

CMIP6: WIP Status Update

A summary of the CMIP6 project to-date from the WGCM Infrastructure Panel (WIP)

Paul J. Durack, Matthew Mizielinski and Karl E. Taylor

WGCM23

9th December 2020 - Virtual



Outline

- WIP background and membership transitions
- Timelines - calibrate where we've come from and where we are
- CMIP6
 - Infrastructure components
 - Successes to date
 - Components and their contribution
 - CVs, Data Request, CMOR, ES-DOC, ESGF Citation service
 - Lessons learned
- Conclusions

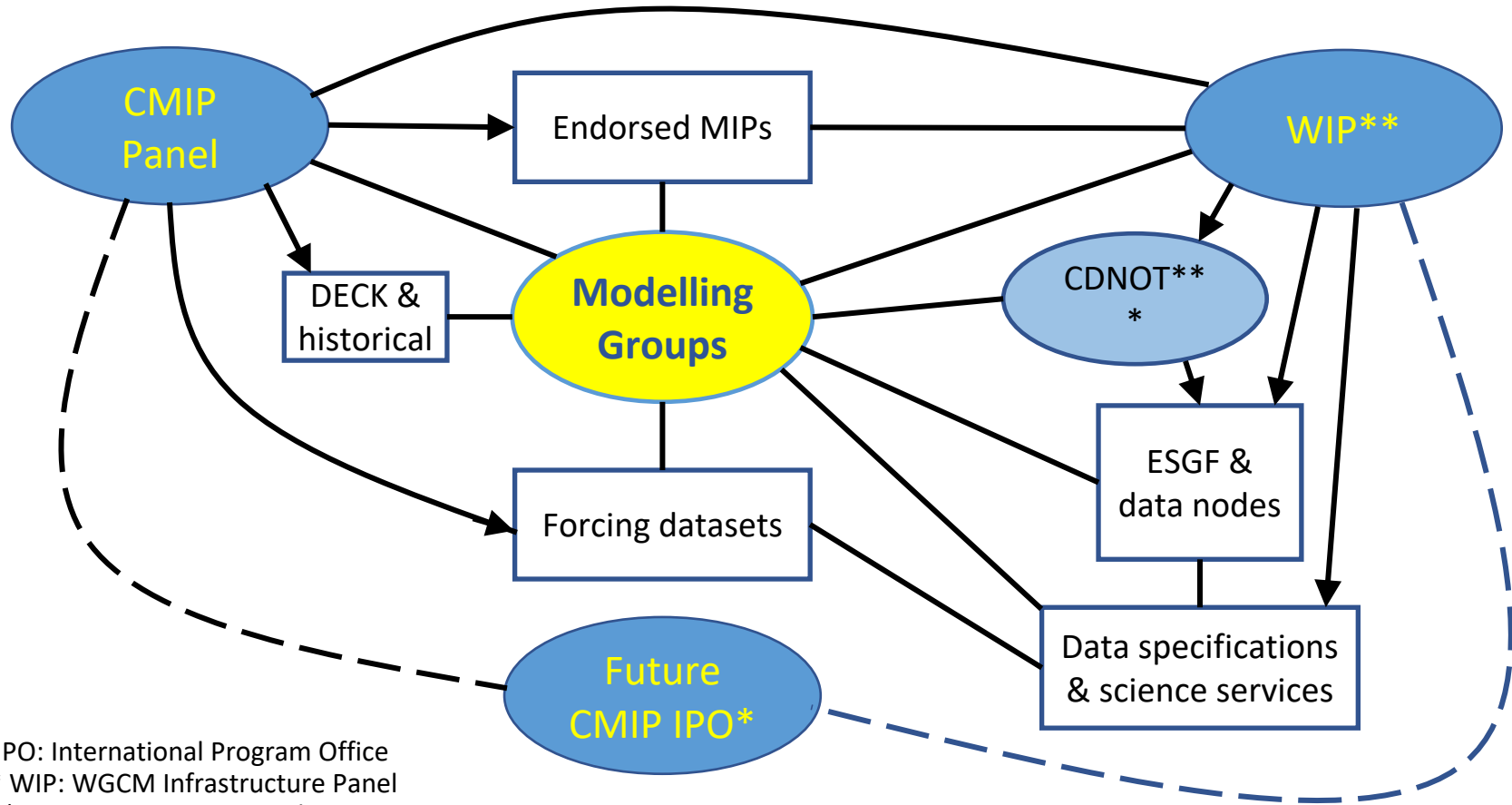
What is the WGCM Infrastructure Panel* (WIP)?

Established by the WGCM in June 2014 to:

- Set requirements ensuring the CMIP infrastructure will serve its purpose
- Write documents defining specifications for the infrastructure and data it hosts.
- Help coordinate development work done by funded infrastructure projects so that the elements work well together
- Communicate and coordinate with data managers at modeling groups via a “CMIP Data Node Operations Team” (CDNOT)

* G. Abdulla, S. Ames, Y. Bai,, P. J. Durack, D. Hassell M. Jukes, S. Kharin, M. Lautenschlager, M. Mizieliński, R. Petrie, M. Stockhouse, K. E. Taylor

WCRP/WGCM Organizational Structure for CMIP



* IPO: International Program Office

** WIP: WGCM Infrastructure Panel

*** CDNOT: CMIP Data Node Operations Team

Arrows indicate direct responsibility or oversight role

WIP membership transitions

- The WIP lost some key members about a year ago:
 - The WIP co-chair and the 2 P.I.s for ES-DOCs resigned
 - Multiple reasons were cited, but fundamentally it reflects the stress placed on the inadequately funded projects that support CMIP infrastructure
 - There was also a breakdown in communication across the CMIP management structure
 - It is a warning sign that all CMIP participants need to include infrastructure planning and support in all phases
 - The CDNOT chair accepted a position at ECMWF
- Ruth Petrie (CEDA) is the new CDNOT chair
- Matthew Mizielski (Hadley Centre) is the new WIP co-chair
- Next month Paul Durack (PCMDI) will replace Karl Taylor as the other WIP co-chair

