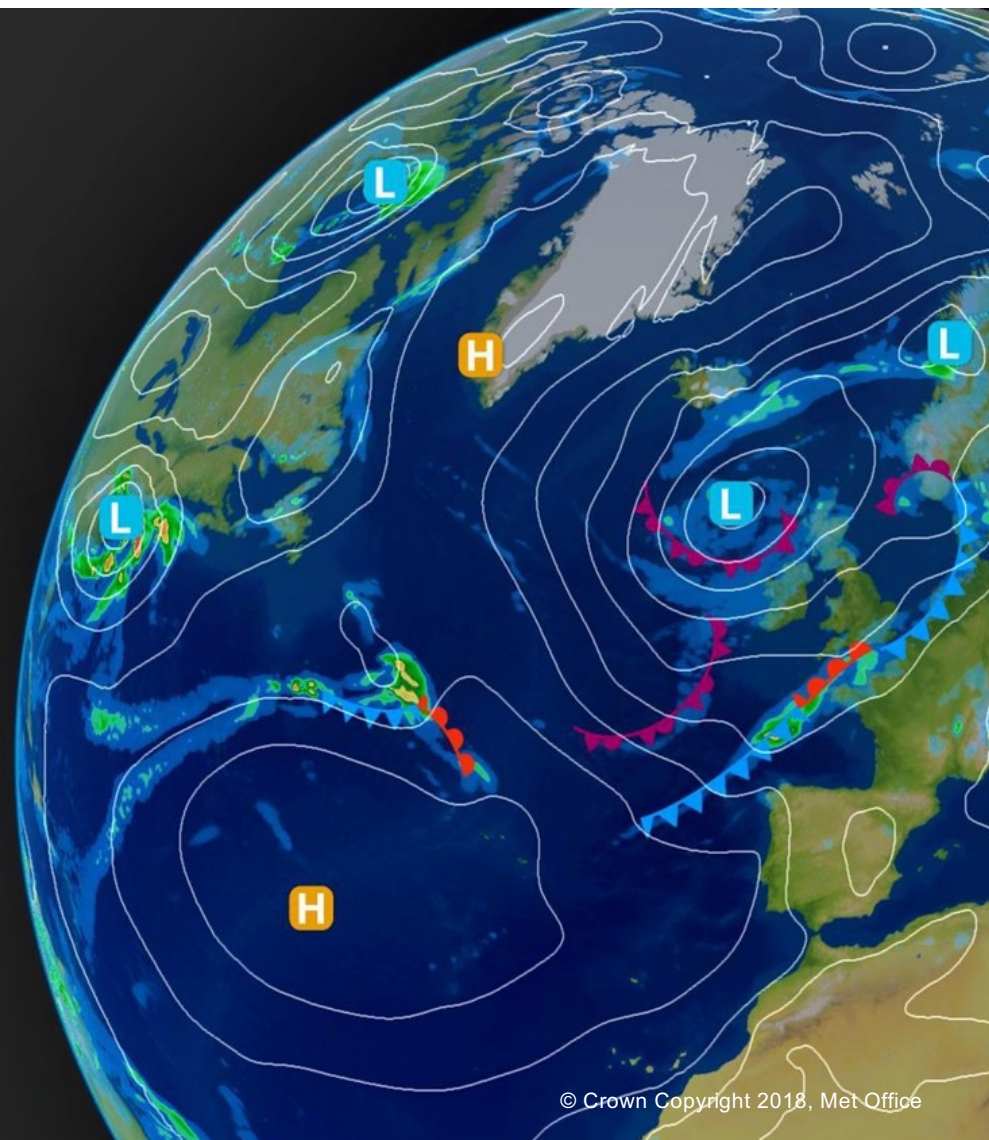


Analysis facilities and Cloud

Matthew Mizielinski
Paul Durack
and the WIP

WGCM-24
9th December 2021 - Virtual



Contributors to CMIP6



<https://pcmdi.llnl.gov/CMIP6/>

NCAR: GLADE & Casper

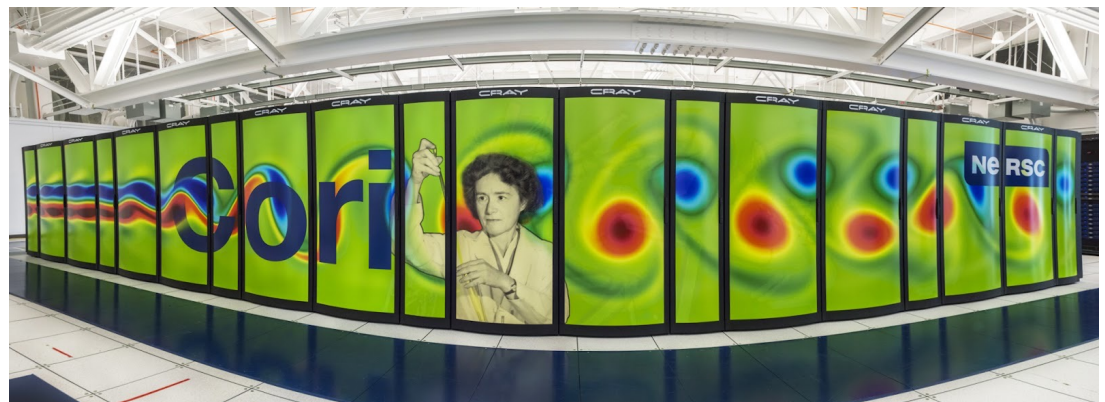
- 100 node Casper cluster including high-memory, GPU and high throughput nodes
- Access to local archive of CMIP, and other project data on GLADE storage facility
- Access available to NSF funded researchers and US university staff and students



<https://www2.cisl.ucar.edu/resources/cmip-analysis-platform>

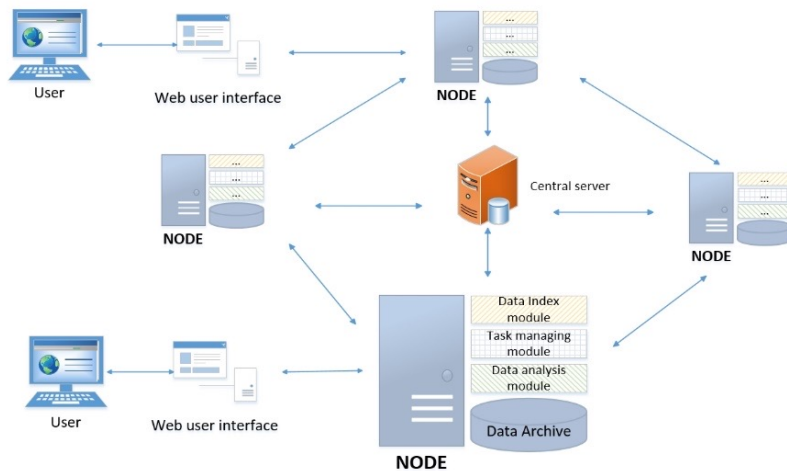
DOE: NERSC Cori

- Cray XC40 HPC with large number of conventional and Xeon Phi nodes
- Large HPSS archives and various types of local storage, including flash
- Open to DoE Office of Science projects, significant use by climate science



Typical climate model analysis facilities in Asia (1)

China: A Collaborative Analysis Framework for distributed gridded Environmental data (CAFÉ)

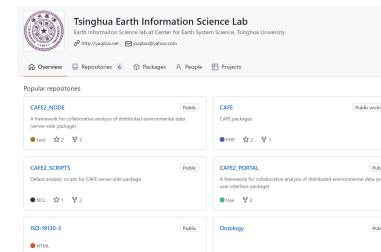


Xu, Hao, Sha Li, Yuqi Bai, Wenhao Dong, Wenyu Huang, Shiming Xu, Yanluan Lin et al.
"A collaborative analysis framework for distributed gridded environmental data."
Environmental Modelling & Software 111 (2019): 324-339.

• A new "ZERO Download" mode

- Multiple data nodes establish a federation.
- Users could submit data analysis tasks and then get the analysis result.
- Data analysis are performed where the data reside.
- Users do not have to download the model data.

