## CMIP6 and its data infrastructure

Advanced School on Earth System Modeling
Indian Institute of Tropical Meteorology, Pune

V. Balaji

NOAA/GFDL and Princeton University

20 July 2016

# Outline

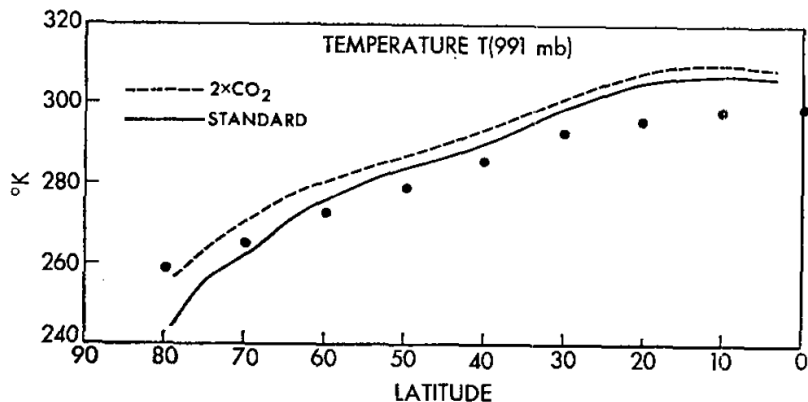# Atmospheric response to doubled $CO_2$



Fig 5 from Manabe and Wetherald (1975), a foundational document of climate modeling.

# The Charney Report (1979)

"Carbon dioxide and climate: A Scientific Assessment."

- Precursor to the IPCC Assessment Reports.
- Based on 5 model runs: 3 from Manabe (GFDL), 2 from Hansen (GISS).
- Conclusions:
    - Direct radiative effects due to doubling of $CO_2$: $\sim 4$ W/m$^2$
    - Feedbacks: water vapor (Clausius-Clapeyron), snow-ice albedo feedback.
    - Cloud effects: "How important the cloud effects are, is, however, an extremely difficult question to answer. The cloud distribution is a property of the entire climate system, in which many other feedbacks are involved."
    - "We believe, therefore, that the equilibrium surface warming will be in the range of 1.5-4.5°C, with the most probable value near 3°C."

Very nice reassessment of the Charney Report: Bony et al (2013).

# World Climate Research Programme established 1980

As a consequence of the Charney Report, the WCRP was established in 1980.

- Under the joint sponsorship of International Council for Science and the World Meteorological Organization, and the Intergovernmental Oceanographic Commission of UNESCO
- Its Joint Scientific Committee (JSC) provides scientific guidance to WCRP – about 18 scientists chosen worldwide.
- WCRP can help inform national and international funding priorities for climate science, including observations and modeling
- Working groups: WGCM (coupled modeling); WGNE (numerical experimentation); WGSIP (seasonal-interannual prediction). Many cross-WG activities.
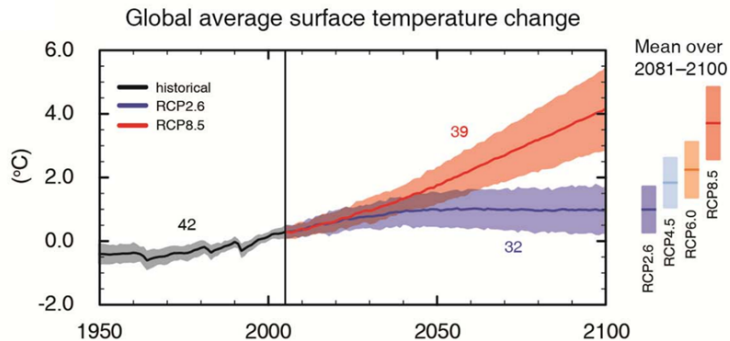
**http://wcrp-climate.org/**

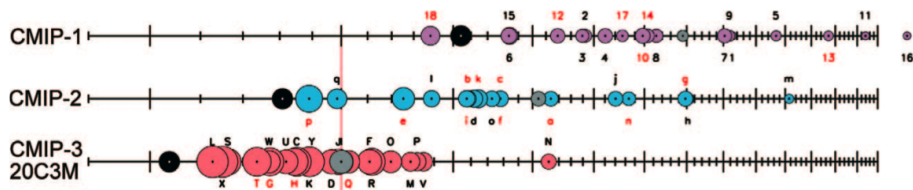Figure SPM.7 from the IPCC AR5 Report. 20th century warming cannot be explained without greenhouse gas forcings.