Outline

- WIP background and membership transitions
- Timelines - calibrate where we've come from and where we are
- CMIP6
 - Infrastructure components
 - Successes to date
 - Components and their contribution
 - CVs, Data Request, CMOR, ES-Doc, ESGF Citation service
 - Lessons learned
- Conclusions

Timeline of the *MIPs

*MIPs temporal evolution and their support infrastructure

Planning begins	1989	1993	1995	1997	2003	2008	2014
Simulations	AMIP1	AMIP2	CMIP1	CMIP2/ CMIP2+	CMIP3	CMIP5	CMIP6
Data volume	1GB	500GB	1GB	500GB	50TB	2PB	>20PB
Host infrastructure	LLNL FTP*	LLNL FTP	LLNL FTP	LLNL FTP	LLNL FTP	ESGF 41# nodes	ESGF 31 nodes
Operations begin	1989	1995	1996	1999	2004	2011	2018

Next slide
focuses
on this
period

*For some groups in addition to data distribution, LLNL computing facilities were used to run contributing simulations

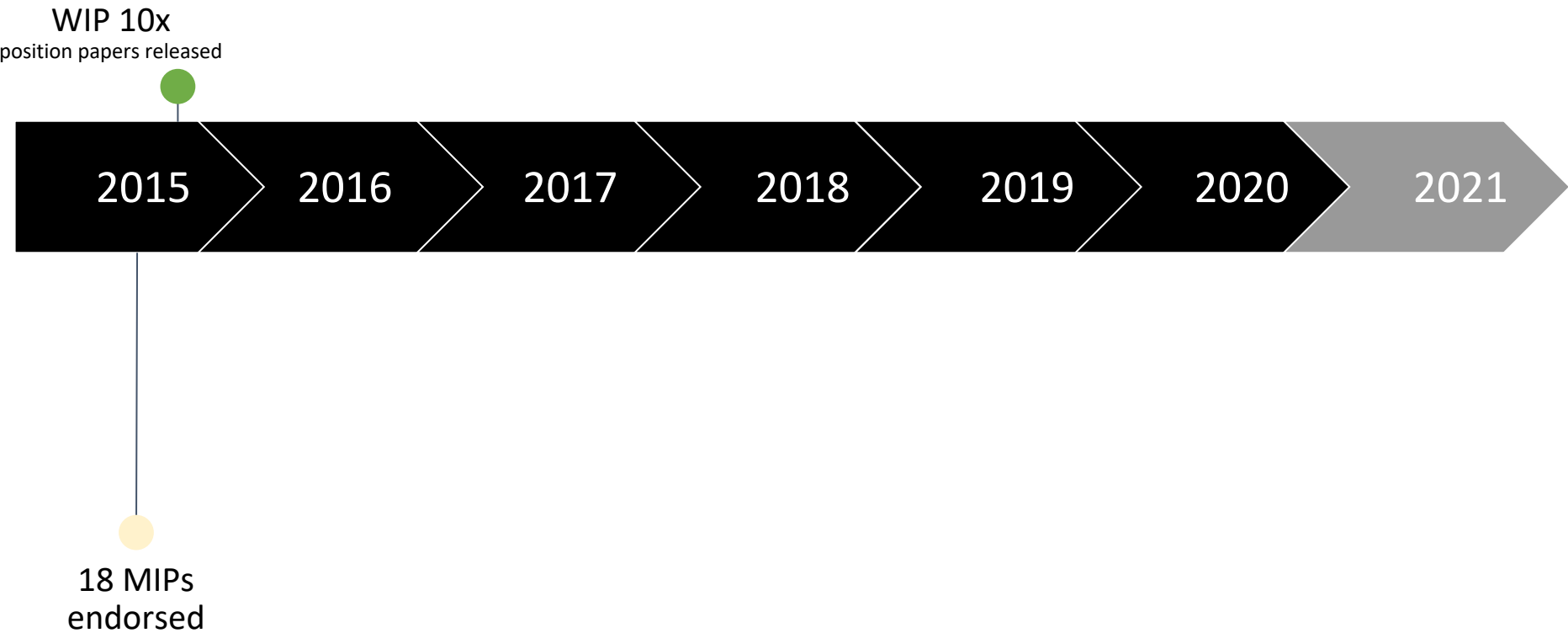
#Just 17 CMIP5 nodes remain to-date, which means a ~50% loss rate from the CMIP5 peak - highlights tier 1 node importance for data preservation