• A GitHub project



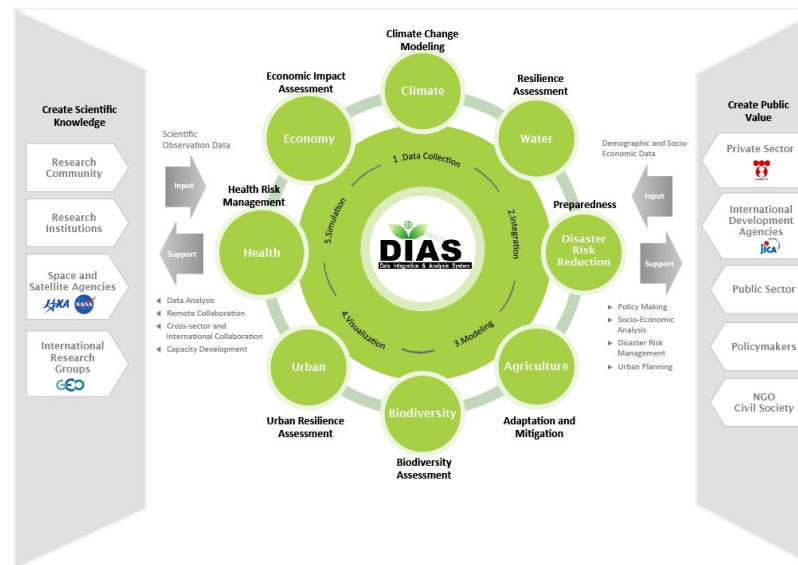
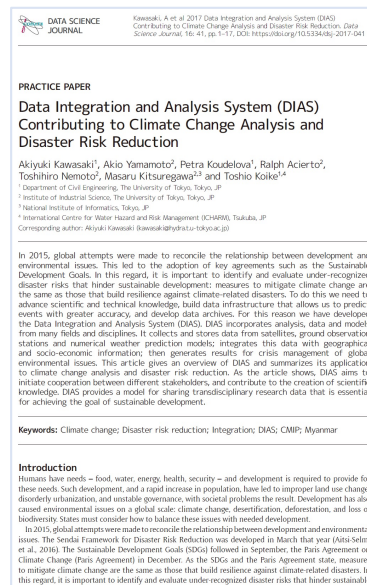
• A Patented Technology



Courtesy Yuqi Bai

Typical climate model analysis facilities in Asia (2)

Japan: Data Integration and Analysis System (DIAS)



Datasets

Datasets available in DIAS (Last Updated: 2021/3/4)

[*] Can be downloaded directly from DIAS. Other datasets can be downloaded from provider institutions database.

[**] Now in preparation, limited access only, or under consultation.

- Earth Observation Satellites
- Greenhouse Gases Observations
- Terrestrial Ecosystems / Carbon Flux Observations
- Weather Observations
- Watershed Observations
- Ocean Observations
- Reanalysis
- Prediction
- Downscaled Data
- Natural Disasters
- Land Use
- Health Hazard

Kawasaki, Akiyuki, Akio Yamamoto, Petra Koudelova, Ralph Acierto, Toshihiro Nemoto, Masaru Kitsuregawa, and Toshio Koike.
 "Data integration and analysis system (DIAS) contributing to climate change analysis and disaster risk reduction."
Data Science Journal 16 (2017).

Courtesy Yuqi Bai

Typical climate model analysis facilities in Asia (3)

Australia: NCI High Performance Computing and High Performance Data Platform

The NCI High Performance Computing and High Performance Data Platform to Support the Analysis of Petascale Environmental Data Collections

Ben Evans¹, Lesley Wyborn¹, Tim Pugh², Chris Allen¹, Joseph Antony¹, Kashif Gohar¹, David Porter², Jon Smillie², Claire Trenham², Jingbo Wang³, Alex Ip¹, and Gavin Bell⁴

¹ National Computational Infrastructure (NCI), Australian National University, Canberra, Australia

² Bureau of Meteorology, Melbourne, Australia

³ Geoscience Australia, Canberra, Australia

⁴ The 9th Column Project, Berlin, Germany

(Ben.Evans, Lesley.Wyborn, Chris.Allen, Joseph.Antony, Kashif.Gohar, David.Porter, Jon.Smillie, Claire.Trenham, Jingbo.Wang)@anu.edu.au, T.Pugh@bom.gov.au, Alex.Ip@ga.gov.au, gavin@9thcolumn.org

Abstract. The National Computational Infrastructure (NCI) at the Australian National University (ANU) has co-located a priority set of over 10 Petabytes (PBytes) of national data collections within a HPC research facility. The facility provides an integrated high-performance computational and storage platform, or a High Performance Data (HPD) platform, to serve and analyse the massive amounts of data across the spectrum of environmental collections – in particular from the climate, environmental and geoscientific domains. The data is managed in concert with the government agencies, major academic research communities and collaborating overseas organisations. By co-locating the vast data collections with high performance computing environments and harmonising these large valuable data assets, new opportunities have arisen for Data-Intensive interdisciplinary science at scales and resolutions not hitherto possible.

