

Status and Outlook for the CMIP Data Request

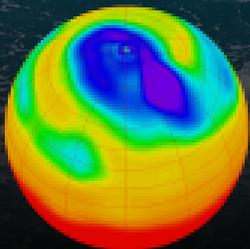
Martin Jukes

Centre for Environmental Data Analysis



1001000
10102001
101010 100101

December 4-7, 2018



**Centre for Environmental
Data Archival**

SCIENCE AND TECHNOLOGY FACILITIES COUNCIL
NATURAL ENVIRONMENT RESEARCH COUNCIL



is-enes

INFRASTRUCTURE FOR THE EUROPEAN NETWORK
FOR EARTH SYSTEM MODELLING

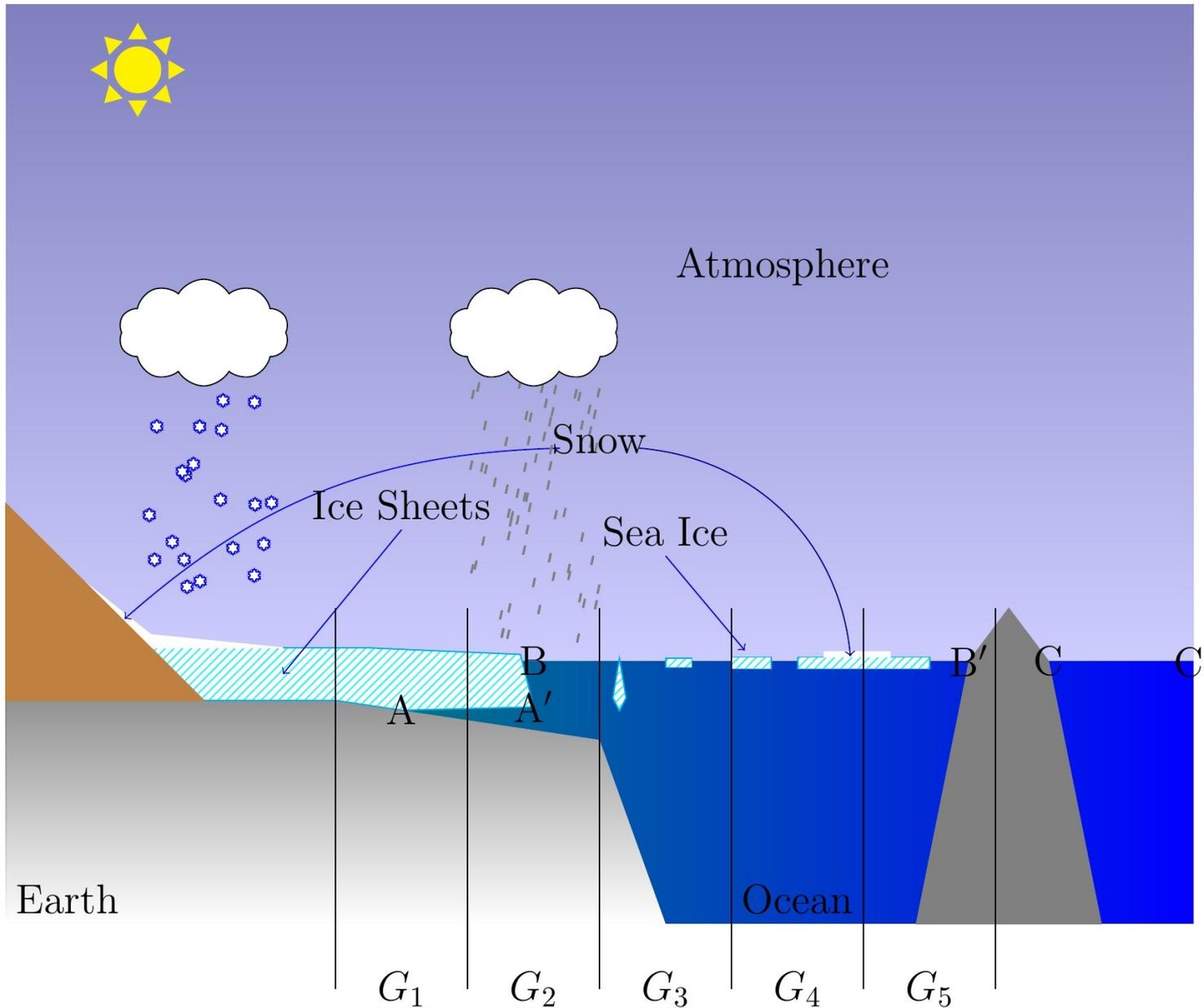
- Martin Juckes
- Karl Taylor
- Matthew Mizielinski
- Alison Pamment
- Sébastien Denvil
- Stéphane Sénési

Outline

- CMIP6 Data Request
 - What is the Data Request? [2]
 - Information Content [5]
 - Components [4]
 - Work-flows [1]
- CMIP7
 - Initial lessons [2]
 - Outlook [2]

- Consolidated specifications of output variables and requirements for 25 MIPs across 287 experiments;
- Reference document, web site and python package;