# WCRP Grand Challenges

The WCRP's organizing principle is a set of evolving grand challenges:

- Clouds, Circulation and Climate Sensitivity
- Melting Ice and Global Consequences
- Climate Extremes
- Regional Sea-level Change and Coastal Impacts
- Water Availability

"... targeted research efforts with the likelihood of significant progress over 5-10 years", at "a specific barrier preventing progress in a critical area of climate science".
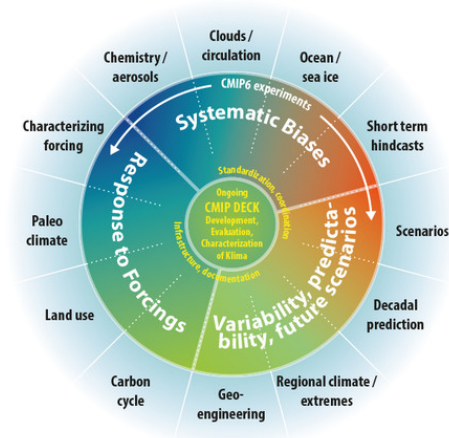
# CMIP established 1990 under the guidance of WGCM



Reichler and Kim (2008), Fig. 1: compare models' ability to simulate 20th century climate, over 3 generations of models.

- Models are getting better over time.
- The ensemble average is better than any individual model.
- Improvements in understanding percolate quickly across the community.

# CMIP6 design: DECK and MIPs



DECK experiments form the core; many specialized MIPs for smaller communities, some 24 of which have been endorsed by CMIP panel. Figure courtesy Meehl et al (*Eos* 2014).

# CMIP6 Scientific Design

**Coupled Model Intercomparison Project Phase 6 (CMIP6): Design and Organization**

**Veronika Eyring, Jerry Meehl, Bjorn Stevens, Ron Stouffer, Karl Taylor** (CMIP Panel)

**Sandrine Bony and Cath Senior** (WGCM Co-chairs)

**V. Balaji** (WGCM Infrastructure Panel co-chair with K. Taylor)

16 January 2015 *(updates to CMIP6 Data Request Timeline on Slide 9)*
Please see the CMIP Panel website for additional information and updates:
http://www.wcrp-climate.org/index.php/wgcm-cmip/about-cmip

**Contact for questions:** CMIP Panel Chair Veronika Eyring (email: Veronika.Eyring@dlr.de)

The final CMIP6 Design, possibly with small modifications to the here presented figures and wording, will be published in a GMD Special Issue together with a description of the CMIP6-Endorsed MIPs and the forcing datasets. This Special Issue will open 30 April 2015.
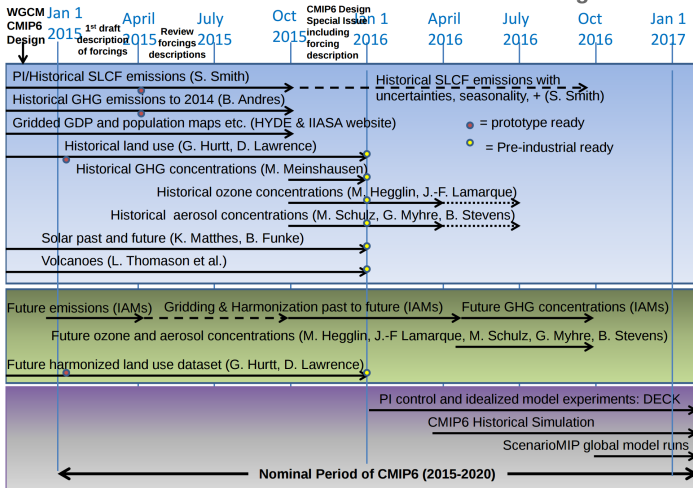
Final experimental design and data request adopted at WGCM19 Meeting, October 2015.
**http://goo.gl/FMYRKe**

# CMIP6 input data (forcings) timeline



Finalize scenario choice, March 2015 (O'Neill, Tebaldi, van Vuuren)