Keywords: high performance computing, high performance data, cloud computing, data-intensive science, scalable data services, data cube, virtual laboratories.

1 Introduction

The National Computational Infrastructure (NCI) at the Australian National University (ANU) has organised a priority set of large volume national environmental data assets on a High Performance Data (HPD) Node within a High Performance

R. Doster et al. (Eds.), IESS 2015, IFIP AICT 448, pp. 569–577, 2015.
© IFIP International Federation for Information Processing 2015

Australia's preeminent high-performance data, storage and computing facility.



The National Computational Infrastructure (NCI) is Australia's leading high-performance data, storage and computing organisation, providing expert services to benefit all domains of science, government and industry.

NCI brings the Australian Government and the Australian research sector together through a broad collaboration involving the largest national science agencies, universities, industry and the Australian Research Council.

NCI empowers government agencies, universities, and industry across multiple domains of research. Our integrated hardware, services and expertise drive high-impact research and groundbreaking outcomes for Australia.

Evans, Ben, Lesley Wyborn, Tim Pugh, Chris Allen, Joseph Antony, Kashif Gohar, David Porter et al.

"The NCI high performance computing and high performance data platform to support the analysis of petascale environmental data collections."

In *International Symposium on Environmental Software Systems*, pp. 569-577. Springer, Cham, 2015.

Courtesy Yuqi Bai

Tiered of downloading Terabytes out of Petabytes of climate model data ?

Large European climate data centers offer the possibility to directly exploit locally available large climate data pools (e.g. CMIP6 data)

Two types of service:

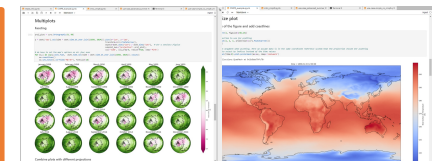
- **“Jump start service”:**
 - + minimal application procedure
 - limited compute resources
- **Analysis platform service:**
 - short project proposal required
 - + guaranteed resource allocation

The offering from DKRZ, IPSL-CNRS, UKRI-CEDA, CMCC:

- **Access to large European climate model data pools** (multi-PByte data collections including CMIP6, CORDEX, ..)
- **Access to associated HPC compute resources**
- **Access to interactive analysis environments** (including jupyter-hub installations at DKRZ, CMCC and STFC)
 - support for e.g. pangeo sw stack (xarray, dask), cdo, ESMValTool and user tailored environments..

Interested? Further information:

- **Climate Analytics service (ECAS):**
<https://portal.enes.org/data/data-metadata-service/climate-analytics-service>
- **Analysis platforms application:**
<https://portal.enes.org/data/data-metadata-service/analysis-platforms>
- **Demos, use-cases, example jupyter notebooks:**
<https://github.com/IS-ENES-Data/Climate-data-analysis-service>



