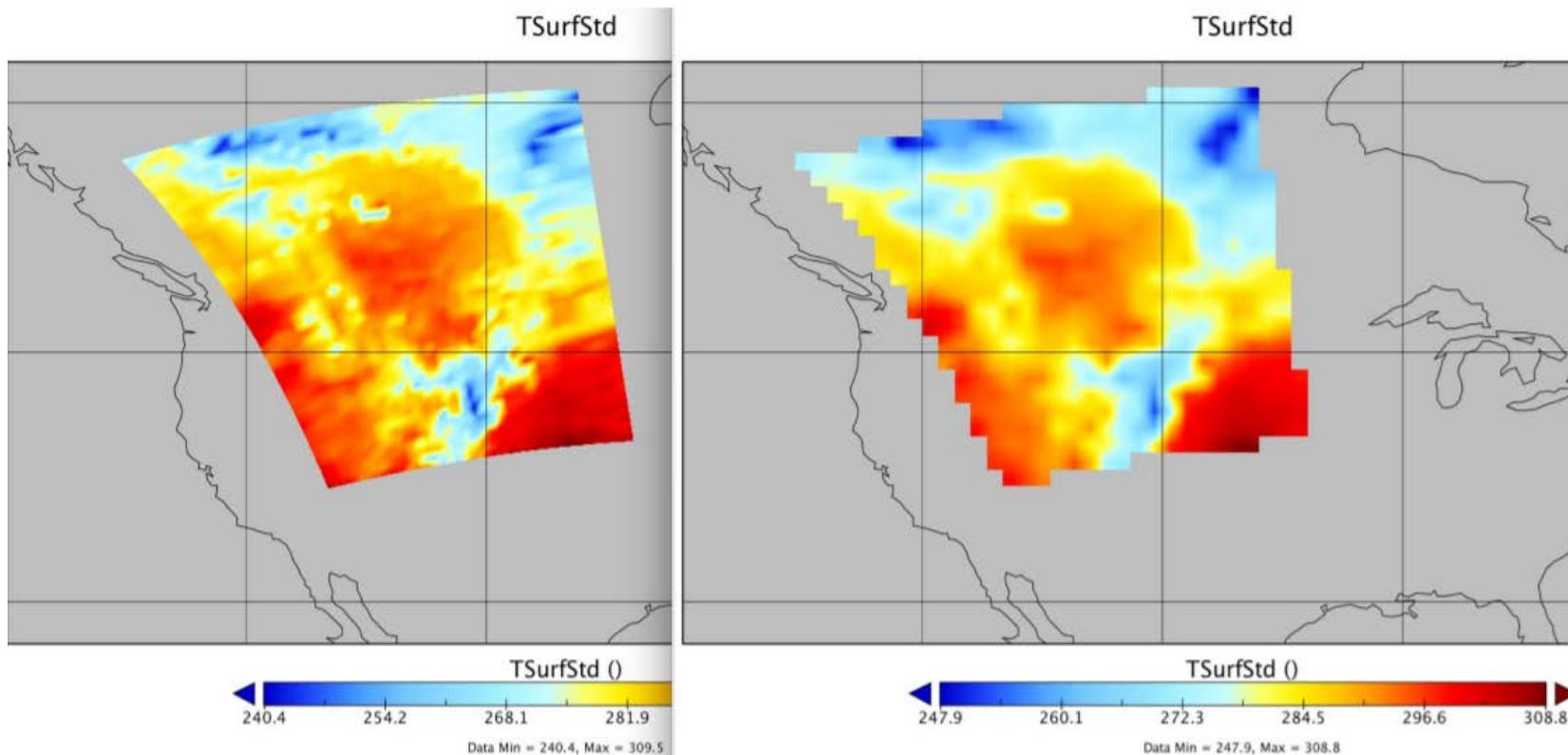
- [ncap2](#) netCDF Arithmetic Processor ([examples](#))
- [ncatted](#) netCDF ATTRIBUTE EDitor ([examples](#))
- [ncbo](#) netCDF Binary Operator (addition, multiplication...) ([examples](#))
- [ncclimo](#) netCDF CLIMatology Generator ([examples](#))
- [nces](#) netCDF Ensemble Statistics ([examples](#))
- [nccat](#) netCDF Ensemble concATenator ([examples](#))
- [ncflint](#) netCDF FiLe INTERpolator ([examples](#))
- [ncks](#) netCDF Kitchen Sink ([examples](#))
- [ncpdq](#) netCDF Permute Dimensions Quickly, Pack Data Quietly ([examples](#))
- [ncra](#) netCDF Record Averager ([examples](#))
- [ncrcat](#) netCDF Record concATenator ([examples](#))
- [ncremap](#) netCDF REMAPer ([examples](#))
- [ncrename](#) netCDF RENAMEer ([examples](#))
- [ncwa](#) netCDF Weighted Averager ([examples](#))

Note that the “averagers” ([ncra](#) and [ncwa](#)) are misnamed because they perform many non-linear statistics as well, e.g., total, minimum, RMS. Moreover, [ncap2](#) implements a powerful domain language which handles arbitrarily complex algebra, calculus, and statistics (using GSL). The operators are as general as [netCDF](#) itself: there are no restrictions on the contents of input file(s). NCO’s internal routines are completely dynamic and impose no limit on the number or sizes of dimensions, variables, and files. NCO is designed to be used both interactively and with large batch jobs. The default operator behavior is often sufficient for everyday needs, and there are numerous command line (i.e., run-time) options, for special cases.



Regrid NASA AIRS Level 2 Swath Data in raw HDF4 format from regional curvilinear 45x30 source grid to equiangular 1x1 degree:

```
% ncremap -s AIRS.2014.10.01.202.L2.RetStd.v6.0.11.0.G14275134307.hdf -d dst_1x1.nc
```



4.12 ncremap netCDF Remapper

SYNTAX

```
ncremap [-3] [-4] [-5] [-6] [-7] [-a alg_typ] [--a2o]
[-D dbg_lvl] [-d dst_fl] [-d2f] [-dpt] [--dpt_fl=dpt_fl]
[--dt_sng=dt_sng] [--esmf_typ=esmf_typ]
[--fl_fmt=fl_fmt] [-G grd_sng] [-g grd_dst]
[-I drc_in] [-i input-file] [-j job_nbr] [-L dfl_lvl]
[-M] [-m map_fl] [--msh_fl=msh_fl]
[--msk_dst=msk_dst] [--msk_out=msk_out] [--msk_src=msk_src] [--mss_val=mss_val]
[-n nco_opt] [--nm_dst=nm_dst] [--nm_src=nm_src]
[--no_cll_msr] [--no_frm_trm] [--no_stg_grd]
[-O drc_out] [-o output-file] [-P prc_typ] [-p par_typ]
[--pdq=pdq_opt] [-ppc=ppc_opt] [--preserve=prs_stt]
[-R rgr_opt] [--rgn_dst] [-rgn_src] [--rnr_thr=rnr_thr]
[--rrg_bb_wesn=bb_wesn] [-rrg_dat_glb=dat_glb] [--rrg_grd_glb=grd_glb]
[--rrg_grd_rgn=grd_rgn] [-rrg_rnm_sng=rnm_sng]
[-s grd_src] [--sgs_frc=sgs_frc] [--sgs_msk=sgs_msk] [--sgs_nrm=sgs_nrm]
[--skl=skl-file] [--stdin] [-T drc_tmp] [-t thr_nbr]
[-U] [-u unq_sfx] [--ugrid=ugrid-file] [--uio]
[-V rgr_var] [-v var_lst[...]] [--version] [--vrb=vrb_lvl]
[--vrt_fl=vrt_fl] [-vrt_ntp=vrt_ntp] [--vrt_xtr=vrt_xtr]
[-W wgt_opt] [-w wgt_cmd] [-x xtn_lst[...]] [--xcl_var]
[--xtr_nsp=xtr_nsp] [--xtr_xpn=xtr_xpn]
[input-files] [output-file]
```

