

NATIONAL COMPUTATIONAL INFRASTRUCTURE

Supporting Data Reproducibility at NCI Using the Provenance Capture System

Jingbo Wang, Nick Car, Wei Si, Edward King, Ben Evans, Lesley Wyborn











ustralian Government Geoscience Australia Australian Government Australian Research Council









More than 70% of researchers have tried and failed to reproduce another scientist's experiments, and more than half have failed to reproduce their own experiments. Those are some of the telling figures that emerged from Nature's survey of 1,576 researchers who took a brief online questionnaire on reproducibility in research.

Source: http://www.nature.com/news/1-500-scientists-lift-the-lid-onreproducibility-1.19970



NCI is high performance computer centre, hosting 10PB research data for sharing among research community. - HPC

NCI is building the National Environmental Research Data Interoperability Platform (NERDIP) to enable data intensive science. - HPD



User generate/ transfer data Data Manager fill DMP and create catalogue

Super computer users

Paper and Data are published

Data visualization Data share and re-use













NCI provides user with Data as a Service

Persistent global parallel filesystem -7.1 PB -4.2 PB -4.2 PB -4.2 PB -4.2 PB -7.6 PB		IPC		NCI Vislab	Web-time analytics software Image: Contract of the second secon
NCI provides fast data storage	Data Management Portal		Data Curation, Publish, Citation		Virtual Desktop Interface, Virtual Laboratory, and other

nci.org.au

services

NCI Proposed solution

Provenance Capture System: to support traceable, reproducible, and machine actionable workflow





The basic classes and relationships of PROV-DM, drawn in accordance with the PROV Ontology.



The information is stored as RDF documents in an RDF graph database, which is then published via a HTTP-based API as a combination of web pages and data.





NCI Proof of concept workflow



NCI Supporting infrastructure: PID

Persistent Identifier Services



NCI Supporting infrastructure: catalogue

NCI GeoNetwork architecture https://geonetwork.nci.org.au





CC

NCI Supporting infrastructure: doc repository

NCI document repository https://cms.nci.org.au

NCI's Document	t Repository
HOME NCI ABOUT PUBLICATION	
ERESEARCH Australasia 2016	<section-header> Image: Second Secon</section-header>
Navigation	NCI's Document Repository
Abstract Articles	



NCI Ocean Colour Data Processing Flow



1.Assembling the data (Level-0 granules); 2.Calibration and geolocation (Level-0 to Level-1b); 3. Generation of ancillary fields (Level-1b to Level-2); 4.Data selection and masking; 5. Atmospheric correction; 6.In-water inversion (deriving concentrations of optically active substances); 7.Regridding (conversion of satellite grids to map projection); 8. Mosaicing (joining multiple satellite images to provide more complete coverage); 9. Data distribution and access.

NCI Provenance report of Ocean Colour Case

An Activity

URI: file:///opt/proms/raijin/g/data/u83/code/ereefs/legacy/proms/promReport_OC/

html | rdf/turtle

Title:	Ocean Colour Process
Started at:	2016-07-01T11:20:56
Ended at:	2016-07-01T11:20:56
Was associated with:	http://orcid.org/0000-0002-6898-2130
Report:	Ocean Colour External Report - 2016-07-01 11:20:56.261893

Neighbours view





The staff resources needed to establish the Ocean Colour reporting were:

•Part-time work over several months to implement the central information services;

•Several days work from three people (an OC product owner and modellers) to model the Ocean Colour process as per PROV-O;

•A week to implement logging within the OC workflow that could be used to establish PROV-O reports from each workflow run;

•storage location URIs generated by dataset registry and PID services

•A day or two to schedule reporting to work continuously.



Thank you

Questions?