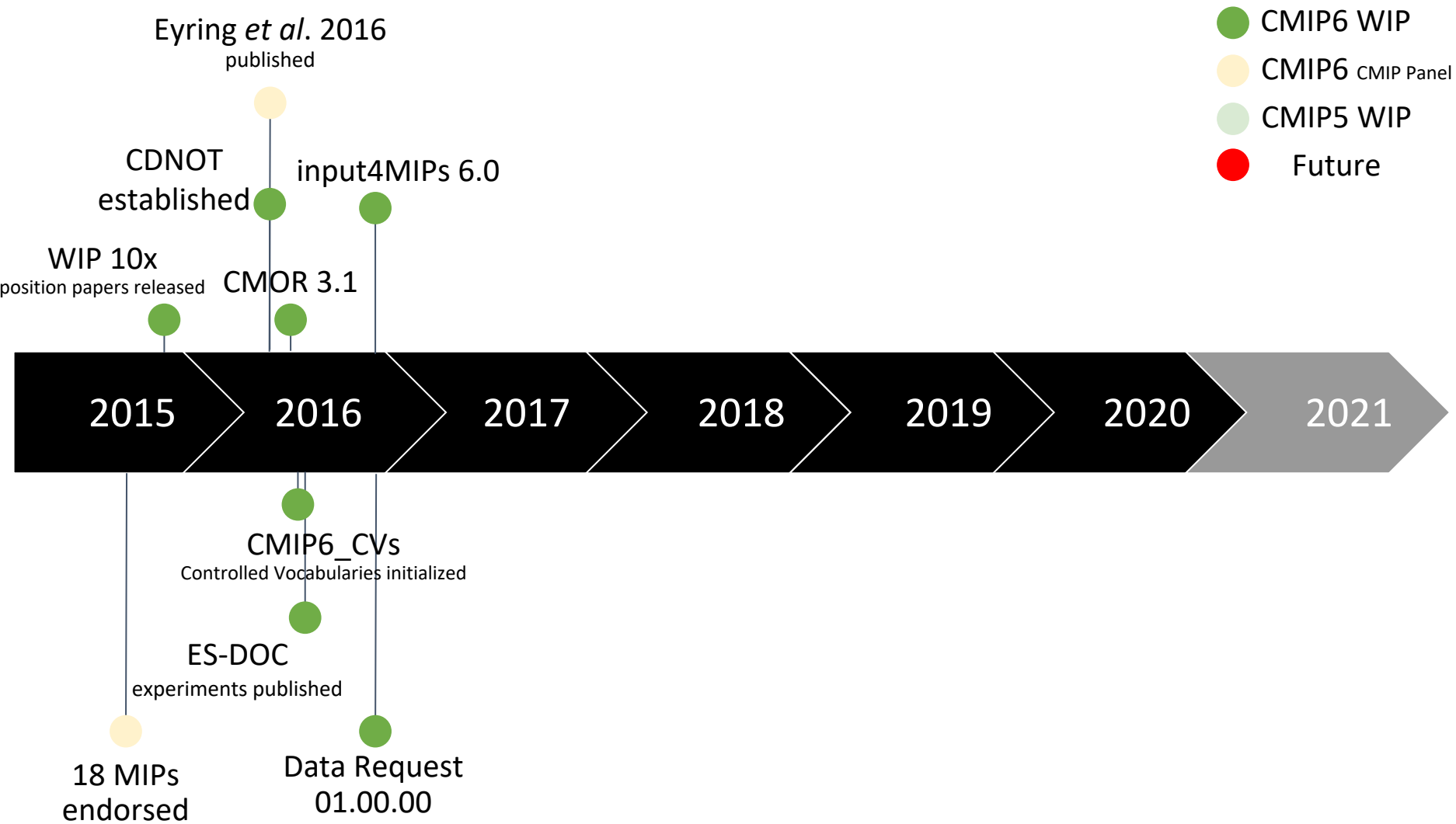
Timeline of CMIP6

- CMIP6 WIP
- CMIP6 CMIP Panel
- CMIP5 WIP
- Future

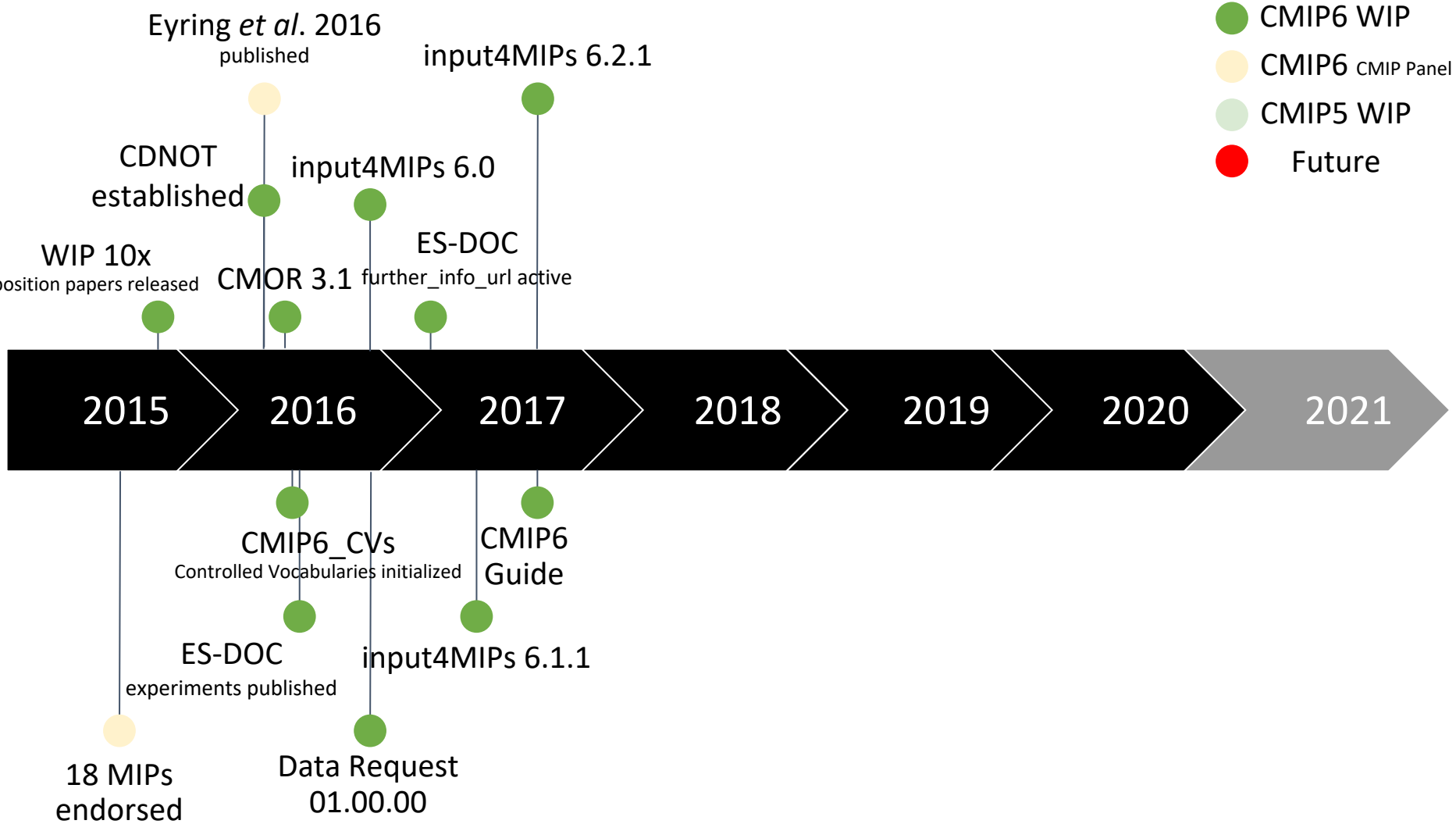


- CMIP6 WIP