DESCRIPTION

`ncremap` remaps the data file(s) in `input-file`, in `drc_in`, or piped through standard input, to the horizontal grid specified by (in descending order of precedence) `map_fl`, `grd_dst`, or `dst_fl` and stores the result in `output-file`(s). If a vertical grid `vrt_fl` is provided, `ncremap` will (also) vertically interpolate the input file(s) to that grid. When no `input-file` is provided, `ncremap` operates in “map-only” mode where it exits after producing an annotated map-file. `ncremap` was introduced to NCO in version 4.5.4 (December, 2015).

`ncremap` is a “super-operator” that orchestrates the regridding features of several different programs including other NCO operators. Under the hood NCO applies pre-computed remapping weights or, when necessary, generates and infers grids, generates remapping weights itself or calls external programs to generate the weights, and then applies the weights (i.e., regrids).

Tutorials:

- <https://code.mpimet.mpg.de/projects/cdo/wiki/Tutorial>
- <http://hannahlab.org/cdo-vs-nco/>

Tools for documenting metadata

Example – modify metadata with **ncatted** (<https://linux.die.net/man/1/ncatted>)

Append the string "Data version 2.0.\n" to the global attribute **history**:

```
ncatted -O -a history,global,a,c,"Data version 2.0\n" in.nc
```

Note the use of embedded C language **printf()**-style escape sequences.

Change the value of the **long_name** attribute for variable **T** from whatever it currently is to "temperature":

```
ncatted -O -a long_name,T,o,c,temperature in.nc
```

Delete all existing **units** attributes:

```
ncatted -O -a units,,d,, in.nc
```

The value of *var_nm* was left blank in order to select all variables in the file. The values of *att_type* and *att_val* were left blank because they are superfluous in *Delete* mode.

Modify all existing **units** attributes to "meter second-1"

```
ncatted -O -a units,,m,c,"meter second-1" in.nc
```

Overwrite the **quanta** attribute of variable **energy** to an array of four integers.

```
ncatted -O -a quanta,energy,o,s,"010,101,111,121" in.nc
```

Tools for documenting metadata

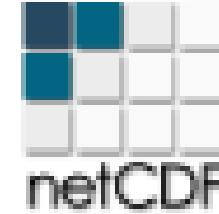
netcdf4 (netcdf4-python)

Main features

- Python package
- Read, write, update NetCDF
- Add/update metadata
- Supports netCDF3 and netCDF4

Documentation

- <https://unidata.github.io/netcdf4-python/netCDF4/index.html>
- http://schubert.atmos.colostate.edu/~cslocum/netcdf_example.html



netCDF4

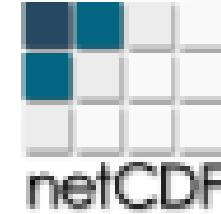
Tools for documenting metadata

netcdf4 (netcdf4-python)

Open file and inspect its structures, (some) variables and attributes

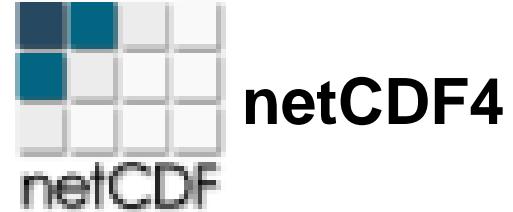
```
from netCDF4 import Dataset
from datetime import datetime
import sys
import numpy as np

filename = "https://thredds.met.no/thredds/dodsC/myocean/siw-tac/siw-metno-
svalbard/2020/05/ice_conc_svalbard_202005151500.nc"
simap = Dataset(filename, "r")
print(simap)
```



netCDF4

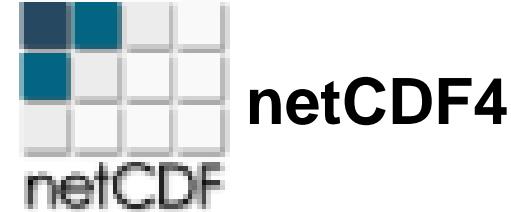
Tools for documenting metadata



Output:

```
<class 'netCDF4._netCDF4.Dataset'>
root group (NETCDF3_CLASSIC data model, file format DAP2):
    title: Arctic Svalbard & Barents Ice Concentration, L4, 1km daily (METNO-ARC-
SEAICE_CONC-L4-NRT-OBS)
    Conventions: CF-1.4
    netcdf_version_id: 3.6.3
    creation_date: 2020-05-15T13:36:20Z
    ...
    dimensions(sizes): maxStrlen64(64), time(1), xc(3812), yc(2980)
    variables(dimensions): int32 time(time), float32 yc(yc), float32 xc(xc), |S1
crs(maxStrlen64), float32 lat(yc,xc), float32 lon(yc,xc), int16
ice_concentration(time,yc,xc), int16 concentration_range(time,yc,xc)
    groups:
```

Tools for documenting metadata



Inspecting metadata for one of the variables:

```
conc = simap.variables['ice_concentration']
print(conc)
```

Output:

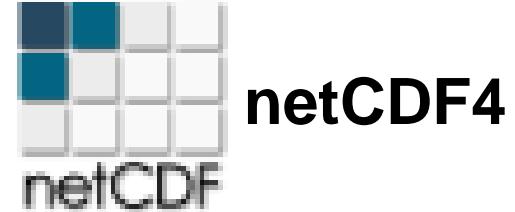
```
<class 'netCDF4._netCDF4.Variable'>
int16 ice_concentration(time, yc, xc) ← The dimensions (axes) of the data variable
    long_name: sea ice concentration
    standard_name: sea_ice_area_fraction
    units: %
    coordinates: lon lat
    grid_mapping: crs
    source: met.no
    _FillValue: -99
    scale_factor: 1.0
    add_offset: 0.0
unlimited dimensions:
current shape = (1, 2980, 3812) ← The actual dimensions of the data variable
filling off
```

