



Data Infrastructure in Climate Sciences

Subtitles: "Challenges" in terms of accessibility and usage

Sébastien Denvil (IPSL/CNRS)





IPSL climate modelling centre (ICMC) http://icmc.ipsl.fr







0.5

0.4

0.3

0.2

0.1

0

-0.1

-0.2

-0.3

-0.4

-0.5

6,0

5.0

÷ 4.0

3.0

2,0

0,0

-1.0

1900

2000

Q

ne 1,0



Node

2100



c) Global (deviations from zonal mean)



IPSL climate model



<u>Climate Simulations</u>

CMIP6 (Coupled Model Intercomparison Projet phase 6), préparation duIPCC AR6 (Intergovernmental Panel Climate Change, Assesment Report 6)





Support to international coordinated experiments : CMIP5

Many, many processes, many,

many communities !

Interconnected communities, all needing access to (some of) the data!

IPCC AR5 variable counts

	1 hour	3 hour	6 hour	daily	month	annual	totals
aerosol	0	0	0	0	81	0	81
atmosphere	75	101	9	86	184	0	455
land	0	3	0	2	59	0	64
land ice	0	0	0	2	13	0	15
ocean	0	1	0	3	116	0	120
biogeochemistry	0	0	0	0	88	71	159
sea ice	0	0	0	4	47	0	51
totals	75	105	9	97	588	71	945

0.5

0.4

0.3

0.2

0.1

0

-0.1

-0.2

-0.3

-0.4

-0.5

6,0

5.0

÷ 4.0

3.0

2,0

0,0

-1.0

1900

2000

Q

ne 1,0

Node

2100

c) Global (deviations from zonal mean)

A one slide guide to CMIP5 from a data perspective

Fifth	World Climate	Original Timing:
Climate	Research Programme	o(2) PB of requested
Model	WCRP- WGCM	output from 20+
Intercomparison	Involves all the	modelling centres
Project	major climate	finished early 2010!
(CMIP5)	modelling <u>centres</u> .	Actual Timing?
		Years late.

Community developed s/w infrastructure for data delivery: Earth System Grid Federation

101 experiments61 model variants59,000 datasets!4.5 million files2 PB in global archive.Unknown PB locally!

Worldwide distributed system

The **Earth System Grid Federation (ESGF)** is a multi-agency, international collaboration of persons and institutions working together to build an <u>open source</u> software infrastructure for the management and analysis of Earth Science data on a global scale

- Software development and project management: ANL, ANU, <u>**BADC</u>**, <u>**CMCC**</u>, <u>**DKRZ**</u>, ESRL, GFDL, GSFC, JPL, <u>**IPSL**</u>, ORNL, LLNL (lead), PMEL, ...</u>
- Operations: tens of data centers across Asia, Australia, Europe and North America

Worldwide distributed system

Storage evolution in 6 years time (from CMIP3 to CMIP5) : a factor x30

- Operational since 2011
- Hundreds of users per month
- Hundreds of To per month
- About 10 000 registered users

CMIP3: centralized

CMIP5: distributed system

- 60 climate models
- 2 PB of data

System Architecture

ESGF is a system of <u>distributed</u> and <u>federated</u> Nodes

that interact <u>dynamically</u> through a Peer-To-Peer (P2P) paradigm

<u>Distributed</u>: data and metadata are published, stored and served from multiple centers ("Nodes")

<u>Federated</u>: Nodes interoperate because of the adoption of common services, protocols and APIs, and establishment of mutual trust relationships

<u>Dynamic</u>: Nodes can join/leave the federation dynamically - global data and services will change accordingly

A client (browser or program) can start from any Node in the federation and discover, download and analyze data from multiple locations as if they were stored in a single central archive.

Software Stack

Internally, each ESGF Node is composed of services and applications that collectively enable metadata discovery, data access, and user management. Software components are grouped into 4 areas of functionality (aka "flavors"):

•<u>Data Node</u> : secure data publication and access

•<u>Index Node</u> :

- metadata indexing and searching (w/ Solr)
- •web portal UI to drive human interaction
- <u>Identity Provider</u> : user authentication and group membership
- <u>Compute Node</u> : analysis and visualization

European coordination

InfraStructure for Earth System modelling

IS-ENES & IS-ENES2 EU projects

1rst phase - 03/2009 to 02/2013 18 partners 2nd phase - 04/2013 to 03/2017 23 partners

Global & Regional climate models Key role of infrastructure : models, data & computing

Recommandations:

1)Access to world-class HPC for climate at least «tailored » for climate up to « dedicated »

2)Develop the next generation of climate models

3)Set up data infrastructure (global and regional models) for large range of users from impact community

4)Improve physical network (e.g. link national archives)

5)Strengthen European expertise and networking

Input to IS-ENES2

ENES

Towards an European Climate Infrastructure Initiative : a sustainable virtual laboratory