- CMIP6 CMIP Panel
- CMIP5 WIP
- Future

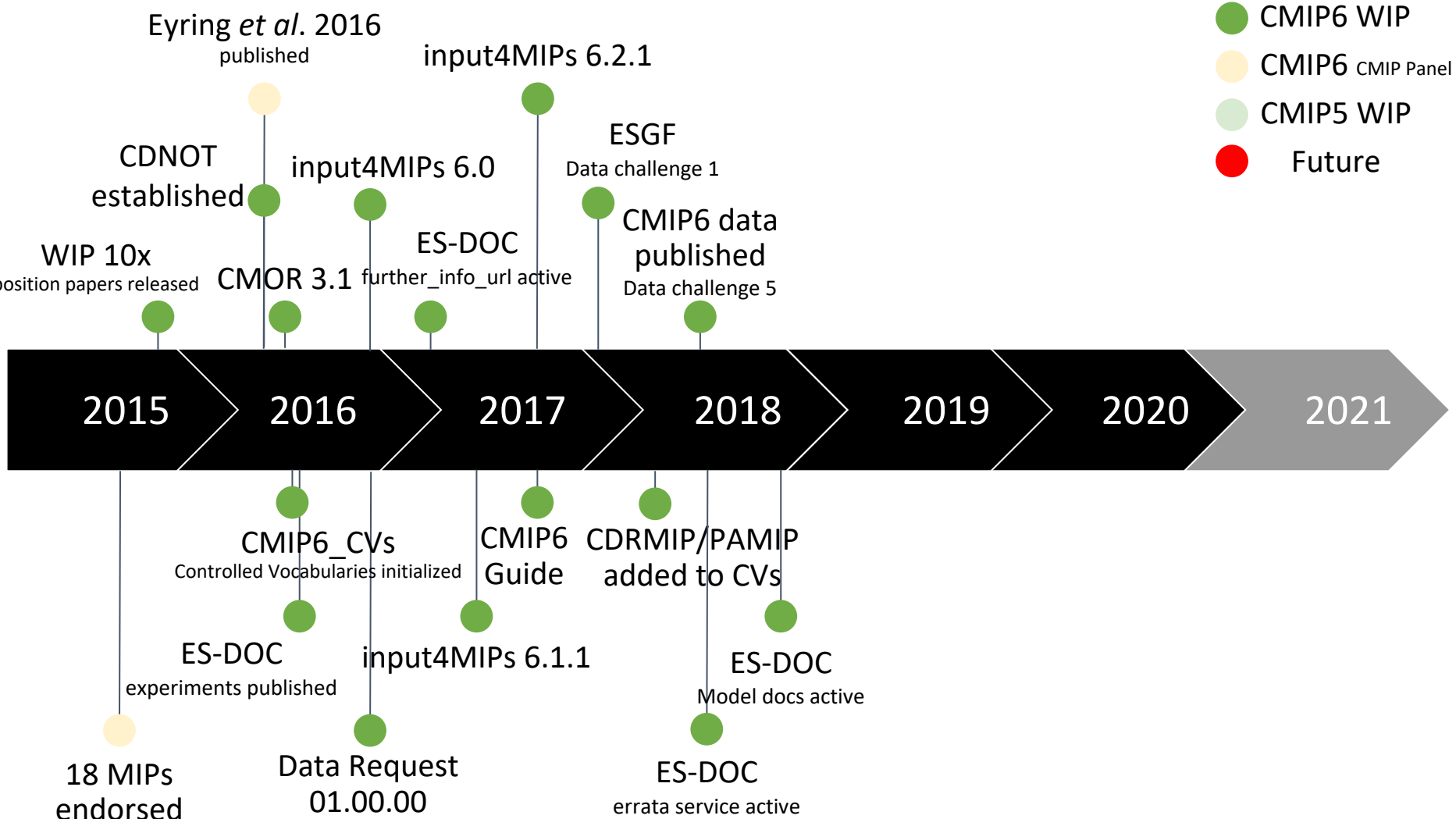




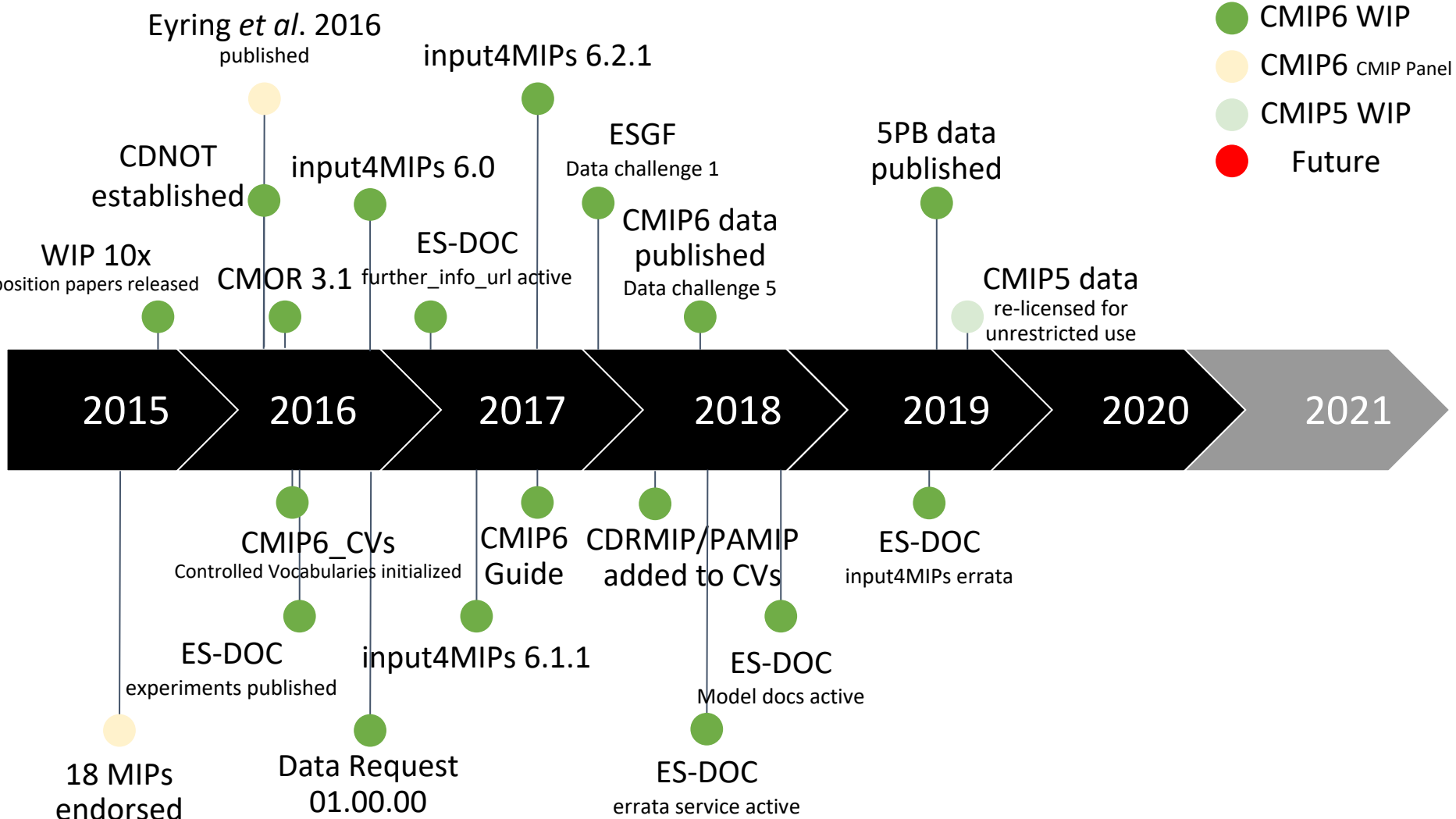
- CMIP6 WIP
- CMIP6 CMIP Panel
- CMIP5 WIP
- Future