The dimensions (axes) of the data variable

The data variable's attributes

The actual dimensions of the data variable

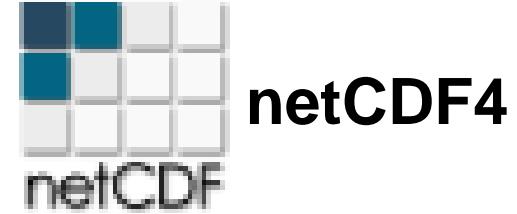
Tools for documenting metadata



Adding a new attribute (metadata) to this variable:

```
dset = Dataset('test-newattr.nc','w', format='NETCDF3_CLASSIC')
dset.createDimension('time',size=simap.variables["ice_concentration"].shape[0])
dset.createDimension('yc',size=simap.variables["ice_concentration"].shape[1])
dset.createDimension('xc',size=simap.variables["ice_concentration"].shape[2])
newconc = dset.createVariable('sea_ice_concentration', np.int16, dimensions=('time', 'yc',
'xc'), fill_value=-99)
newconc[:] = conc[:] ← Copy all data
for name in conc.ncattrs():
    if name != '_FillValue':
        newconc.setncattr(name, getattr(conc,name)) } Copy all variable's attributes from original file
newconc.setncattr('source_full_name', 'Meteorological Institute of Norway')
print(newconc) → Create a new attribute
```

Tools for documenting metadata

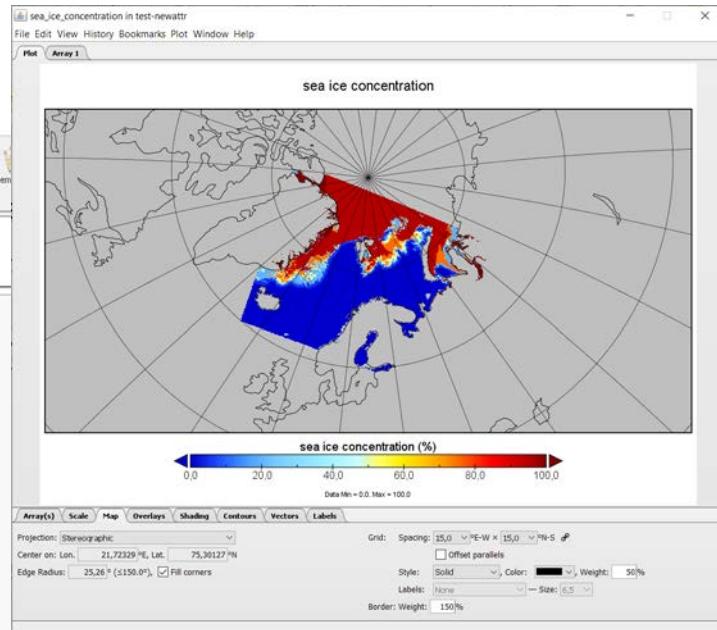
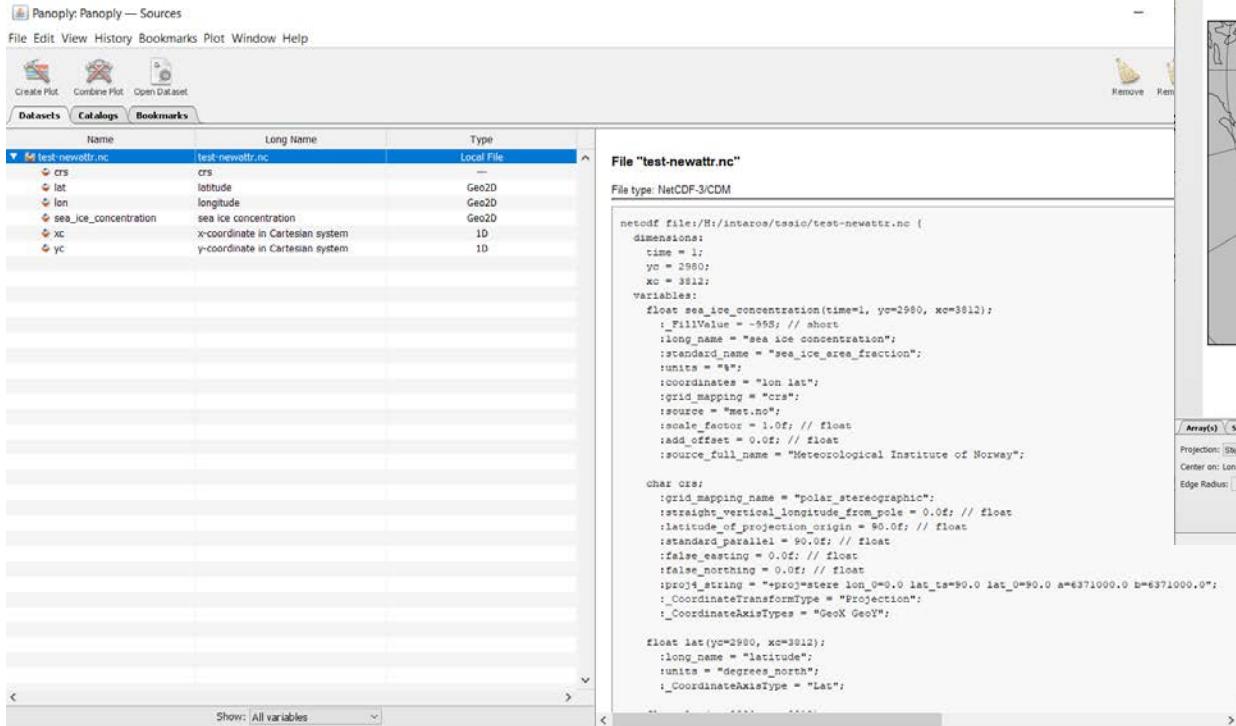


Output:

```
<class 'netCDF4._netCDF4.Variable'>
int16 sea_ice_concentration(time, yc, xc)
    _FillValue: -99
    long_name: sea ice concentration
    standard_name: sea_ice_area_fraction
    units: %
    coordinates: lon lat
    grid_mapping: crs
    source: met.no
    scale_factor: 1.0
    add_offset: 0.0
    source_full_name: Meteorological Institute of Norway<----- New attribute added
unlimited dimensions:
current shape = (1, 2980, 3812)
filling on
```

Tools for documenting metadata

- After copying other variables & metadata



Tools for documenting metadata



xarray

Main features

- Python NetCDF file I/O
- Advanced analysis tools
- Multi-dimensional, labelled arrays
- Split-Apply-Combine
- Supports parallel computing