<http://w3id.org/cmip6dr>

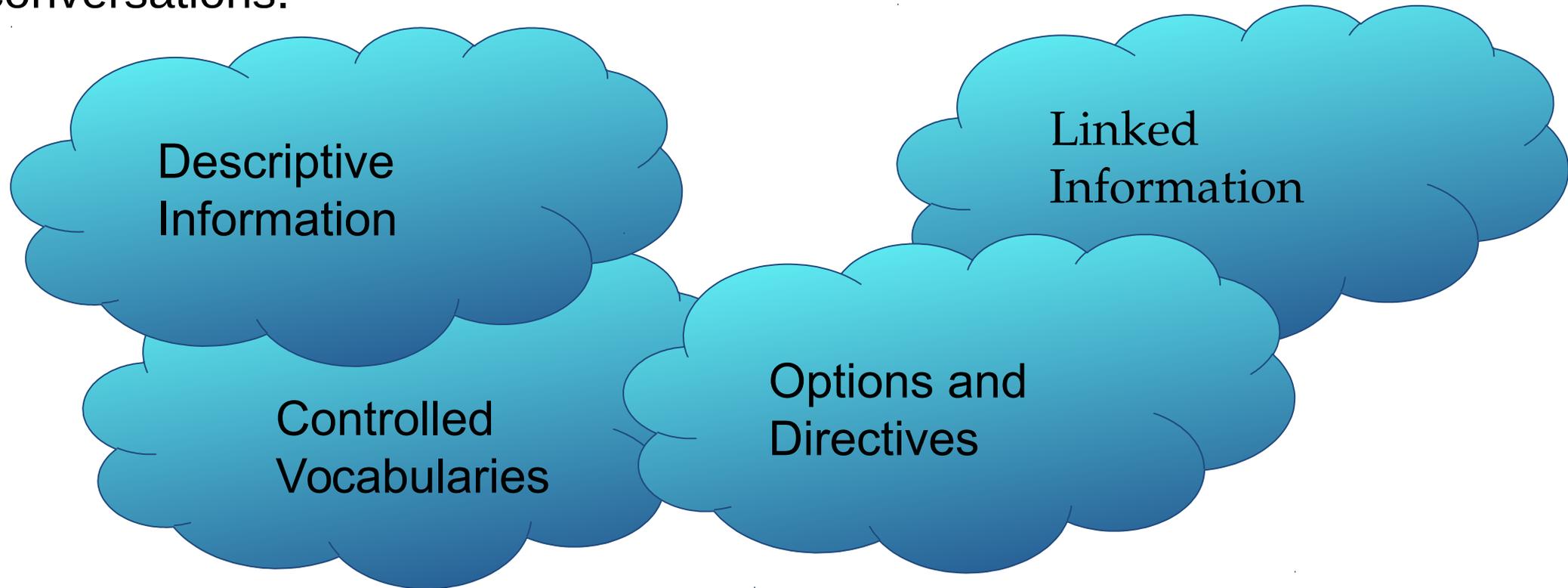


CMIP6 Data Request

- Defining the diagnostics to be archived;
- Many new surface categories, especially in the cryosphere (e.g. `floating_ice_shelves`,

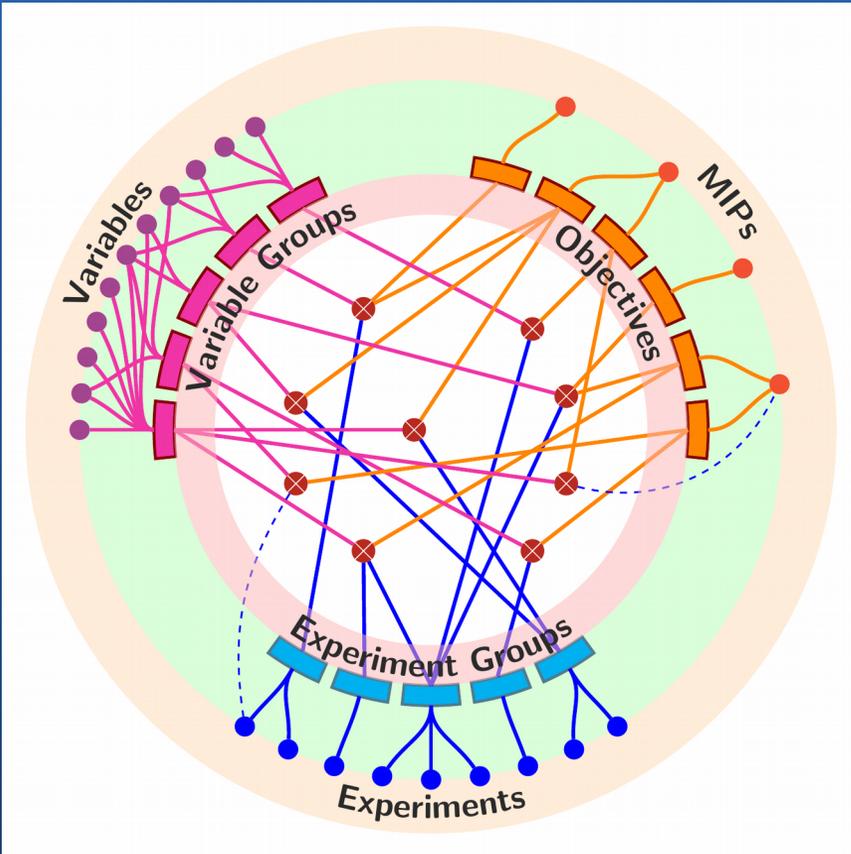
Information Content of the CMIP6 Data Request

The information in the Data Request comes in many forms, each requiring a different range of conversations.



Linked Content

- The **variable group** DCPP-Amon (36 variables) is requested by DCPP from the historical and DCPP-A **experiments** to support the “hindcasts” **objective** (i.e. “To assess the decadal predictability and forecast skill of forced and internally generated climate”).



Descriptive Content

Example

- **Label:** fNproduct
- **Title:** Deforested or Harvested Biomass as a Result of Anthropogenic Land Use or Change
- **Description:** When land use change results in deforestation of natural vegetation (trees or grasslands) then natural biomass is removed. The treatment of deforested biomass differs significantly across models, but it should be straight-forward to compare deforested biomass across models.