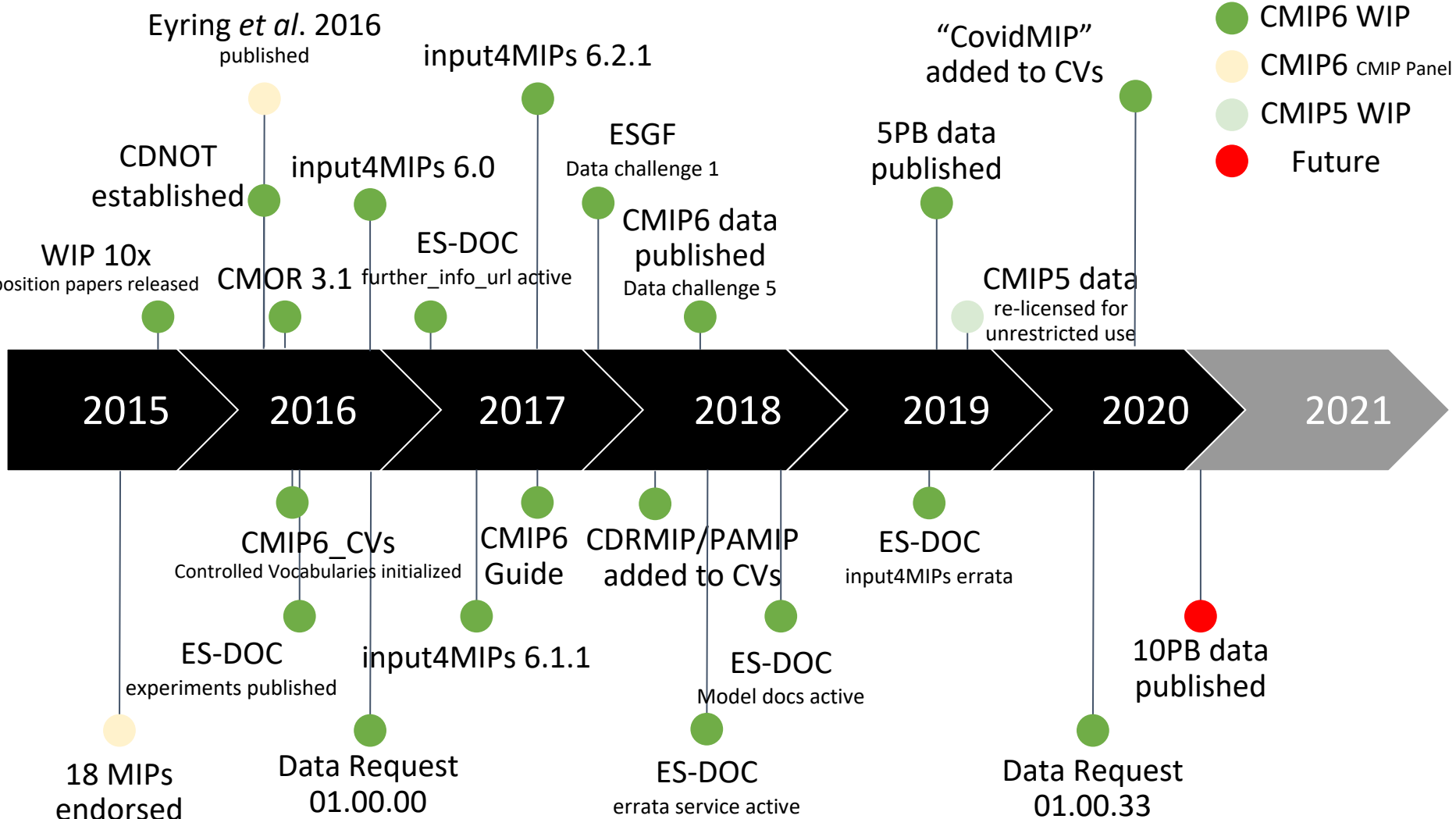
- CMIP6 WIP
- CMIP6 CMIP Panel
- CMIP5 WIP
- Future