Documentation

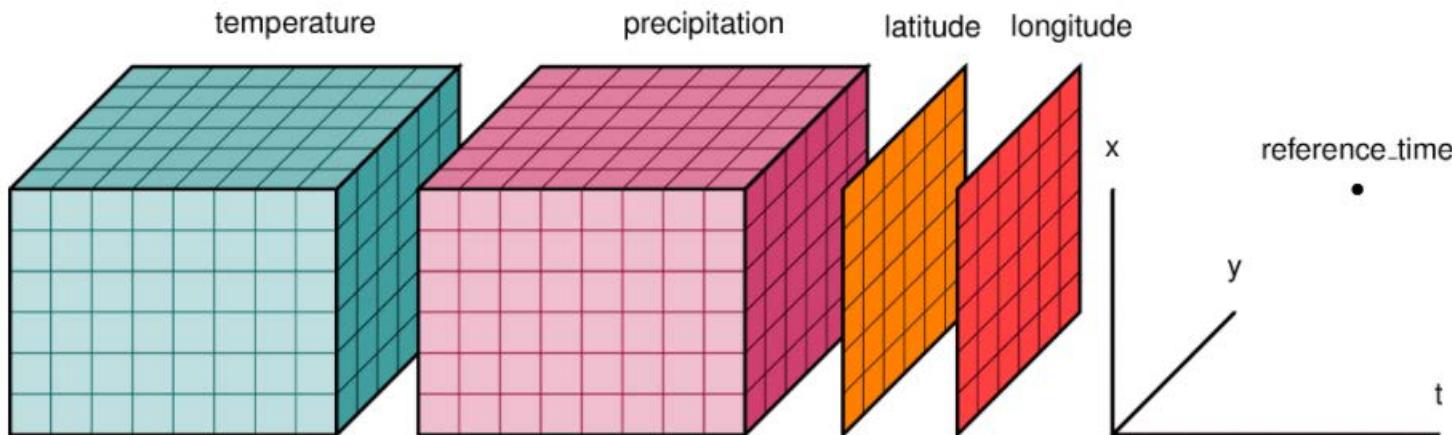
- <http://xarray.pydata.org/en/stable/>

Tools for documenting metadata



xarray

Tutorial: https://rabernat.github.io/research_computing/xarray.html



Tools for documenting metadata



Install need packages and import them

```
import xarray  
import netCDF4
```

Get the data file

```
curl -O http://www.ledo.columbia.edu/~rpa/argo_float_4901412.npz
```

Load data and list names of variables (data) and coordinates

```
argo_data = np.load('argo_float_4901412.npz')  
argo_data.keys()
```

Output:

```
['S', 'T', 'levels', 'lon', 'date', 'P', 'lat']
```

Then, extract the variable you want (salinity)

```
S = argo_data.f.S
```

Tools for documenting metadata



Check the dimensions of the salinity variable

```
print(S.shape)
```

Output:

```
(78, 75)
```

Make a DataArray with coordinates

```
da_salinity = xr.DataArray(S, dims=['level', 'date'],  
                           coords={'level': levels,  
                                   'date': date},)
```

Print (a summary) of your new array

```
da_salinity
```

Output:

Tools for documenting metadata



```
array([[ 35.638939,  35.514957,  35.572971, ... ,  35.820938,  35.777939,
       35.668911],
       [ 35.633938,  35.521957,  35.573971, ... ,  35.810932,  35.583897,
       35.667912],
       [ 35.681946,  35.525959,  35.572971, ... ,  35.795929,  35.662907,
       35.665913],
       ... ,
       [ 34.915859,  34.923904,  34.923904, ... ,  34.934811,  34.940811,
       34.946808],
       [ 34.915859,  34.923904,  34.921906, ... ,  34.932808,  34.93681 ,
       34.94381 ],
       [ 34.917858,  34.923904,  34.923904, ... ,         nan,  34.93681 ,
       nan]])
```

Coordinates:

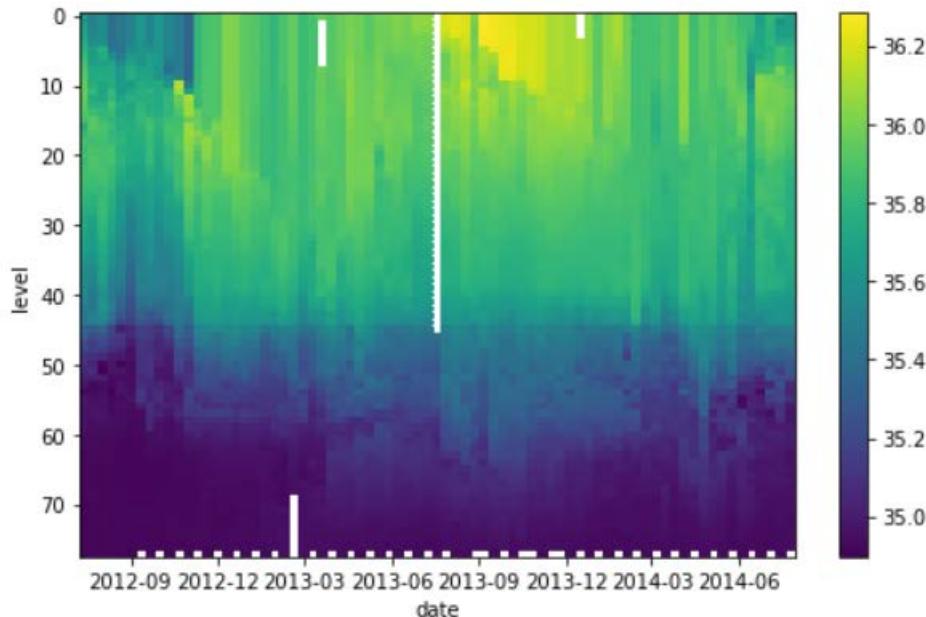
```
* level      (level) int64 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 ...
* date       (date) datetime64[ns] 2012-07-13T22:33:06.019200 ...
```

Tools for documenting metadata



Make a plot

```
da_salinity.plot(yincrease=False)
```



Tools for documenting metadata



Add some attributes (required by CF)

```
da_salinity.attrs['units'] = 'PSU'  
da_salinity.attrs['standard_name'] = 'sea_water_salinity'  
da_salinity
```

Output:

```
array([[ 35.638939,  35.514957,  35.572971, ...,  35.820938,  35.777939,  
       35.668911],  
       ...
```

Coordinates:

```
* level      (level) int64 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 ...  
* date       (date) datetime64[ns] 2012-07-13T22:33:06.019200 ...
```

Attributes:

```
units:          PSU  
standard_name: sea_water_salinity
```

Validating metadata

What can be validated?

- Conformance vs Climate and Forecast (CF) Convention, e.g.
 - Inclusion of mandatory attributes
 - Use of standard names for parameters, labels from controlled vocabularies
 - Inclusion of units for each variable
 - Correct data type of each element

What is not covered by validation

- Content of your title, summary, and other text attributes
- Valid contact information (e.g. e-mail)
- Individual validator tools may differ
 - e.g. <http://pumatest.nerc.ac.uk/cf-checker-status.html>

Validation tools for NetCDF/CF

NorDataNet validator: https://www.nordatanet.no/en/dataset_validation/form

- Supports CF 1.6 - CF 1.7 + ACDD
- Builds on [IOOS compliance checker](#)
- Generates a “score” and a detailed report

PUMA validator: <http://pumatest.nerc.ac.uk/cgi-bin/cf-checker.pl>

- Supports CF 1.0 - CF 1.7
- Developed from scratch
- Generates a detailed report (Errors, Warnings, Information)

Validation tools

ACDD

(Attribute Conventions
for Data Discovery)

Global Attributes

Highly Recommended

Attribute	Description
title	A short phrase or sentence describing the dataset. In many discovery systems, the title will be displayed in the results list from a search, and therefore should be human readable and reasonable to display in a list of such names. This attribute is also recommended by the NetCDF Users Guide and the CF conventions .
summary	A paragraph describing the dataset, analogous to an abstract for a paper.
keywords	A comma-separated list of key words and/or phrases. Keywords may be common words or phrases, terms from a controlled vocabulary (GCMD is often used), or URLs for terms from a controlled vocabulary (see also "keywords_vocabulary" attribute).
Conventions	A comma-separated list of the conventions that are followed by the dataset. For files that follow this version of ACDD, include the string 'ACDD-1.3'. (This attribute is described in the NetCDF Users Guide .)