**Next
deadline:
20.12.2021**

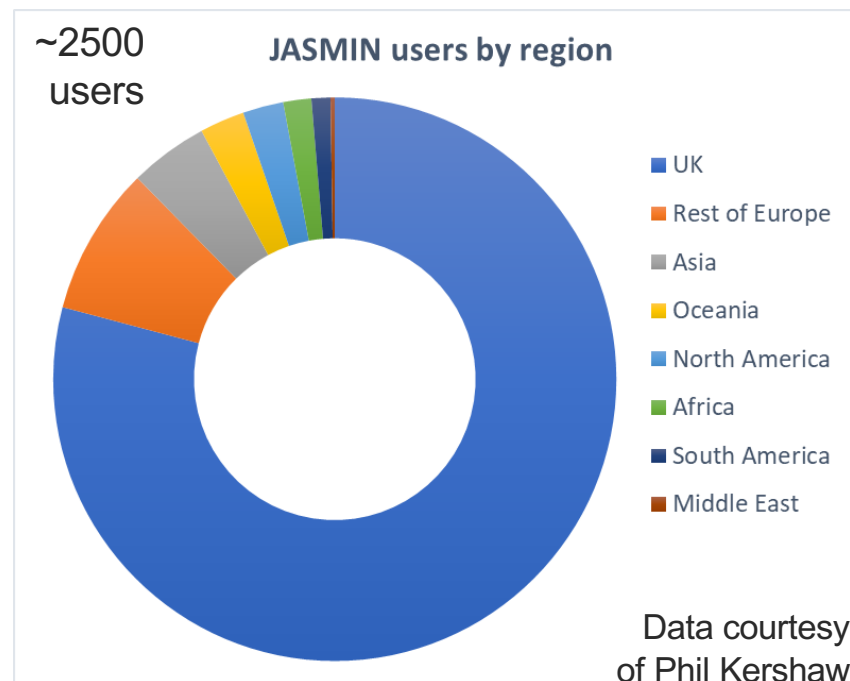


The IS-ENES3 project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 824084

**Stephan
Kindermann**

Access to analysis facilities

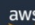
- Most facilities are restricted access
 - Local users
 - Funded collaborations
- Gap in provision for Africa and South America



Pangeo and commercial cloud

- AWS providing hosting for a subset of CMIP6 data in Zarr format along with compute through ASDI
 - Some free access for users
 - Many other datasets available
- Pangeo providing and managing data
 - Collaborating with GFDL and CEDA
- Change in data storage, and compression, has implications for analysis and archiving options

PANGEO

Registry of Open Data on AWS 

Coupled Model Intercomparison Project 6

[agriculture](#) [atmosphere](#) [climate](#) [earth observation](#) [environmental](#) [model](#) [oceans](#) [simulations](#) [weather](#)

Description

The sixth phase of global coupled ocean-atmosphere general circulation model ensemble.

This [application](#) is one of several possibilities to find CMIP6 data citations. Alternative tools to find CMIP6 data references are described in [this blog post](#). General information on the Citation Service is available at: [cmip6cite.wdc-climate.de](#).


Update Frequency

Core CMIP6 datasets are added as soon as they are available.

1. Introduction

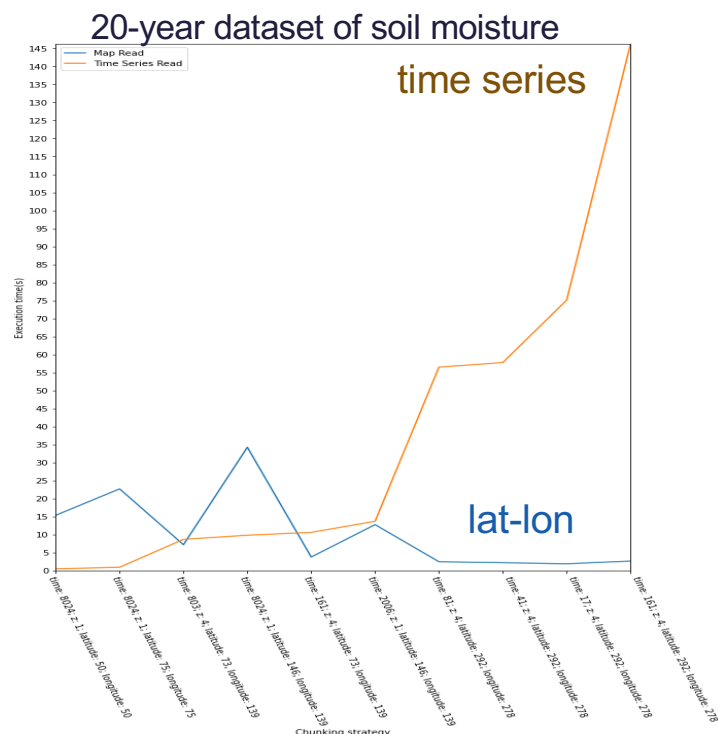
Amazon Sustainability Data Initiative

The Amazon Sustainability Data Initiative (ASDI) seeks to accelerate sustainability research and innovation by minimizing the cost and time required to acquire and analyze large sustainability datasets. ASDI supports innovators and researchers with the data, tools, and technical expertise they need to move sustainability to the next level.



Object Store: Different storage strategies showed radically different performance

Adapted from
slide by Phil
Kershaw



- Experiment with different storage chunking arrangements
- Way in which data is written and stored has significant impact on performance when used
- Rewriting with alternate chunking is fast, needs planning

Summary

- Host of national and regional facilities available for local science communities
 - Some outreach for specific projects
 - Existing collaborations
- Need to advertise these facilities and encourage their use
 - Common tool sets appearing based around core set used by Pangeo
- Commercial cloud spreading general access to other regions
 - Cost and funding
 - Access is not universal