

## The Emergence of Reference Datasets

Increasingly, large data/HPC facilities now manage major reference collections (e.g., satellite Earth observation, geophysics) in response to a growing research demand at larger geographical scales and/or over longer time periods. Data from multiple surveys and/or progressive time series acquisitions are aggregated into co-located reference collections: some are further standardised and structured into High Performance Datasets (HPD) for use in Data-Intensive Science. Reproducibility requires a capability to be able to access the extract used in research, but at-scale, multiple data extracts cannot be stored indefinitely and there are new requirements from publishers to include persistent identifiers and landing pages for citation.

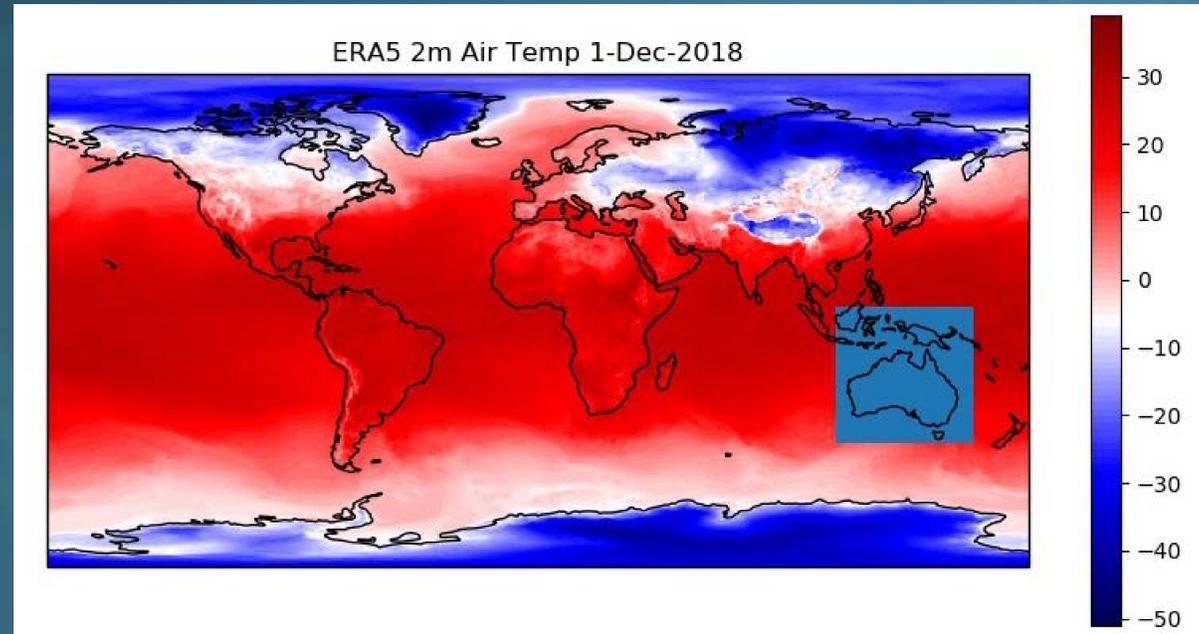
## The New Publisher Requirements

The revised Commitment Statement by the Coalition for Publishing Data in the Earth and Space Sciences (COPDESS) requires research outputs related to publications to be FAIR and have unique persistent identifiers assigned to each generated/produced dataset. This is easy for small datasets developed by a particular research project over a defined period: outputs are made available for download as files that relate to a publication. Any repository distributing such datasets acts more as a traditional library: with no changes/enhancements made to the dataset.

## How Do YOU Cite the Data Derivatives?

Supporting consistent approaches to citation for derivative datasets from large, often dynamic, reference datasets requires community agreement and consideration of costs and technical solutions. Developing such publication standards will not be easy: they have to be consistent across multiple Earth and environmental centres internationally, and include not just research repositories but also government repositories. Increasing dependency by researchers on commercial cloud data storage services such as Google, Amazon, etc., means they might also need to be engaged in this process.