Recommended

Attribute	Description
id	An identifier for the data set, provided by and unique within its naming authority. The combination of the "naming authority" and the "id" should be globally unique, but the id can be globally unique by itself also. IDs can be URLs, URNs, DOIs, meaningful text strings, a local key, or any other unique string of characters. The id should not include white space characters.
naming_authority	The organization that provides the initial id (see above) for the dataset. The naming authority should be uniquely specified by this attribute. We recommend using reverse-DNS naming for the naming authority; URLs are also acceptable. Example: 'edu.ucar.unidata'.
history	Provides an audit trail for modifications to the original data. This attribute is also in the NetCDF Users Guide : 'This is a character array with a line for each invocation of a program that has modified the dataset. Well-behaved generic netCDF applications should append a line containing: date, time of day, user name, program name and command arguments.' To include a more complete description you can append a reference to an ISO Lineage entity; see NOAA EDM ISO Lineage guidance .

Validation tools for NetCDF/CF

Example:
Sea ice chart
from met.no
delivered as part
of CMEMS

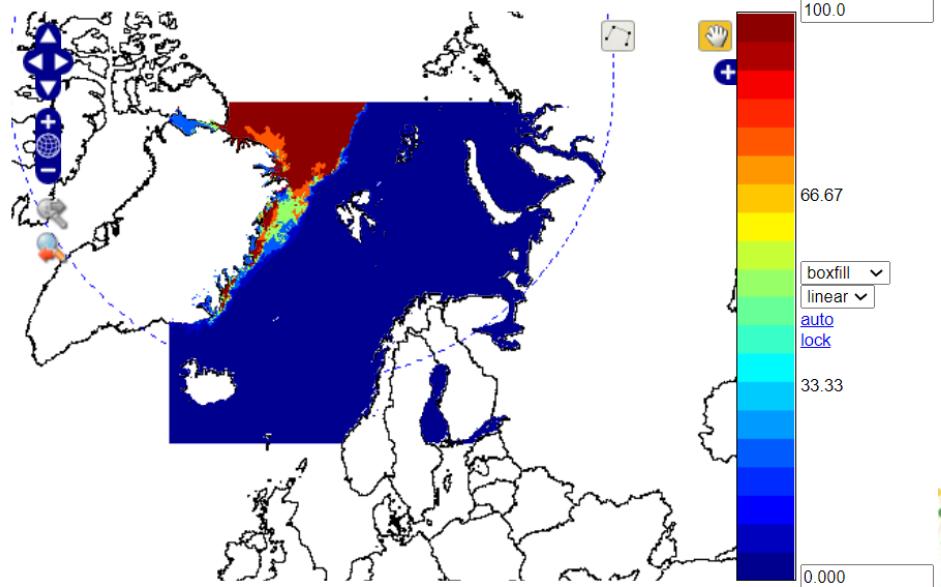
 MET Norway Thredds Service
 Arctic Svalbard & Barents Ice
Concentration, L4, 1km daily (METNO-ARC-
SEAICE_CONC-L4-NRT-OBS)
 - sea_ice_area_fraction
 - concentration range
 - latitude
 - longitude

[User guide](#)



Date/time: UTC [first frame](#) [last frame](#)

[Fit layer to window](#)



[test image](#) [Open in Google Earth](#)

Overlay opacity:

Validation tools for NetCDF/

Example: PUMA validator

Finds CF version

Checking all variables

Summary report - ALL PASS

```
Checking against CF version auto...
Check another file | NetCDF format | CF Convention.
```

File name: ice_conc_svalbard_202009041500.nc

Output of CF-Checker follows...

```
CHECKING NetCDF FILE: /tmp/13135.nc
-----
Using CF Checker Version 3.1.1
Checking against CF Version CF-1.4
Using Standard Name Table Version 74 (2020-08-04T14:43:55Z)
Using Area Type Table Version 10 (23 June 2020)
Using Standardized Region Name Table Version 4 (18 December 2018)

-----
```

Checking variable: time

```
-----
```

Checking variable: yc

```
-----
```

Checking variable: xc

```
-----
```

Checking variable: lat

```
-----
```

Checking variable: lon

```
-----
```

Checking variable: crs

```
-----
```

Checking variable: ice_concentration

```
-----
```

Checking variable: concentration_range

```
-----
```

ERRORS detected: 0
WARNINGS given: 0
INFORMATION messages: 0



Validation tools for NetCDF/CF

Example: NorDataNet validator
(also checks ACDD compliance)

Suggests improvements

Only supports CF 1.6 - 1.7

File uses CF 1.4

(Latest CF convention: 1.9)

IOOS Compliance Checker Report

Acdd:1.3
Corrective Actions

Highly Recommended | 6

Name	Reasoning
<i>Global Attributes</i>	<ul style="list-style-type: none">• keywords not present• summary not present• Conventions does not contain 'ACDD-1.3'
<i>variable "concentration_range" missing the following attributes:</i>	<ul style="list-style-type: none">• coverage_content• standard_name
<i>variable "ice_concentration" missing the following attributes:</i>	<ul style="list-style-type: none">• coverage_content
<i>variable "time" missing the following attributes:</i>	<ul style="list-style-type: none">• standard_name
<i>variable "xc" missing the following attributes:</i>	<ul style="list-style-type: none">• standard_name
<i>variable "yc" missing the following attributes:</i>	<ul style="list-style-type: none">• standard_name

Validation tools for NetCDF/CF

Example: NorDataNet validator

Follows CF 1.6

ALL TEST PASS

You are testing your dataset "S1A_EW_GRDM_1SDH_20150703T082545_20150703T082645_006644_008DE8_5388.nc" against CF-1.6 convention

Congratulations! Your dataset is compliant with the required test.

IOOS Compliance Checker

IOOS Compliance Checker Report

Cf:1.6
Corrective Actions

Errors 0

NorDataNet
Norwegian Scientific Data Network



Tools for documenting metadata

Rosetta

- Web-based application
 - Can read metadata from header(s), and add user defined metadata
 - Saves “setup” as templates for future use
 - Open source tool, written in Java
 - Customised version for NMDC, NorDataNet and SIOS
- <http://tomcat.nersc.no/rosetta/>
- [Rosetta User Manual](#)



Rosetta

This specific version of Rosetta has been tailored for NMDC, NorDataNet and SIOS.

Welcome to Rosetta, a data transformation tool. Rosetta is a web-based service that provides an easy, wizard-based interface for data collectors to transform their datalogger generated ASCII output into Climate and Forecast (CF) compliant netCDF files. These files will contain the metadata describing what data is contained in the file, the instruments used to collect the data, and other critical information that otherwise may be lost in one of many dreaded README files.

In addition, with the understanding that the observational community does appreciate the ease of use of ASCII files, methods for transforming the netCDF back into a user defined CSV or spreadsheet formats is planned to be incorporated into Rosetta.

We hope that Rosetta will be of value to the science community users who have needs for transforming the data they have collected or stored in non-standard formats.

Rosetta is currently under continued further development, and ready for beta testing.