- Descriptive information contains “plain language” (sometimes with specialised scientific terminology) describing, for example, a diagnostic

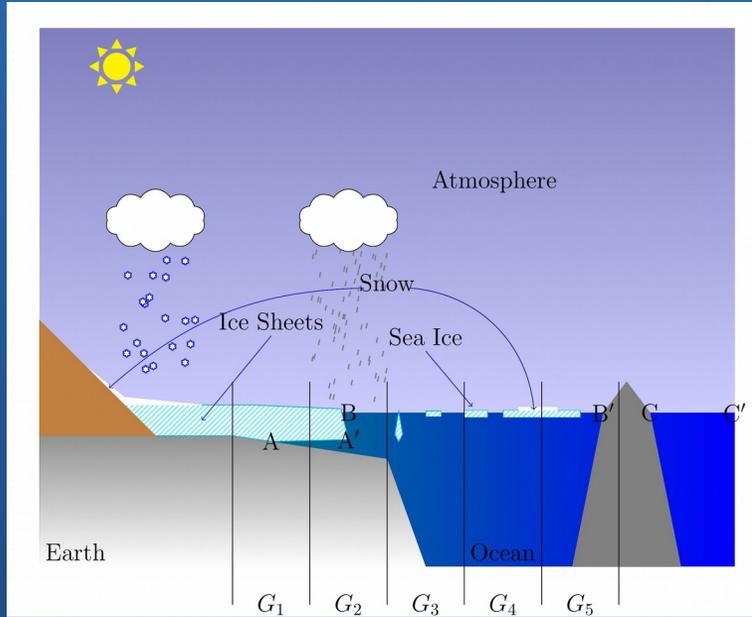
Controlled Vocabularies

- Activity
- Experiment
- Frequency
- Grid Label
- Institution
- Nominal Resolution
- Realm
- Source (model)
- Source Type
- Sub-experiment
- Table

- The request relies heavily on external vocabularies, such as CF standard names, CMIP Controlled Vocabularies, CMOR directives.
- Following the evolutions of these vocabularies while consolidating input from multiple sources is a significant challenge.



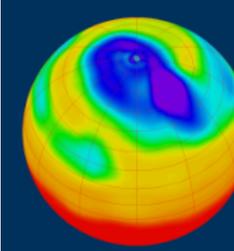
Options and Directives



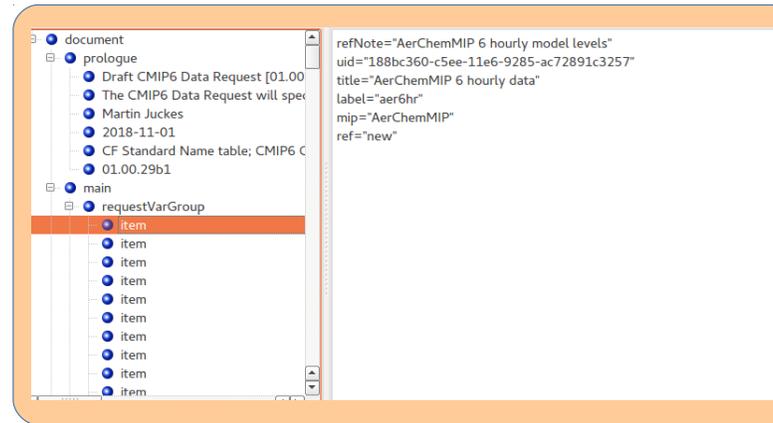
Some attributes represent values that need to be inserted into file metadata, other indicate options or instructions.

- `cell_measures = "--MODEL"` --> data should be archived on cell position used in model (e.g. vertex, edge or centre) rather than being interpolated to centre.
- `dimensions = "latitude, longitude"` --> data should be archived on a spatial grid appropriate to the parameter in question.

Accessing the data request



- The request is delivered through a range of formats
- The primary reference XML document is around 20,000 records
- Python library with command line and API
- Web pages



Reference Document

Web pages

CMIP6 Data Request

Home

3.1 Request variable group: a collection of request variables: [scenarioMipBaseline] Baseline set of variables for ScenarioMIP experiments

[Home](#) → [3.1 Request variable group: a collection of request variables section index](#)

- refNote: rvg1
- uid: x09999_6554f543-81af-11e8-b252-1c4d70487308
- title: Baseline set of variables for ScenarioMIP experiments
- label: scenarioMipBaseline
- mip: [mip] CMIP [CMIP]
- ref: auto

[Usage summary for Baseline set of variables for ScenarioMIP experiments \[scenarioMipBaseline\]](#)

Links from other sections

1.4 Request variable (carrying priority and link to group)	3.3 Request link: linking a set of variables and a set of experiments
--	---

1.4 Request variable (carrying priority and link to group)

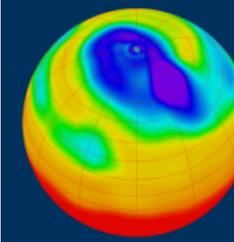
- abs550aer.AERmon [1]: ambient aerosol absorption optical thickness at 550 nm {groups: 7, vars: 1}
- agesno.Limon [1]: Snow Age {groups: 3, vars: 2}
- agesco.Omon [1]: Sea Water Age Since Surface Contact {groups: 7, vars: 2}
- arag.Oyr [1]: Aragonite Concentration {groups: 8, vars: 2}
- aragpfc.Rp10: Field Call Area for Atmospheric Variables {groups: 9, vars: 4}

Modeling centres have their own code to extract the information they need from the request, using the XML document directly, the python library, or parsing the output of the command line tool.