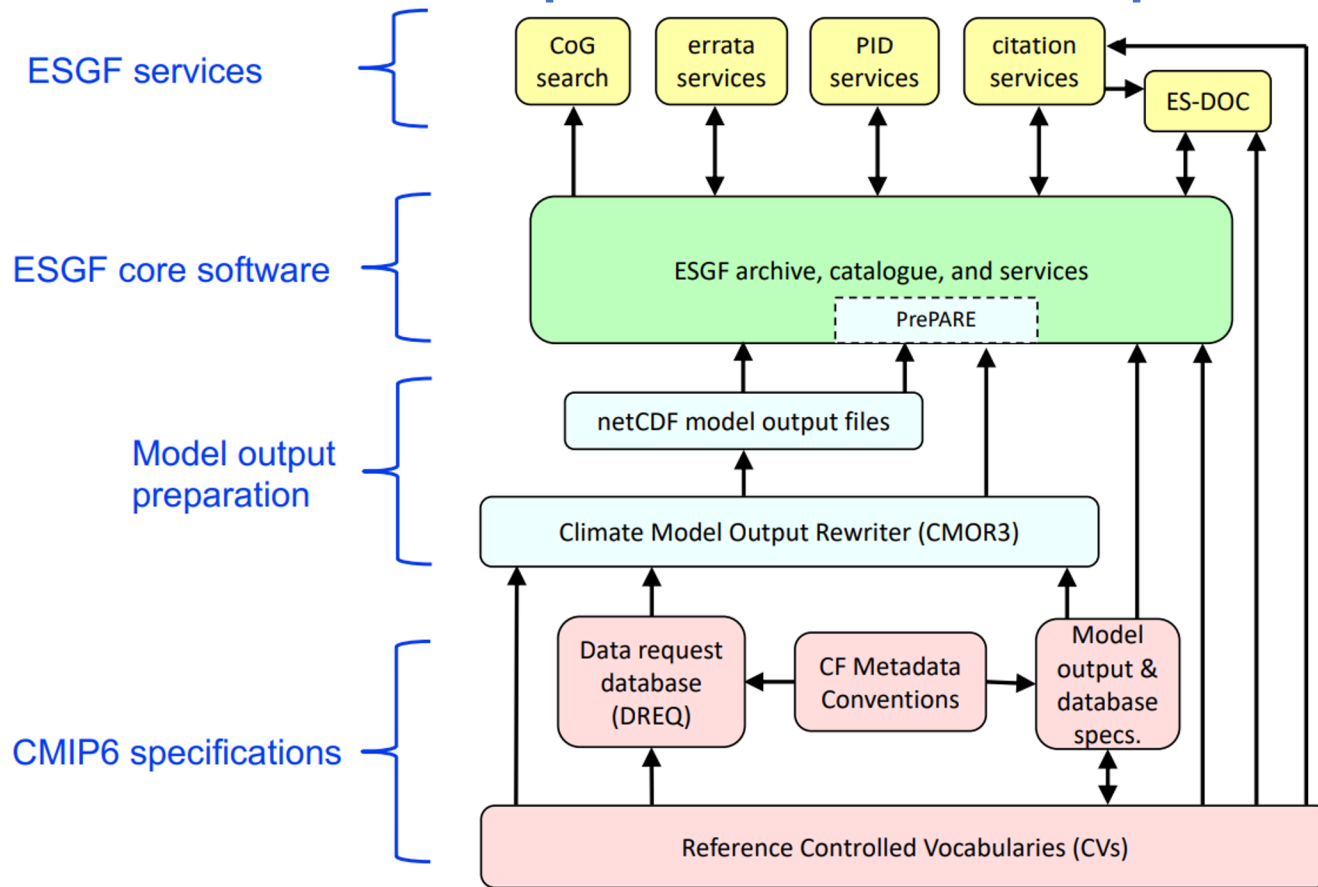
- CMIP6 WIP
- CMIP6 CMIP Panel
- CMIP5 WIP
- Future



- CMIP6 WIP
- CMIP6 CMIP Panel
- CMIP5 WIP
- Future

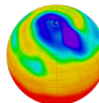


Infrastructure components and dependencies



CMIP infrastructure support

- PCMDI
 - DOE has provided 31-years of *MIP support
- ESGF
 - Originated by U.S. DOE
 - More recent major contributions from numerous others
- IS-ENES
 - European contribution to ESGF & CMIP infrastructure
- Numerous other projects and institutions, including DKRZ, IPSL, CEDA, ES-DOC, NASA, NOAA, ...
- 31+ ESGF nodes and 52 Modeling institutions around the world representing 26 countries ...



Centre for Environmental Data Analysis
SCIENCE AND TECHNOLOGY FACILITIES COUNCIL
NATURAL ENVIRONMENT RESEARCH COUNCIL