What would you like to do?

[Convert a file to the netCDF format and create a new template](#)

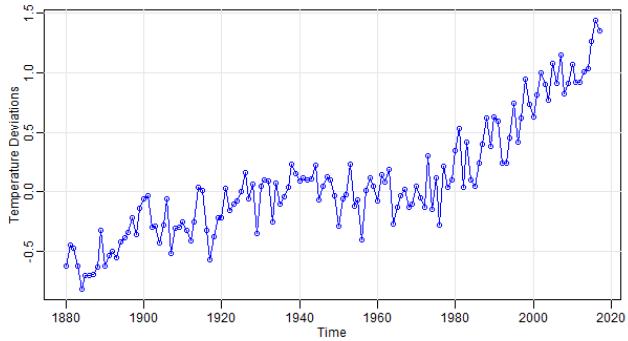
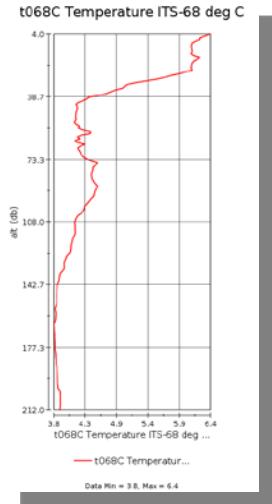
[Upload, modify, and use an existing template](#)

[Rosetta User Manual](#)

Tools for documenting metadata

Rosetta

- Profiles, time series, trajectory



- Template concept and benefits

Example - Rosetta

Steps

1. Start with template or build one
2. Select type of observation platform
3. Specify headerlines in the datafile
4. Specify delimiter and decimal separator
5. Specify variable attributes
6. Specify variable attribute details
7. Specify site specific information
8. Specify general information
9. Download resulting NetCDF file and new template



Example - Rosetta

Select to start without or with template:

What would you like to do?

Convert a file to the netCDF format and create a new template

Upload, modify, and use an existing template

Select observation Platform

Select Observation Platform



Single Station or
Tower (timeSeries)



Moored Buoy
(profile)



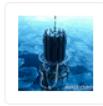
Radiosonde
(trajectory)



Wind Profiler
(profile)



Aircraft (trajectory)



Single CTD/XBT
cast (profile)

Upload file(s)

browse and upload (template and) text file

Example - Rosetta

Specify headerlines

All lines before data begins must be selected here

Specify Header Lines

Indicate which lines are header (i.e. not data) lines, or select 'No Header Lines' if there are none.

No header lines available in this file

	#	Line Data
<input checked="" type="checkbox"/>	315	# binavg_excl_bad_scans = yes
<input checked="" type="checkbox"/>	316	# binavg_skipover = 0
<input checked="" type="checkbox"/>	317	# binavg_surface_bin = no, min = 0.000, max = 0.000, value = 0.000
<input checked="" type="checkbox"/>	318	# Derive_date = Aug 12 2016 11:06:22, 7.21f [derive_vars = 5]
<input checked="" type="checkbox"/>	319	# Derive_in = C:\Seabird\CTDData\sta0591.cnv C:\Seabird\SBEDataProcessing-Win32\Seasoft.com
<input checked="" type="checkbox"/>	320	# file_type = ascii
<input checked="" type="checkbox"/>	321	*END*
<input type="checkbox"/>	322	1566 4.000 6.3774 3.307919 6.3800 3.307177 96.48 2.5116 0.2472 65.6037 2.3088e+02 7.34149 33.1770 33.1662 26.0649 1474.00 1474.00 0.0000...
<input type="checkbox"/>	323	1587 5.000 6.2846 3.305076 6.2843 3.306049 96.41 2.5115 0.2767 65.4606 1.8441e+02 7.34798 33.2354 33.2465 26.1228 1473.72 1473.94 0.0000...
<input type="checkbox"/>	324	1606 6.000 6.2037 3.296740 6.1752 3.294673 96.41 2.5105 0.3296 65.3523 1.4661e+02 7.34028 33.2209 33.2257 26.1217 1473.40 1473.85 0.0000...
<input type="checkbox"/>	325	1631 7.000 6.1971 3.295611 6.2059 3.296281 96.41 2.5074 0.3811 64.9476 1.1881e+02 7.34579 33.2143 33.2131 26.1172 1473.38 1473.78 0.0000...
<input type="checkbox"/>	326	1655 8.000 6.1298 3.291956 6.1733 3.294838 96.41 2.5087 0.4114 64.8805 9.7768e+01 7.38343 33.2388 33.2284 26.1450 1473.16 1473.71 0.0000...



Example - Rosetta

pecify delimiter and decimal separator:

Specify Delimiters

Please specify delimiter(s) used.

Tab Comma Whitespace Colon
 Semicolon Single Quote Double Quote Other

Please specify decimal separator used.

Point Comma



Example - Rosetta

Specify variable attributes

Specify Variable Attributes

Click on each column and specify the information asked for.

Specify 'Do not use this column of data' for all columns that are not to be saved in the netCDF file.

All columns must have a green tickmark before you can continue.

#	<input checked="" type="checkbox"/> Do Not Use	<input checked="" type="checkbox"/> sea_water_p	<input checked="" type="checkbox"/> sea_water_t	<input checked="" type="checkbox"/> sea_water_e	<input checked="" type="checkbox"/> sea_water_t	<input checked="" type="checkbox"/> sea_water_e	<input checked="" type="checkbox"/> height_abov	<input checked="" type="checkbox"/> volume_frac	<input checked="" type="checkbox"/> fluorescenc	<input checked="" type="checkbox"/> beam_transm	<input checked="" type="checkbox"/> downwelling	<input checked="" type="checkbox"/> volume_frac
+ 0 * Sea-Bird SBE 9 Data File:												
322	1566	4.000	6.3774	3.307919	6.3800	3.307177	96.48	2.5116	0.2472	65.6037	2.3088e+02	7.34149
323	1587	5.000	6.2846	3.305076	6.2843	3.306049	96.41	2.5115	0.2767	65.4606	1.8441e+02	7.34798
324	1606	6.000	6.2037	3.296740	6.1752	3.294673	96.41	2.5105	0.3296	65.3523	1.4661e+02	7.34028
325	1631	7.000	6.1971	3.295611	6.2059	3.296281	96.41	2.5074	0.3811	64.9476	1.1881e+02	7.34579
326	1655	8.000	6.1298	3.291956	6.1733	3.294838	96.41	2.5087	0.4114	64.8805	9.7768e+01	7.38343
327	1680	9.000	6.0743	3.289393	6.0743	3.289479	96.41	2.5076	0.4601	64.7450	8.0697e+01	7.35942
328	1704	10.000	6.0675	3.289681	6.0684	3.289333	96.41	2.5057	0.5427	64.3239	6.6608e+01	7.34023
329	1728	11.000	6.0673	3.290823	6.0681	3.290466	96.41	2.5023	0.5681	63.9660	5.5489e+01	7.33149
330	1752	12.000	6.0709	3.293776	6.0715	3.293283	96.41	2.5007	0.6037	64.0387	4.6529e+01	7.34642
331	1777	13.000	6.0775	3.295214	6.0781	3.295042	96.41	2.5003	0.6467	64.0487	3.9042e+01	7.35050
		11.000	6.0770	3.295000	6.0770	3.295700	96.40	2.5000	0.6701	64.0400	3.9042e+01	7.35015

Example - Rosetta

Specify variable
attribute details

Enter Variable Attributes

What would you like to do with this column of data?