XML reference document

- The full request is provided in a single structured document, ensuring portability and enabling use in data production environments.
- Flat “table” structure: 35 sections, each consisting of a list of “item” records;
- Constrained to enforce typing of data values (e.g. float, integer, integer list, etc.
- Schema largely fixed since early 2016.

Web Pages



Web pages support searching and browsing of content.

The web site provides access to tables of variables for each MIP and experiment.

The search page allows variables to be filtered on variable name, long name, standard name, units and description.

Units: [kelvin] Kelvin

[Home](#) → [Units section index](#)

- uid[i]: fd70554e-3468-11e6-ba71-5404a60d96b5
- title[i]: Kelvin
- text[i]: K
- label[i]: kelvin

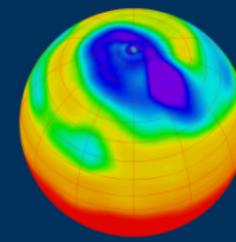
Links from other sections

1.2 MIP Variable

1.2 MIP Variable

Variable	Standard name	Long name	Units	Description	uid
<input type="text" value="\$t"/>	<input type="text"/>	<input type="text" value="flux "/>	<input type="text"/>	<input type="text"/>	<input type="text"/>
tau_	magnitude_of_surface_downward_stress	[Momentum flux] (N m-2) (module of the momentum lost by the atmosphere to the surface.)		<590e6e4c-9e49-11e5-803c-0d0b866b59f3>	
tendacabf_	tendency_of_land_ice_mass_due_to_surface_mass_balance	[Total surface mass balance flux] (kg s-1) (The total surface mass balance flux over land ice is a spatial integration of the balance flux)		<4146943a-4f40-11e6-a814-ac72891c3257>	
tendlibmassbf_	tendency_of_land_ice_mass_due_to_basal_mass_balance	[Total basal mass balance flux] (kg s-1) (The total basal mass balance flux over land ice is a spatial integration of the basal mass balance flux)		<414699d0-4f40-11e6-a814-ac72891c3257>	
tendlicalf_	tendency_of_land_ice_mass_due_to_calving	[Total calving flux] (kg s-1) (The total calving flux over land ice is a spatial integration of the calving flux)		<41469f3e-4f40-11e6-a814-ac72891c3257>	

Python Library



Available in svn or Python Package Index: `pip install dreqPy`

Use in python code:

```
from dreqPy import scope
sc = scope.dreqQuery()
v1 = sc.volByMip2( 'C4MIP', pmax=2 )
v2 = sc.volByMip2( {'C4MIP', 'LUMIP'}, pmax=2 )
```

From the LINUX command line:

```
drq -m C4MIP,LUMIP -p 1 -t 1 → 4.20Tb
drq -m HighResMIP -p 1 -t 1 → 29.0Tb
drq -m HighResMIP:DiurnalCycle -p 1 -t 1 → 3.9Tb
```

Workflows

Input from:

- MIPs;
- Modelling groups;
- Teams working on CMOR, ES-DOC

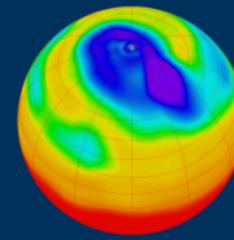
Communication through:

- spreadsheets (excel and google);
- Emails;
- Discussion forum (years 1 & 2);
- Github issues;

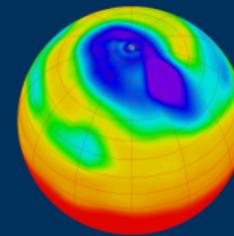
Since June 2018:

- Beta releases and reviews to improve stability of new releases.

Successes

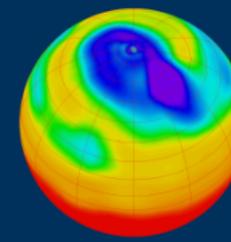


- Schema allowing consolidation of multiple MIP requirements;
 - Variable definitions re-used (927 standard names, 1272 MIP variables, 2063 CMOR variables) enabling robust definitions;
 - Permits selection of variables based in scientific objectives, priorities, and experiments;
- Re-usable meta-data structures, e.g. cell methods, units, spatial structures;
 - robust validation of syntax;
- Multiple viewing formats giving flexibility to users;
 - e.g. web site for search and browsing, XML for embedding in work flows, python API for systematic exploration;



- 1. Lack of stability early on;**
- 2. Slow trickle of requirements;**
 - e.g. revisions to reference vocabularies;**
- 3. Too much complexity;**
- 4. Divergent approaches from MIPs;**
 - developing the data request at the same time as the new “endorsed MIP” process was being introduced added confusion and stress.**
- 5. Time to solution too long;**

File names



*Sub-expt.
(optional)*

For CMIP6:

tas_Amon_BCC-CSM2-HR_1pctCO2_r1i1p1f1_gn_198001-198412.nc

Variable

Table

Source
(model)

Experiment

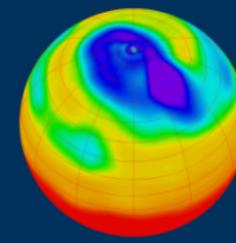
Member

Grid id.

Time
range

The file naming scheme ensures unique file names for each data product. The name can be split into 3 categories of information: (1) the parameter (green), (2) the experiment (blue) and (3) implementation details (red).