ESGF Preparation

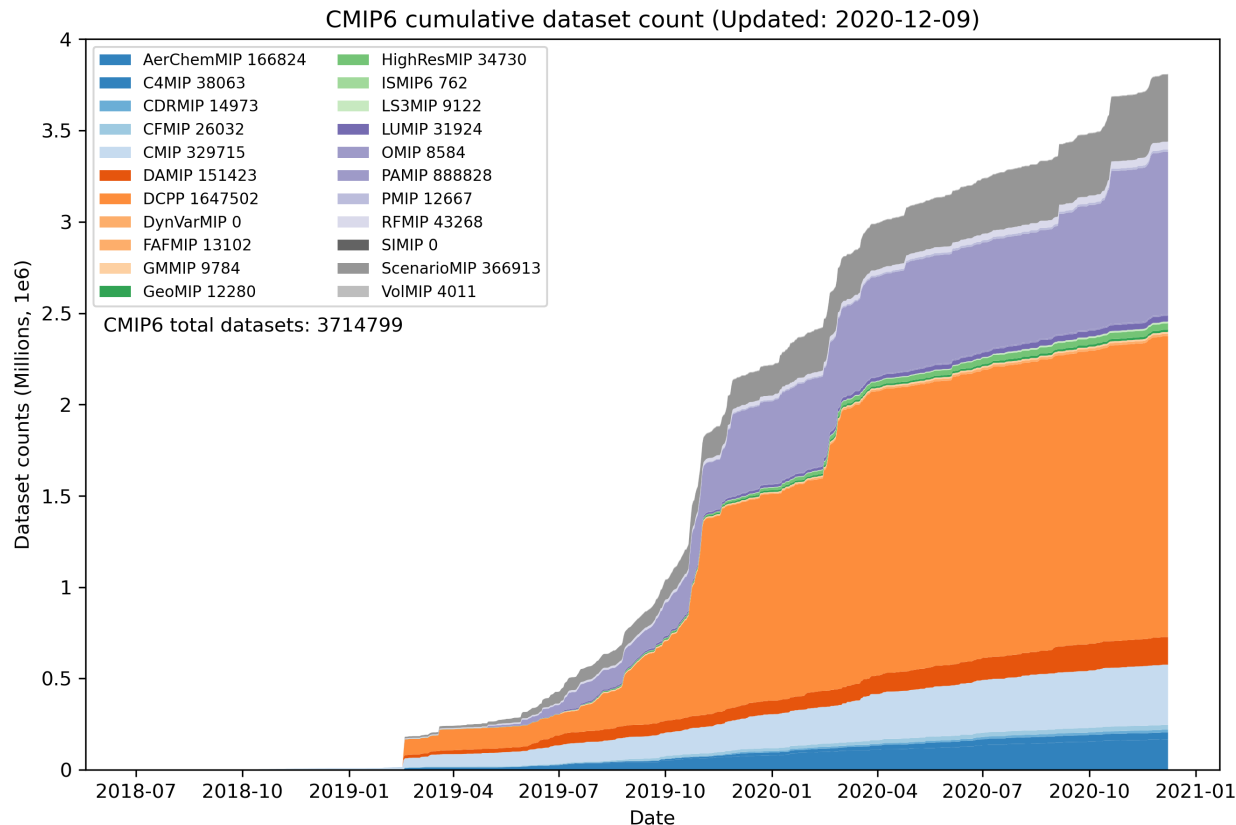
- In addition to software/hardware development/configuration
- Data challenges 1-5
 - Tested infrastructure to ensure a smooth/seamless roll out of CMIP6 for users
- Plans began at ESGF Face-to-Face (F2F) December 2017
- Began January 2018
- Finished (challenge 5) with first CMIP6 data published June 2018

Tasks 1-16 vs Challenges 1-5	1	2	3	4	5
1. Install (or update) the ESGF software stack	✓	✓	✓	✓	✓
2. Run quality control on primary data	✓	✓	✓	✓	✓
3. Publish primary data	✓	✓	✓	✓	✓
4. Publish replica data	✓	✓	✓	✓	✓
5. Verify search and download is functional	✓	✓	✓	✓	✓
6. Register data with PID assignment service		✓	✓	✓	✓
7. Verify Citation Service registers DOIs for published data		✓	✓	✓	✓
8. Populate “further_info_url” through ES-DOC scanning		✓	✓	✓	✓
9. Replicate published data			✓	✓	✓
10. Apply the “test suite”			✓	✓	✓
11. Verify the metrics collection for the dashboard			✓	✓	✓
12. Register an errata with the Errata Service			✓	✓	✓
13. Retract a version of the data				✓	✓
14. Publish a new version of the data				✓	✓
15. Ensure homogeneity across ESGF CoG sites				✓	✓
16. Move testing to production environment				✓	✓

Table 1. Table of the data challenge tasks and the challenge at which they were implemented. Note the grey tick of Task 11 in Challenge 5 indicates that the task was optional.