**CMIP6 Forcing Timeline**

PI/Historical SLCF emissions (S. Smith)

Historical SLCF emissions with uncertainties, seasonality, + (S. Smith)

Historical GHG emissions to 2014 (B. Andres)

Gridded GDP and population maps etc. (HYDE & IIASA website)

● = prototype ready

○ = Pre-industrial ready

Historical land use (G. Hurtt, D. Lawrence)

Historical GHG concentrations (M. Meinshausen)

Historical ozone concentrations (M. Hegglin, J.-F. Lamarque)

Historical aerosol concentrations (M. Schulz, G. Myhre, B. Stevens)

Solar past and future (K. Matthes, B. Funke)

Volcanoes (L. Thomason et al.)

Future emissions (IAMs)  Gridding & Harmonization past to future (IAMs)  Future GHG concentrations (IAMs)

Future ozone and aerosol concentrations (M. Hegglin, J.-F Lamarque, M. Schulz, G. Myhre, B. Stevens)

Future harmonized land use dataset (G. Hurtt, D. Lawrence)

PI control and idealized model experiments: DECK

CMIP6 Historical Simulation

ScenarioMIP global model runs

**Nominal Period of CMIP6 (2015-2020)**

WGCM CMIP6 Design — Jan 1 2015 — 1st draft description of forcings — April 2015 — Review forcings descriptions — July 2015 — Oct 2015 — CMIP6 Design Special Issue including forcing description — Jan 1 2016 — April 2016 — July 2016 — Oct 2016 — Jan 1 2017

## http://goo.gl/FMYRKe
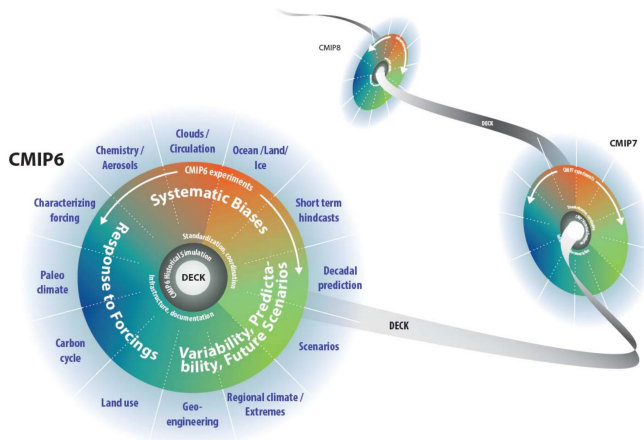
# CMIP evolution

DECK is designed to evolve slowly or not at all.



- IPCC Assessment Reports are snapshots of the "state of the science", but not directly linked to CMIP.
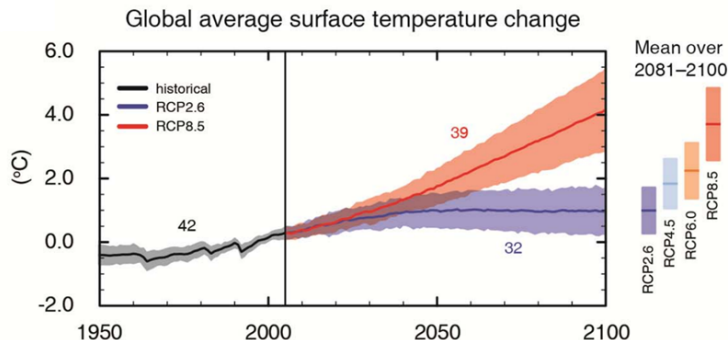
# IPCC Timeline

All dates in red are official dates from IPCC plenary in Nairobi, 2016-04, and IPCC XC meeting 2016-05-19.

- 2022-09: AR6 Synthesis Report
- 2021-02: WG1 Report Approved
- 2020-12: Final WG1 Draft and SP goes to inter-governmental review
- 2020-06: 4th Lead Author meeting
- 2020-02: Post 3rd LA meeting, second-order draft sent out for expert review. Any citations here will have to have been submitted for peer review by this date. First-order draft can use pre-citation material.
- 2019-03: Data in public domain.
- 2018-09: GFDL runs complete.

Earlier special reports (1.5C, cryosphere, land) not based on CMIP6.