Variable names



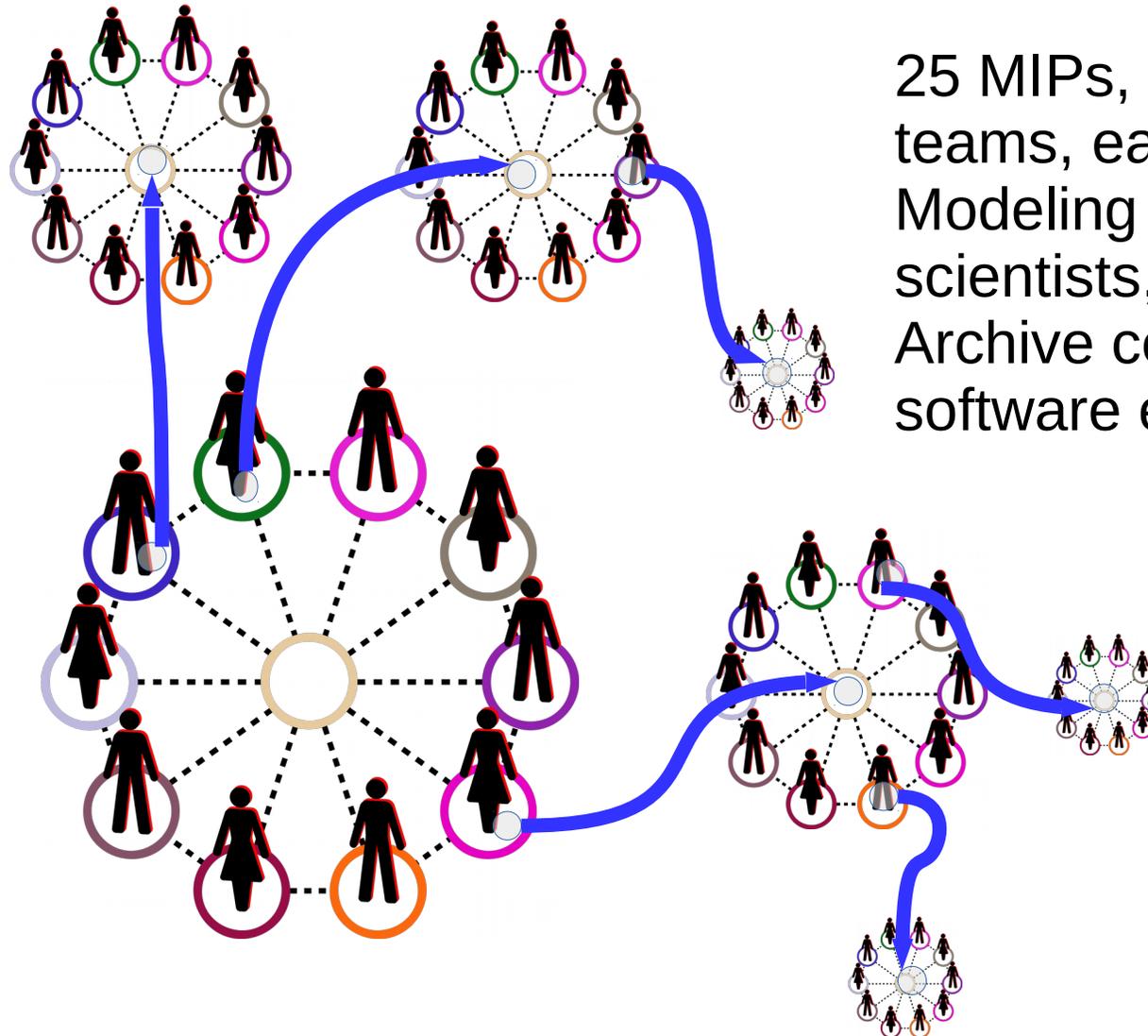
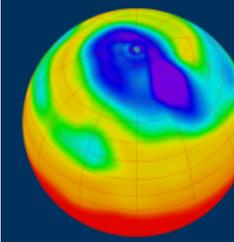
Variable names must be consistent with unique file names;
Precise rules for variables depend on the file naming convention;
Discussion of variable names starts early file naming
convention is finalised later in the process;

However:

Rules for variables only depend on part of the file naming
convention (the “variable name + table” pair) which could be
finalised earlier

tas_Amon_BCC-CSM2-HR_1pctCO2_r1i1p1f1_gn_198001-198412.nc

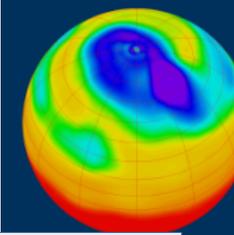
Communication: the challenge



25 MIPs, each involving multiple research teams, each team with multiple specialists;
Modeling centres: managers, environmental scientists, data scientists, software engineers;
Archive centres: standards, data scientists, software engineers.

Different communication strategy needed for different user groups.

Communication: approaches



Early CMIP6

Emails (lists for distribution:
individual for queries);
Forum site for discussion;
Spreadsheets;

Late CMIP6

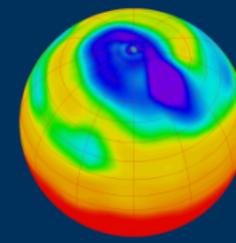
Emails (lists for distribution:
individual for queries);
Github for issues;
Spreadsheets + google sheets;

CMIP7 ideas

Emails (lists: general for distribution
+ **MIP specific for queries**);
Github for issues;
Google sheets? (*with python
integration?*);

The screenshot shows a forum interface with a navigation bar at the top containing 'Data Request', 'Discussions', and 'Activity'. The user 'martinjukes' is logged in. The main content area is titled 'dreq01' and is described as a 'Trial forum for community discussion about the CMIP6 data request'. It features a 'New Discussion' button and a list of discussion threads. Each thread includes a title, a status (e.g., 'Closed'), the author, the date, and the number of views and comments.