Fig 1. The ERA-5 2m air temperature dataset. ERA-5 provides hourly information at a horizontal resolution of around 31 km globally and at 137 levels in the vertical. The entire dataset from 1950s will be available in late 2019 and is expected to grow to 20 Petabytes. (Graphic from Kate Snow, NCI)



# How do YOU preserve your dataset extracts to ensure reproducibility and meet new publisher requirements?

Community discussion is being organised through the Research Data Alliance Data Versioning Working Group here:



WG Homepage  
<https://bit.ly/2U0Vxyd>

Please comment on the White Paper being developed here:

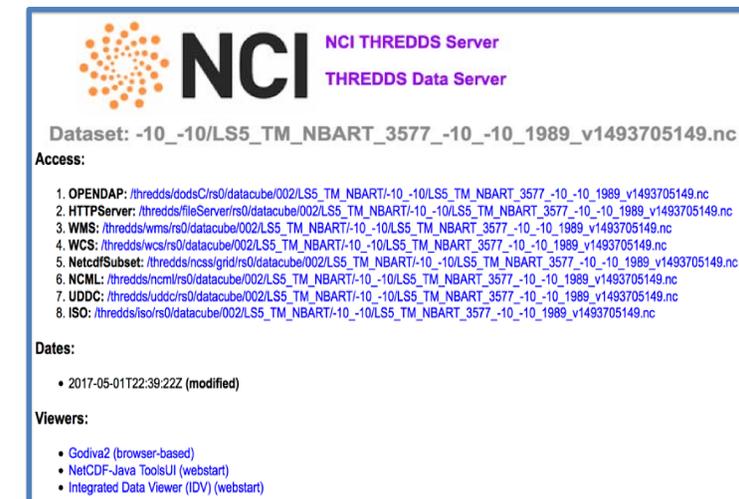


WG White Paper  
<https://bit.ly/2w7220>

## Perspectives on a Major NCI Collection

Landsat	5	NBAR Pixel Quality	25m v2
Landsat	5	Surface Reflectance NBAR	25m v2
Landsat	5	Surface Reflectance NBAR-T	25m v2
Landsat	7	NBAR Pixel Quality	25m v2
Landsat	7	Surface Reflectance NBAR	25m v2
Landsat	7	Surface Reflectance NBAR-T	25m v2
Landsat	8	NBAR Pixel Quality	25m v2
Landsat	8	Surface Reflectance NBAR	25m v2
Landsat	8	Surface Reflectance NBAR-T	25m v2

Table 1. Multiple derivatives of the Landsat Datasets at NCI



Dataset: -10\_-10/LS5\_TM\_NBART\_3577\_-10\_-10\_1989\_v1493705149.nc

Access:

- OPENDAP: /thredds/dodsC/rs0/datacube/002/LS5\_TM\_NBART/10\_-10/LS5\_TM\_NBART\_3577\_-10\_-10\_1989\_v1493705149.nc
- HTTPServer: /thredds/fileServer/rs0/datacube/002/LS5\_TM\_NBART/10\_-10/LS5\_TM\_NBART\_3577\_-10\_-10\_1989\_v1493705149.nc
- WMS: /thredds/wms/rs0/datacube/002/LS5\_TM\_NBART/10\_-10/LS5\_TM\_NBART\_3577\_-10\_-10\_1989\_v1493705149.nc
- WCS: /thredds/wcs/rs0/datacube/002/LS5\_TM\_NBART/10\_-10/LS5\_TM\_NBART\_3577\_-10\_-10\_1989\_v1493705149.nc
- NetcdfSubset: /thredds/netcdf/rs0/datacube/002/LS5\_TM\_NBART/10\_-10/LS5\_TM\_NBART\_3577\_-10\_-10\_1989\_v1493705149.nc
- NCML: /thredds/netcdf/rs0/datacube/002/LS5\_TM\_NBART/10\_-10/LS5\_TM\_NBART\_3577\_-10\_-10\_1989\_v1493705149.nc
- UDDC: /thredds/uddc/rs0/datacube/002/LS5\_TM\_NBART/10\_-10/LS5\_TM\_NBART\_3577\_-10\_-10\_1989\_v1493705149.nc
- ISO: /thredds/iso/rs0/datacube/002/LS5\_TM\_NBART/10\_-10/LS5\_TM\_NBART\_3577\_-10\_-10\_1989\_v1493705149.nc

Dates:

- 2017-05-01T22:39:22Z (modified)

Viewers:

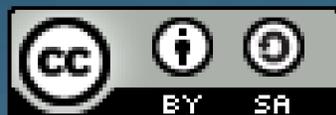
- Godiva2 (browser-based)
- NetCDF-Java ToolsUI (webstart)
- Integrated Data Viewer (IDV) (webstart)

Fig 2. Multiple ways to access and view just one derivative

## Potential approaches from NCI

To reproduce a derivative dataset from a national/global reference dataset requires a dual approach:

1. Researchers need to publish references to all dependent datasets as well as the recipe on how the dataset is derived; and
2. Repositories need to record and make publicly accessible (preferably in machine readable ways) any changes to the dataset as well as referencing the data service(s) that provide data access.



Australian National University