# Multi-model ensembles for climate projection



- Critically depends on software, metadata, and data standards: the Earth System Grid Federation (`http://esgf.org`): a 3 PB federated archive.
- Key technical issues like replication, versioning, subsetting, QC, citation.

# The global data infrastructure underpinning MIPs

- MIPs, and in general any science involving cross-model comparisons, critically depend on the global data infrastructure – the "vast machine" (Edwards 2010) – making this sort of data-sharing possible.
- Infrastructure should not be a research project.
- Infrastructure should be treated as such by the national and international research agencies, but it is instead funded piecemeal, as a soft-money afterthought. This places the system at risk (NRC 2012: "A National Strategy for Advancing Climate Modeling", ISENES-2 Infrastructure Strategy document, 2012.)

# Role of WGCM and its infrastructure panel

- Provide scientific guidance and requirements for the GDI; exert greater influence over its design and features.
- Provide standards governance allowing for orderly evolution of standards.
- Provide design templates (e.g CMOR extensions) for groups designing MIPs and work to ensure their conformance to standards.
- Work with academies and publishers to require adequate data citation and recognition for data providers.
- Intercede with national agencies to provision data infrastructure with adequate and stable long-term funding.

# WIP: The WGCM Infrastructure Panel formed 2014

- Chaired by V. Balaji (Princeton/GFDL) and K. Taylor (PCMDI).
- Strategy to develop a series of "position papers" on global data infrastructure and its interaction with the scientific design of experiments. These will be presented to WGCM annual meeting.
  - protocol document for the "endorsed MIPs" delivered. Working with CMIP panel and MIP sponsors on CMIP6 data request.
  - data access policies: would open access simplify the technical design of the infrastructure?
  - data citations. Developing and promoting a path to data citations using DOIs and the emerging data journals, such as ESSD, Nature Scientific Data.
  - projected data volumes for CMIP6, strategies for managing the growth path
- Close involvement of the WIP and CMIP panel (e.g. joint papers)
- Interest from other WCRP working groups! (WGSIP, WGNE)
- Covers not only ESGF requirements but also other tools: ESDOC, CMOR, CF Conventions, ..

## WIP Membership

- V. Balaji (co-chair): GFDL
- Karl Taylor (co-chair): PCMDI
- Luca Cinquini: NASA JPL
- Cecelia DeLuca: NOAA
- Sebastien Denvil: IPSL
- Mark Elkington: MOHC
- Francesca Guglielmo, LSCE
- Eric Guilyardi: IPSL
- Martin Juckes: BADC
- Slava Kharin: CCCma
- Michael Lautenschlager: DKRZ
- Bryan Lawrence : NCAS, BADC
- Dean Williams: PCMDI

a blend of computer and climate scientists representing data centers and modeling groups: rotating membership with overlapping 2-year cycles

## Why not carry on as in the past?

- Heavy reliance on a few individuals worked O.K. for CMIP5, but may fail for the distributed management envisioned for CMIP6

- Need a procedure for evolving the infrastructure in a coordinated way so that the many groups and projects developing it can be responsive to the scientific needs.

- A panel with broad expertise may more nimbly respond to future needs than relying on a few individuals to poll community experts and build a consensus.

- Modeling groups are tasked with meeting the MIP requirements and deserve formal input to define them.

- Anything done to ensure that standards are as uniform as possible across all MIPs will reduce the burden.

- Membership on an official panel might help individual members to fund their work in this area.

## WIP Mission

"to promote a robust and sustainable global data infrastructure in support of the scientific mission of the WGCM"

- Establish standards and policies for sharing climate model output
- ensure consistency across WGCM activities
- Extend standards as needed to meet evolving needs
- Review and provide guidance on requirements of the infrastructure (e.g. level of service, accessibility, level of security)
- Oversee
  - file formats, structure and metadata
  - controlled vocabularies, name spaces, and naming conventions
  - protocols for interfacing components of the infrastructure
  - URL and catalog standards
  - protocols for data publication (including version identification), node management and data harvesting
  - standardized descriptions of models and simulations
  - security protocol for authentication and authorization query formats.

Covers ESGF, DRS, CoG, CMOR, ESDOC, ...