Discussion Title	Status	Author	Date	Category	Views	Comments
Tuning overview tables wrt to set of MIPs and simulated experiments	Closed	martinjukes	October 2016	CMIP6 Data Request: Schema	1	2
Status of xlsx doc on CMOR_variables	Closed	martinjukes	July 2016	CMIP6 Data Request: Schema	1	2
No DECK data request ?	Closed	martinjukes	July 2016	CMIP6 Data Request: Schema	1	2
dreqPy (Data request python API) : Need for basic examples	Closed	martinjukes	July 2016	CMIP6 Data Request: Schema	1	1
dimension definitions missing for some ocean variables.	Closed	martinjukes	July 2016	CMIP6 Data Request: Schema	1	1
01.beta.30	Closed	martinjukes	July 2016	CMIP6 issues	1	1
missing standard names.	Closed	nadeau1	June 2016	CMIP6 Data Request: Schema	2	0
standard_name issues for same variable	Closed	nadeau1	June 2016	Structure of the CMIP6 data request	1	0
Release 01.beta.27	Closed	martinjukes	May 2016	CMIP6 issues	1	1
Release 01.beta.26	Closed	martinjukes	April 2016	CMIP6 issues	11	2
Organisation of files generated by CMOR	Closed	martinjukes	April 2016	CMIP6 issues	1	3



Initial meeting in June 2018 (see http://bit.ly/dreq_rev2018);
Discussed finalisation of CMIP6 request and outlook for the future.

- Rationalisation: e.g. fewer grid specification options;
- Clearer schema;
- Clearer interfaces to other components;

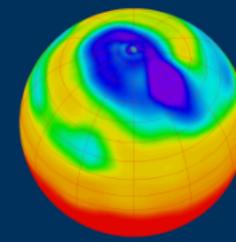
Funding for a further meeting is allocated in IS-ENES3;



is-enes

INFRASTRUCTURE FOR THE EUROPEAN NETWORK
FOR EARTH SYSTEM MODELLING

Recommendations



Data Request Decisions

- Rationalise;
- Streamline interactions with CF;
- Improve tracking of discussions (all discussions in a tracked environment);

Request to WIP

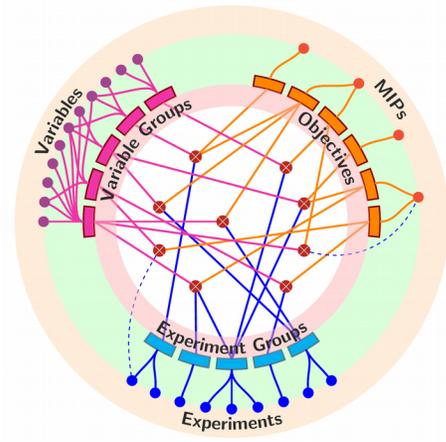
- Agree framework for parameter component of file names for CMIP7 early (i.e. 2019), so that variable naming rules are clear;
- Agree a communication framework (e.g. github, forum, email lists: thisMIP@dreq.org) and insist on registration by interested MIPs;
- Establish meaningful metrics of progress;

Request to WGCM/CMIP panel

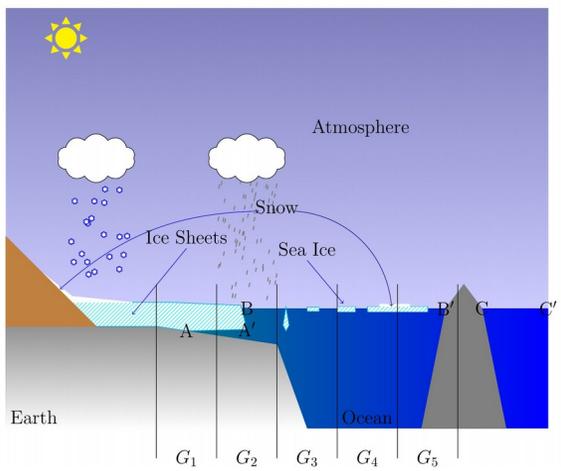
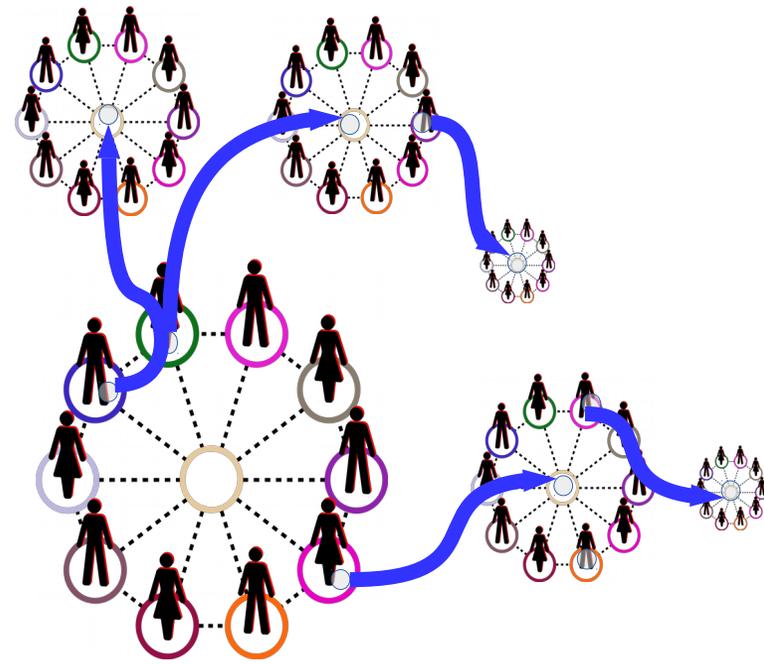
- Clarity around responsibilities of MIPs;
- Encourage coordination among related MIPs (perhaps through core projects);