0.5

0.4

0.3

0.2

0.1

0

-0.1

-0.2

-0.3

-0.4

-0.5

6,0

5.0

÷ 4.0

3.0

2,0

0,0

-1.0

1900

2000

Q

ne 1,0

Node

2100

c) Global (deviations from zonal mean)

To keep in mind

"the potential to interpret, compare and reuse climate information results is strongly related to the quality of their description"

Computation useless if results cannot be stored/distributed/read

Earth System Documentation

A climate simulation

http://earthsystemcog.org/projects/es-doc-models/

<u>CIM</u>

- The CIM is **intentionally very general**. It can be customized for particular user communities through the addition of specific Controlled Vocabularies.
- A Controlled Vocabulary defines the content that can be used within a CIM document. For example, in the case of climate models, the CIM schema (structure) allows a ModelComponent to have a child ModelComponent. And each of those components can have "types."
- A CV is required to list the permitted types. For example, the CMIP5 CV allows an "atmosphere" model to have a child "advection" model, but not a child "ocean" model. Thus, in order to be valid a CIM document must conform both to the CIM schema and to a particular set of CVs.

ES-DOC tools Guided Tour

ES-DOC

Tools - Document Viewer

CMIP5 Model - HadGEM2-ES

Model Simulation Experiment Platform

Overview	Citations	Contacts	Components	Grids
Project	CM	P5		
Short Name	Had	GEM2-ES		
Long Name	Had	ley Global Envir	onment Model 2 - I	- Earth System
Institute	UK	Met Office Hadle	ey Centre	
Funder	UK	Met Office Hadle	ey Centre	
Principal Inve	stigator Chr	is Jones		
Release Date	200	9-11-26 00:00:0	0	
Language				
Description	The (lea Had bias Thr syst sub con asp glob tree Aus are well com imp (DV schi The nutr of D add che trop for o mod	HadGEM2-ES I ding to HadGEM (GEM2-AO proje- ues. The latter had bugh focussed were ematic errors in stantially improve timental warm bi- ects of ENSO are al climate indices s, and the produ- tralia though mare good and agree with the C4MIP sponent carbon roved. The ocea IS) emissions fro- eme is an impro- se have differen- ients. The HadC MS. The diat-Ha- titions of a tropose mistry and sulpho- ospheric ozone limate forcing. I fel.	model was a two st M2-AO) and the add act targeted two key ad a particularly hig working groups a nu HadGEM1, such a red mean SSTs and as in HadGEM1 ha re improved. Overa es. [2] In HadGEM2 uctivity is better tha ay cause problems well with observed fluxes validate bette an biology (HadOCC om phytoplankton.) vement over the st nt processes for rem DCC scheme perfor adOCC scheme ha spheric chemistry s nate aerosols have distribution and the including interactive	stage development from HadGEM1, representing improvements in the physical model didion of earth system components and coupling (leading to HadGEM2-ES). [1] The ay features of performance: ENSO and northern continent land-surface temperature igh priority in order for the model to be able to adequately model continental vegetation. Number of mechanisms that improved the performance were identified. Some known as the Indian monsoon, were not targeted for attention in HadGEM2-AO. HadGEM2-AO in wind stress and improved tropical SST variability compared to HadGEM1. The northern as been significantly reduced. The power spectrum of El Nino is made worse, but other all there is a noticeable improvement from HadGEM1 to HadGEM2-AO when comparing I2-ES the vegetation cover is better than in the previous HadCM3LC model especially for an in the non-interactive HadGEM2-AO model. The presence of too much bare soil in s for the dust emissions scheme. The simulation of global soil and biomass carbon stores ded estimates except in regions of errors in the vegetation cover. HadGEM2-ES compares fels. The distribution of NPP is much improved relative to HadCM3LC. At a site level the tter against observations and in particular the timing of the growth season is significantly CC) allows the completion of the carbon cycle and the provision of di-methyl sulphide . DMS is a significant source of sulphate aerosol over the oceans. The diat-HadOCC standard HadOCC scheme as it differentiates between diatom and non-diatom plankton. emoving carbon from the surface to the deep ocean, and respond differently to iron orms well with very reasonable plankton distributions, rates of productivity and emissions ias slightly too low levels of productivity which requires further tuning to overcome. The scheme, new aerosol species (organic carbon and dust) and coupling between the e significantly enhanced the earth system capabilities of the model. This has improved the he distributions of aerosol species compared to observations, both of which are impor

ES-DOC

Tools – Document Comparator step 1

es-doc Earth System Documentation

Project CMIP5 - Comparator Model Component Properties -

Open

Support

Step 1 : Select Model Component Properties

1. Select Models	
ACCESS1.0	view
ACCESS1.3	view •
BCC-CSM1.1	view
CCSM4	view 🔍
CFSV2-2011	view
CMCC-CESM	view
смсс-см	view =
смсс-смѕ	view
CNRM-CM5	view
CSIRO-MK3.6.0	view 🔍
EC-EARTH	view
GFDL-CM2P1	view
GFDL-CM3	view 🔍
GFDL-ESM2G	view
GFDL-ESM2M	view
GFDL-HIRAM-C180	view
GFDL-HIRAM-C360	view
GISS-E2-H	view