# Position Paper: Formation of CDNOT

- WIP recommended to the WGCM and CMIP panel the formation of a technical consortium charged with operationalizing the CMIP6 ESGF Federation: the CMIP6 Data Node Operations Team (CDNOT).
- Distinct bodies (with overlapping membership) responsible for requirements (WIP), software development (ESGF, ESDOC, ...), and operations (CDNOT)
- Formation approved by WGCM and CMIP, June 2015.
- Sébastien Denvil appointed Chair of CDNOT.
- Many sites have proposed members: if you are planning to operate a CMIP6 data node, please contact Sébastien right away! CDNOT operations are imminent (as soon as ESGF 2.0 is released).

# Position paper: CMIP6 Data Request

Led by Martin Juckes, STFC.
Highlights:

- Data request now available in machine-readable formats, including XLS and XML.
- A python API to allow the building of workflow tools that can work directly with the data request (e.g for setting switches in the model or post-processing).
- Endorsed MIPs have provided input on how the data will be used and analyzed.
- Actions needed from MIPs: develop and share analytic capabilities related to data request.
- Actions needed from modeling groups: review data request and provide feedback re feasibility.

## CMIP6 Data Request: Complexity and Volume

- the CMIP6 Data Request is made available by Martin Juckes via the **drq** tool and data stored in XML. Python tools can be used to mine this in many ways, e.g:

$$\boxed{\texttt{drq -m GMMIP -xls}} \tag{1}$$

  (assumes Tier 1 experiments and Priority 1 variables only).

- Assuming we follow our plans as stated, GFDL will generate 1-1.5 PB of public data on ESGF. This assumes the use of netCDF4 lossless data compression (0.7 atm, 0.4 ocn).

- Total number of Priority 1 variables across MIPs and DECK: $\sim$1300 (this is the number of variables that will have to be QCed before publication).

## Position paper: Data reference structure: syntax, vocabularies, filenames and global attributes

Highlights: mostly follows CMIP5 with some additional items:

- Allows easier grouping and selection. For instance, runs distinguished only by forcings will now be seen as an ensemble (extension of `rip` to `ripf`).
- Notation for data regridded on standard grids (e.g 1x1, see below).
- Improved association of data across multiple files, e.g auxiliary `cell_measures` such as `volcello` (ocean cell volume).
  - `external_variables` now adopted in the CF convention as a general mechanism for variables in other files.
  - CMIP6 will use `further_info_url` to locate external variables.
- DCPP extensions to allow additional forecast lead time coordinate.
- More sophisticated tracking of datasets (see discussion of PIDs below, `tracking_id` is now deprecated)
- When these papers are finalized and released, modeling groups can incorporate these into their workflows.

# Other data and metadata standards

(Not in any WIP position paper as of now, but are WIP recommendations).

- The WIP recommends the use of netCDF4 with lossless compression as the data format for CMIP6.
  - Lossless compression from `zlib` (settings `deflate=2` and `shuffle`) expected to generate roughly 2X decrease in data volumes (varies depending on data entropy or noisiness).
  - Requires upgrading entire toolchain (data production and consumption) to netCDF4.
  - Potential loss in performance during data creation.
- The WIP recommends the use of standard grids for datasets where native-grid data is not strictly required. For example: the Clivar OMDP may request the use of WOA standard grids ($1° \times 1°$, $0.25° \times 0.25°$) as the target grid of choice.
- No progress on adoption of standard calendars.

# Grid diversity may increase in CMIP6



Downstream communities may not wish to deal with novel grids, but specialist communities are likely to insist on it for their own research.

# Model metadata

(Not in any WIP position paper as of now, but are WIP recommendations).

- ESDOC documents of model metadata are a required element in quality control and DOI generation.
- Considerably simplified questionnaire relative to CMIP5.
- Command-line tools (e.g `py-esdoc`) will be made available to make it easy to generate, clone, share CIM documents.
- Forcing documentation in the works.
- Should we include tuning documentation? cf. Hourdin et al BAMS paper, "The Art and Science of Model Tuning", soon to be published in BAMS. Outcome of 2014 tuning workshop.

# ESDOC Comparator Tool



ES-Doc comparator tool (CMIP5)

**http://compare.es-doc.org/**

## Position papers: Replication, versioning and errata

Main requirement is for end users to know if they are working with the right dataset in a federation where data is replicated multiple times, may have been retracted or superseded. Highlights:

- Extend use of persistent identifiers (PIDs) for dataset tracking (replaces `tracking_id` from CMIP5).
- Lists of PIDs can be used as supplementary citation information in papers.
- PID-based query system to see if errata have been reported, or data have been superseded.
- Proposal to ESGF working teams on how PIDs can be incorporated into replication workflow.