Outlook

- Funding for further work through IS-ENES3;
- Plan for more MIPs in CMIP7 (e.g. hydrology, fire);
- Explore options for supporting more projects that want to use



The End



Successes

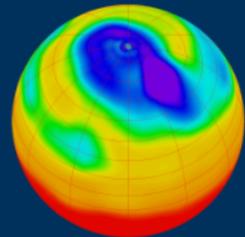
- Schema allowing consolidation of multiple MIP requirements;
- Re-usable meta-data structures;
- Multiple viewing formats;
- Automated tests of units, cell methods, attribute types and much more;

Problems

- 1. Lack of stability early on;**
- 2. Slow trickle of requirements;**
- 3. Too much complexity;**
- 4. Divergent approaches from MIPs;**
- 5. Revisions to reference vocabularies;**
- 6. Time to solution too long;**

Status and Outlook for the CMIP Data Request

Martin Juckes, Center for
Environmental Data Archival

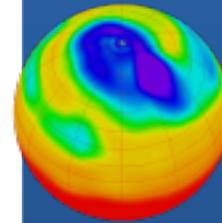


**Centre for Environmental
Data Analysis**

SCIENCE AND TECHNOLOGY FACILITIES COUNCIL
NATURAL ENVIRONMENT RESEARCH COUNCIL



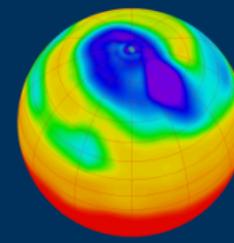
is-enes
INFRASTRUCTURE FOR THE EUROPEAN NETWORK
FOR EARTH SYSTEM MODELLING



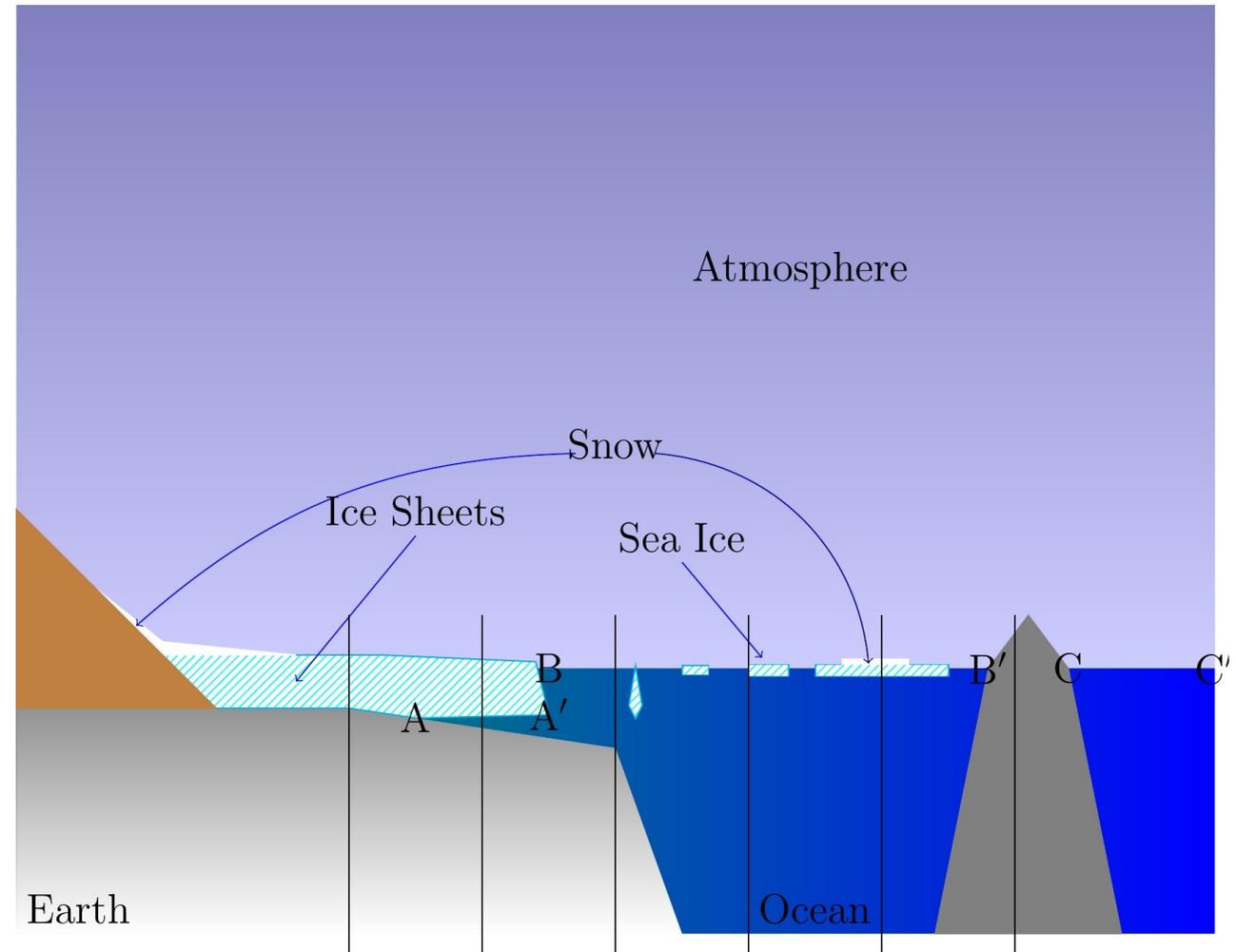
**Centre for Environmental
Data Archival**