Assign a variable name Do not use this column of data
sea_water_pressure

use metadata from another column?

Is this variable a coordinate variable? (examples: latitude, longitude, time)

Yes No

What type of coordinate variable?

altitude

Specify variable data type:

Integer Float (decimal) Text

Required Metadata:

Variable Description **prDM Pressure Digiquartz db**

Units **db**
 show unit builder

Recommended Metadata:

CF Name **sea_water_pressure**

Additional Metadata:

[+] [-] Calendar Type

done

Enter Variable Attributes

What would you like to do with this column of data?

Assign a variable name Do not use this column of data
sea_water_temperature

use metadata from another column?

Is this variable a coordinate variable? (examples: latitude, longitude, time)

Yes No

Specify variable data type:

Integer Float (decimal) Text

Required Metadata:

Instrument Description **Temperature SensorID 55 Serial#**

Missing Value **99**

Variable Description **t068C Temperature ITS-68 deg C**

Units **degree_Celsius**
 show unit builder

Recommended Metadata:

Instrument Height (negative for depths)

Instrument Height Unit

Maximum Value (Calibrated)

Minimum Value (Calibrated)

CF Name **sea_water_temperature**

Additional Metadata:

[+] [-] Calibration Range

done

Example - Rosetta

Specify site specific information

- what information that appears here is dependent on type of observation platform

Specify Site Specific Information

* denotes required field

*Station or Platform Name?

is a regex

F/F Håkon Mosby

*Station or Platform Date and Time?

is a regex

2016-08-12

*Station Latitude?

is a regex

78.05

degrees_north ▾

*Station Longitude?

is a regex

13.533

degrees_east ▾

Example - Rosetta

Specify general information

Specify General Information

* denotes required field

* Title ?

is a regex

CTD station collected near Isfjorde

* Naming Authority ?

is a regex

UNIS

* ISO Topic Category ?

oceans

* Keywords ?

is a regex

EARTH SCIENCE, OCEANS, OCEAN

* License ?

is a regex

CC-BY

* Publisher Name ?

is a regex

Ragnheid Skogseth

* Summary ?

is a regex

CTD data collected near Isfjorden,

* ID ?

is a regex

UNIS-CTD-20160812-qc-binned

* Keywords Vocabulary ?

GCMD Science Keywords

* Data Assembly Center ?

is a regex

UNIS

* Processing Level ?

is a regex

Unknown

Publisher Email ?

is a regex

Ragnheid.Skogseth@unis.no

Example - Rosetta

Download resulting files.

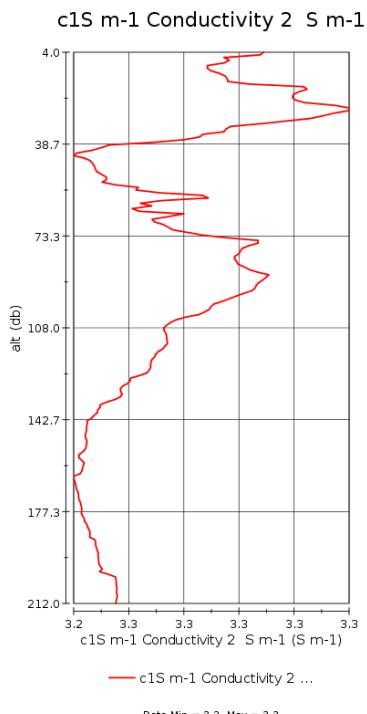
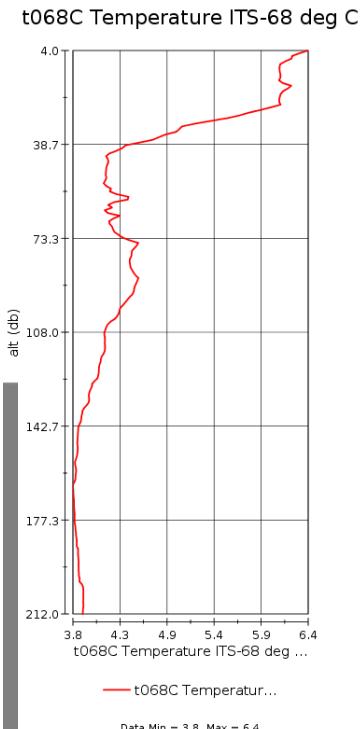
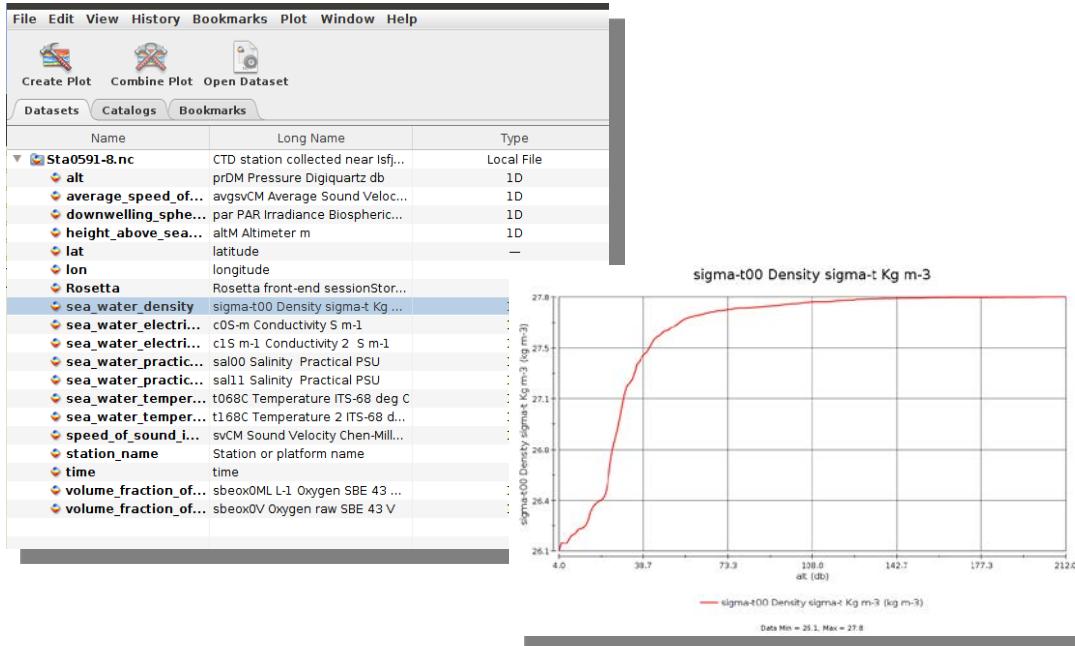
Download Converted File

 [netCDF Data File](#)

 [Sta0591-Rosetta_2020-06-21_232216.template](#)

Example - Rosetta

Inspect your new NetCDF file, e.g. download Panoply:
<https://www.giss.nasa.gov/tools/panoply/download/>