# Position paper: Data citation and long-term access

Highlights:

- Main requirement: ensure proper citation of data used in a study to acknowledge contributions by modeling groups.
- Automated QC mechanisms to ensure adherence to metadata and data quality standards.
- Commitment to long-term archival by at least some data centers.
- Links connecting datasets to model and experiment documentation (ESDOC/CIM)
- DOI generation at the granularity of *model* and *simulation*.
- Action needed from WGCM: endorse the requirement of data citation as part of the terms of use of CMIP6 model output.
- Recommendation to modeling groups: generate citations in the emerging data science journals e.g., Nature Scientific Data or ESSD. Possibly approach for special issue?

# Position paper: Data licensing and access control

- Main requirement: simplified access control on ESGF, data license applicable even when data is found in non-ESGF repositories.
- For CMIP6 data licenses will be embedded in the data files (netCDF global attribute). There will be choice of two different licenses (Creative Commons "share-alike" and "non-commercial share-alike")
- Recognition that many users will (and did) use data from secondary ("dark") repositories. Embedded license implies that user is subject to the terms of use no matter where they retrieved the data.
- Required action from WGCM: endorse the new WIP-recommended licensing policy.
- Required action from modeling groups: choose a license consistent with your own institutional policies and record in global file attributes. Let us know if the two recommended licenses are both unacceptable.

## Position paper: CMIP6 Data Volume

CDNOT member institutions and ESGF require realistic data volume estimates for hardware planning.

- A number of current estimates are based on an assumption of geometric progression (straight line on a log scale!) drawn through CMIP3 and CMIP5.
- Based on known growth in number of models, years simulated, and increase in resolution, the actual growth will likely be less.
- Some centres (e.g UKMO and GFDL) are developing tools to allow us (and possibly others) to make accurate data volume estimates based on Martin's data request documents, model resolution, experiment planning.
- WIP will release in 2016 best estimates of CMIP6 total data volume.

# CMIP6 Data Request: preliminary analysis

# WIP Position Papers: Current Status

- Recommended formation of CMIP6 Data Node Operations Team (CDNOT: Sébastien Denvil, Chair).
- Recommended use of netCDF4 lossless compression for CMIP6
- Data Citation and Long Term Access: DOIs issued for quality-controlled data at the granularity of model and simulation.
- Recommended use of Persistent Identifier (PID) at the dataset level. Allows tracking for datasets for replication, versioning, and errata.
- Simplified licensing and data access: licenses embedded in files (two options: open access and non-commercial use)
- Recommended use of standard grids (e.g 1x1) of limited set of high-value data.
- Standard format machine-readable data request for DECK and MIPs.
- Finalizing Data Reference Structure and Syntax (paths and controlled vocabularies) and netCDF attributes.
- Data volume estimates to be released after data request finalized.

## Summary: CMIP6

- the CMIP multi-model ensemble is an organizing principle and fundamental tool of climate science
- Addresses WCRP Grand Challenges through WGCM-led experimental design
- Design allows both for continuity across CMIP generations and new experiments in response to evolving science
- provides input to IPCC and other Assessment Reports
- WGCM Infrastructure Panel translates CMIP experimental design into requirements for the global data infrastructure
- Close involvement of WIP with ESGF-XC (overlapping membership)
- Communities: `cmip6-modelgroups-sci` and `cmip6-data-request`

## Summary: global data infrastructure

- WGCM Infrastructure Panel translates CMIP experimental design into requirements for the global data infrastructure
- Governance at different stages of infrastructure: requirements (WIP), software development (ESGF, ESDOC, CoG, CMOR, ...), CMIP6 implementation and operations (CDNOT).
- Close involvement of WIP with ESGF-XC and CDNOT (overlapping membership)
- WIP has produced 11 position papers (out of the promised 4), available on the WIP website. Data volume estimate paper soon to follow.
- Communities: WIP, CDNOT, `esgf-devel`

`https://www.earthsystemcog.org/project/wip/resources/`