. Select Components	υN
Aerosols	•
Emission & Concentration	•
Model	•
Transport	•
Atmosphere	••
Convection Cloud Turbulence	••
Cloud Scheme	••
Cloud Simulator	•
Dynamical Core	••
Advection	••
Orography & Waves	••
Radiation	
Atmospheric Chemistry	•
Emission & Concentration	•
Gas Phase Chemistry	•
Heterogen Chemistry	•
Stratospheric	•
Tropospheric	•
Photo Chemistry	•
- /	_

		Help	Reset	Next		
	3. Select	Properties				
	Scientif	ic Properties				
	Aero	osol Scheme				
		Bin Framework				
Bin Species						
		Bulk Species				
		Framework				
		Modal Framework				
		Modal Species				
		Scheme Characte	ritics			
		Scheme Type				
		Species				
	Cou	pling With				
	Gas					
	Ocean Biogeochemical Coupling					
	Processes					
	Veg	etation Model Cou	ıpling			
	Standar	d Properties				
	Cita	tions				
		Location				
		Title				
	Description					
Long Name						
	PLE	nail Address				

ES-DOC

Tools - Document Comparator step 2

es-doc Earth System Documentation

Project CMIP5 - Comparator Model Component Properties -

Help

Open

CSV

Support

Back

Step 2 : View report table

??? = Incomplete documentation. N/A = Not applicable (model did not realize component).

Component	Aerosols > Model	Atmosphere > Radiation	Atmosphere > Radiation
Property	Scientific Properties > Aerosol Scheme > Bulk Species	Scientific Properties > Aerosol Types	Scientific Properties > Longwave > Number Of Spectral Intervals
ACCESS1.0	BC (black carbon / soot) Dust POM (particulate organic matter) SOA (secondary organic aerosols) Sea salt Sulphate	BC (black carbon / soot) Dust POM (particulate organic matter) SOA (secondary organic aerosols) Sea salt Sulphate	9
ACCESS1.3	BC (black carbon / soot) Dust POM (particulate organic matter) SOA (secondary organic aerosols) Sea salt Sulphate	BC (black carbon / soot) Dust Nitrate POM (particulate organic matter) SOA (secondary organic aerosols) Sea salt Sulphate	9
3CC-CSM1.1	N/A	BC (black carbon / soot) Dust Organic POM (particulate organic matter) Sea salt Sulphate	8
CCSM4	BC (black carbon / soot) Dust Organic Sea salt Sulphate	BC (black carbon / soot) Dust Organic Sea salt Sulphate	8
CFSV2-2011	N/A	BC (black carbon / soot) Dust Organic Sea salt Sulphate	16
CMCC-CESM	N/A	Sulphate	16
CMCC-CM	N/A	Sulphate	16
CMCC-CMS	N/A	Sulphate	16
CNRM-CM5	N/A	BC (black carbon / soot) Dust loe Organic POM (particulate organic matter) Sea salt Sulphate	16
CSIRO-MK3.6.0	BC (black carbon / soot) Dust Organic POM (particulate organic matter) SOA (secondary organic aerosols) Sea salt Sulphate	BC (black carbon / soot) Dust Ice POM (particulate organic matter) Sea salt Sulphate	10
EC-EARTH	N/A	BC (black carbon / soot) Dust Organic Sulphate	6
GFDL-CM2P1	N/A	BC (black carbon / soot) Dust Organic Sea salt Sulphate	10
GFDL-CM3	BC (black carbon / soot) Organic POM (particulate organic matter) SOA (secondary organic aerosols) Sulphate	BC (black carbon / soot) Dust Organic Sea salt Sulphate	10
GFDL-ESM2G	N/A	BC (black carbon / soot) Dust Organic Sea salt Sulphate	10
GFDL-ESM2M	N/A	BC (black carbon / soot) Dust Organic Sea salt Sulphate	10
			10

Models = 42, Components = 2, Properties = 3

High Performance Data (HPD) Analysis Environment

Mutualized

Jointly *delivered* by →IPSL laboratories. Joint *users* (initially): →IPSL community Joint *users* (target): →French Academic community

Access services to ESGF System Users don't have to find, download, and keep up to date the data they need

CMIP5, CORDEX Reanalysis, Obs4MIPs

Analysis capabilities Environmental Data Compute Service Web Service Provision for : →Climate Science →Earth Observation →Environmental studies

Big DATA Platform

Collaboration Environment

- \rightarrow Access to Curated Archive.
- → Large shared "Group Work Spaces"

 \rightarrow climate analysis enabled system \rightarrow + 1 PB of high performance disk coupled to hundreds of cores configured for analysis

Thank you for your attention