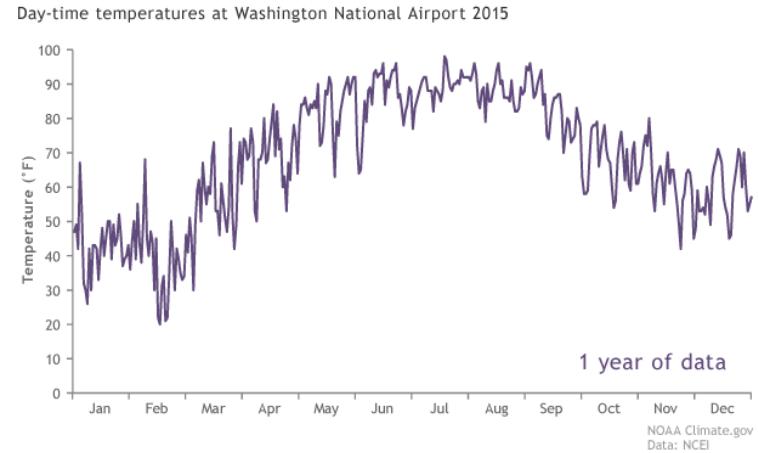


Benefits of standardisation

- Metadata
 - Document all important aspects of your data so you can reuse it! ...and know what you are working with...
 - Become visible as a scientist producing valueable data
 - Support reproducible science - by enabling other scientist to reuse your data
- Formats
 - Large collections of libraries and tools you can use
 - Eases collaboration with other scientists - share data *and* code
 - Metadata can be included in the same file as data (e.g. NetCDF)

Assignment 1: Documenting and formatting YOUR own dataset

- Data:
 - Bring your data file, e.g. in CSV format
 - Metadata
 - Compile associated information (e.g. notes from data collection, reports from field experiments, scientific papers, emails from colleagues)
 - Include enough information for both yourself as well as other scientists to reuse it!
 - Formats
 - Use Rosetta, NCO/CDO or Python (own script) to prepare your dataset in NetCDF/CF format



Assignment 1: Documenting and formatting YOUR dataset

- Objectives:

- Hands-on preparing and documenting data for use in analysis

- Tasks:

- Find description needed to use the data, you & other scientists
 - Prepare this for documenting data
 - Data conversion to a standard format

- Submission of assignment

- Deadline: 19 Oct 2021, 21:00 Bergen time
 - How to submit: [Google Drive Course Folder](#)
 - Work individually or in pairs

- Deliverables:

- NetCDF file
 - Short presentation

- Presentation (max 3 slides):

- Name, email
 - Title or your dataset
 - Description of your dataset
 - Documentation sources for data
 - Data conversion process
 - Assessment of FAIRness
 - Original format
 - Generated NetCDF
 - What you learnt

Assignment 1: Documenting and formatting the BMW dataset

- Data:
 - What it contains (met data)
 - Where to find it (a: download the data yourself in CSV format, b: get data from course organisers)
- Metadata
 - Find description of parameters, units, etc. that you will need to use the data in assignment 2
 - Include enough information for both yourself as well as other scientists to reuse it!
- Formats
 - Use Rosetta, NCO/CDO or Python (own script) to prepare the BMW dataset in NetCDF/CF format

The screenshot shows the 'Last ned data' (Download data) section of the Været i Bergen website. It includes fields for 'Tidslinje' (Timeline) with start and end dates, a 'Stasjon' (Station) dropdown set to 'Ulriken', and a 'Parametere' (Parameters) list with checkboxes for various meteorological variables. Below this is a 'Filformat' (File format) section with options for 'Excel (.xlsx)' and 'CSV'. At the bottom, there are links for 'Siste degn' (Last week), 'Historiske data' (Historical data), 'Varinfo' (Variable info), and 'Universitet i BERGEN'.

<https://veret.gfi.uib.no/?action=download>

Assignment 1: Documenting and formatting the BMW dataset

- Data:
 - Parameters to include: ALL
 - Period: 10.10.2003 - YESTERDAY
- Metadata
 - Find description of parameters, units, etc. that you will need to use the data in assignment 2 – own reuse!
 - Include enough information for OTHER SCIENTISTS! to reuse it
- Formats
 - NetCDF/CF – Use the NorDataNet validator to check your file

Global metadata:

- Title
- Summary
- Creator(name, email, ...)
- License
- ...

Parameter metadata:

- Parameter name
- Parameter standard name
- Unit
- Valid min, valid max
- Fill value
- ...

Assignment 1: Documenting and formatting the BMW dataset

- Objectives:
 - Hands-on preparing and documenting data for use in analysis
- Tasks:
 - Find description needed to use the data, you & other scientists
 - Prepare this for documenting data
 - Data conversion to a standard format
- Submission of assignment
 - Deadline: 19 Oct 2021, 21:00 Bergen time
 - How to submit: [Google Course Folder – Assignment 1](#)

Deliverables:

- NetCDF file
- Short presentation

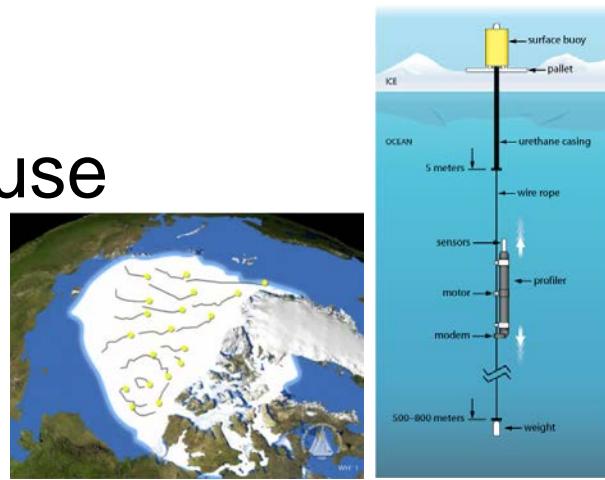
Presentation (max 3 slides):

- Name, email
- Title or your dataset
- Description of your dataset
- Documentation sources for data
- Data conversion process
- Assessment of FAIRness
 - Original format
 - Generated NetCDF
- What you learnt

Assignment 1:

Other datasets you can use

- Datasets
 - Woods Hole Oceanographic Institution (Ice Thethered Platform data)
 - EUMETSAT Ocean and Sea Ice SAF ice extent indicator
 - A test setup with multiple surface irradiance measurements at MET (a combined dataset with K&Z CNR4 and Apogee net radiation instruments)
- Tasks:
 - Find description needed to use the data, you & other scientists
 - Prepare this for documenting data
 - Data conversion to a standard format
- How to get the data
 - Open the [Course Data Folder](#) in Github



Resources

Rosetta: <http://tomcat.nersc.no/rosetta/> ; <https://github.com/Unidata/rosetta>

NCO/CDO: <http://nco.sourceforge.net/>

netCDF4: <https://unidata.github.io/netcdf4-python/netCDF4/index.html>

xarray: <http://xarray.pydata.org/en/stable/>

DataOne: <https://old.dataone.org/education-modules>

ESIP: <https://commons.esipfed.org/node/1422> ; [Data management training](#)