SCIENCE AND TECHNOLOGY FACILITIES COUNCIL
NATURAL ENVIRONMENT RESEARCH COUNCIL

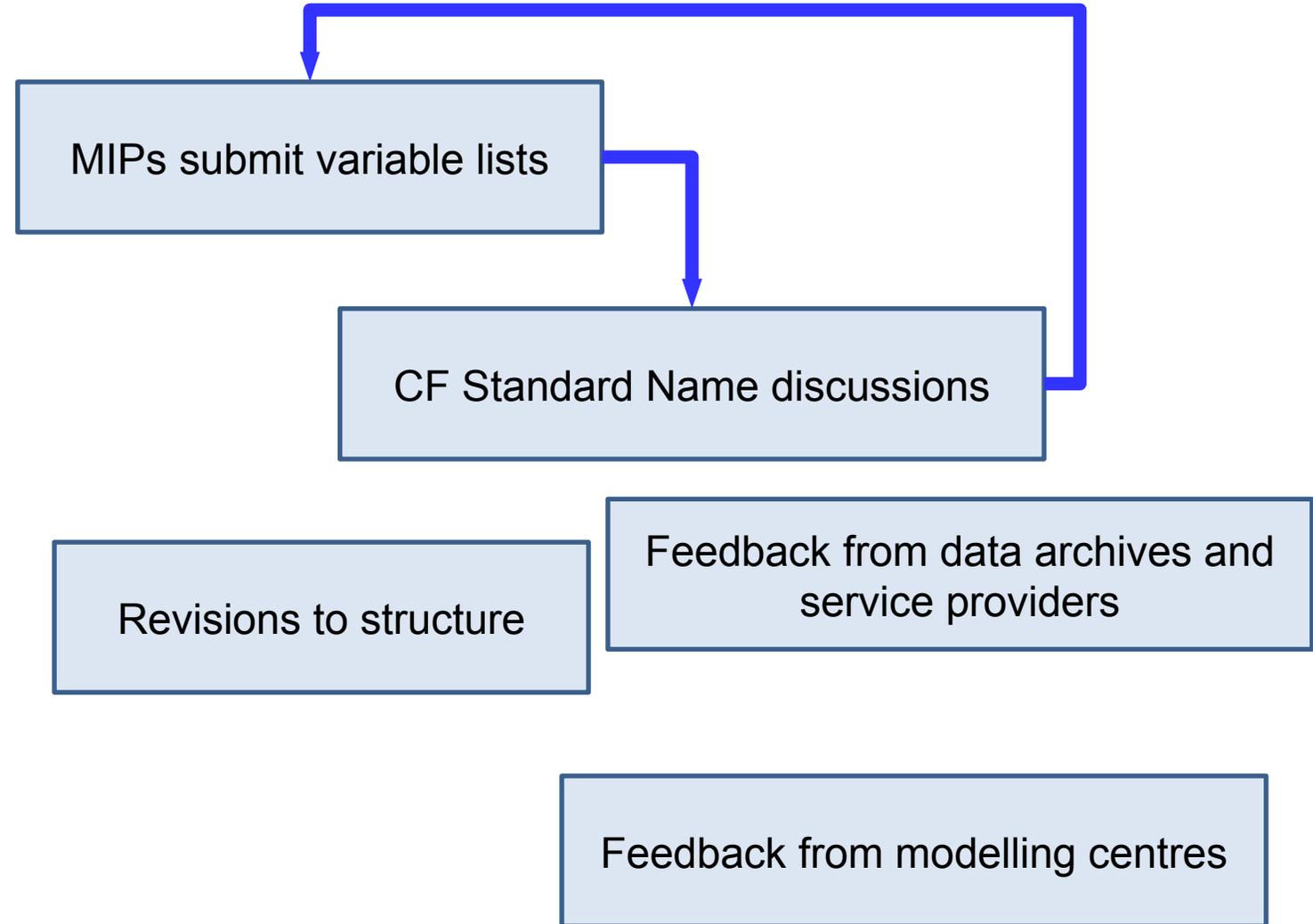
<http://w3id.org/cmip6dr>



New surface types to consider,
e.g.:
floating ice shelves,
grounded ice sheets
snow on sea ice

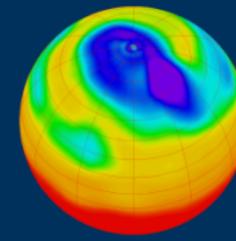


Workflows



Aerosol	177	153	129
Atmosphere	753	325	256
Atmospheric Chemistry	18	13	11
Land	365	291	191
Landice	174	78	57
Ocean	260	165	148
Ocean Biogeochemistry	287	214	151
Sea Ice	104	92	79

Lessons ...



Lessons:

Streamline interactions with CF;

MIPs need a technical point of contact who understands the relevant standards;

Need a clear framework for organising content early on;

For many scientists, email is the most effective means of communication;

Lessons:

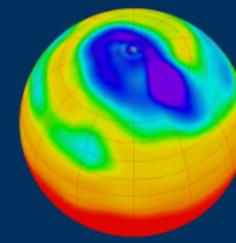
Streamline interactions with CF;

MIPs need a technical point of contact who understands the relevant standards;

Need a clear framework for organising content early on;

For many scientists, email is the most effective means of communication;

Scope of CMIP7



It is usual for the climate modelling community to start talking about limiting the scope of CMIP, especially when the stress of dealing with the current phase is high.

But the world is warming ... CMIP5 provided scenario runs with coupled carbon cycles, CMIP6 adds nitrogen cycles, and greater variability in the biosphere and cryosphere.

CMIP7: expect greater complexity in models, e.g. resolving food productivity, and moving from science questions (what is the environmental impact of forcing X?) to policy questions (what is societal impact of policy decision Y?)

Lessons:

- Need to simplify the request: remove complex grid preferences; clarify linkages;
- Variable definitions can be re-used by other MIPs;