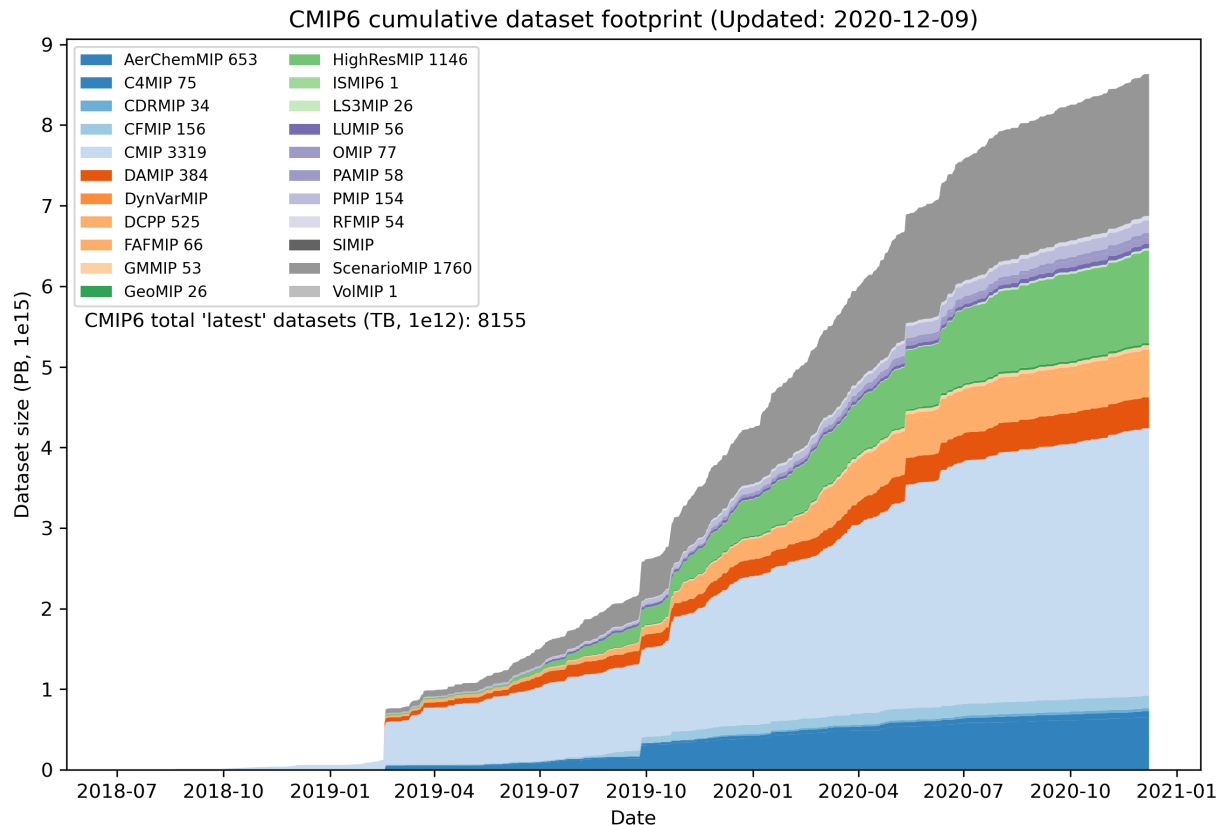
ESGF Published data

- Over 3.7 million datasets on ESGF across all CMIP6 activities/MIPs
- Delivery has been seamless – thanks to data challenges and ESGF stability testing
- Datasets - unique variable collections per experiment R1P1F



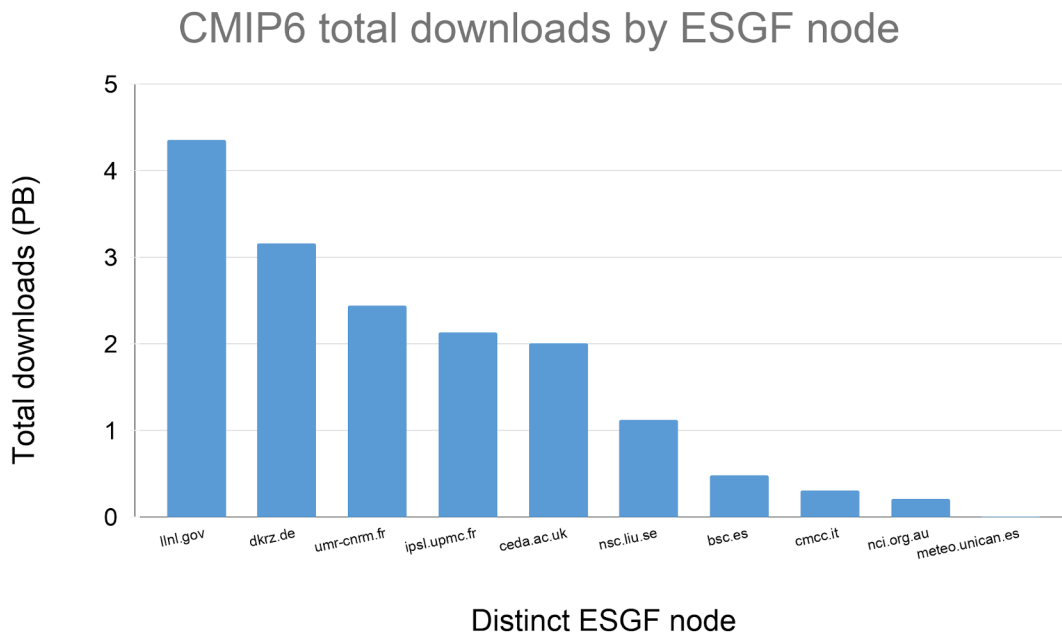
ESGF Published data

- Over 3.7 million datasets on ESGF across all CMIP6 activities/MIPs
- Delivery has been seamless – thanks to data challenges and ESGF stability testing
- Datasets - unique variable collections per experiment R1P1F
- Footprint – storage units in PBs



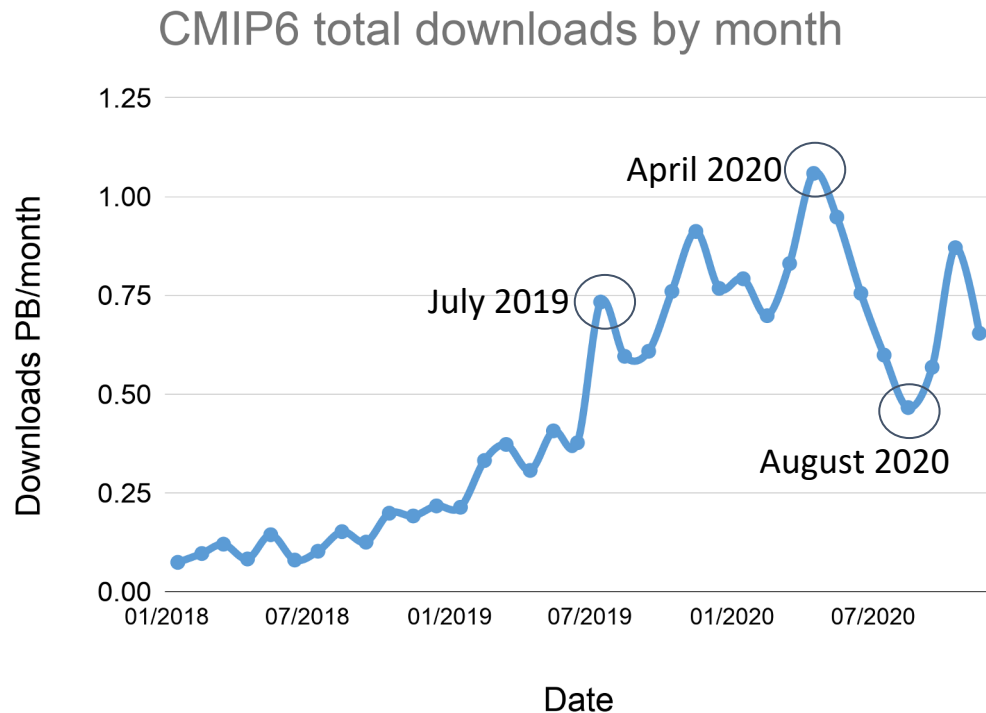
ESGF publication and replication

- 17.7 PB CMIP6 data available including ~10 PB of unique and ~8 PB of replicated data
- 16.2 PB of total CMIP6 downloads (to November 2020)
- LLNL delivered 27% of downloads to date
- DKRZ delivered 19%
- CNRM delivered 15%
- IPSL delivered 13%
- CEDA delivered 12%



Downloaded data

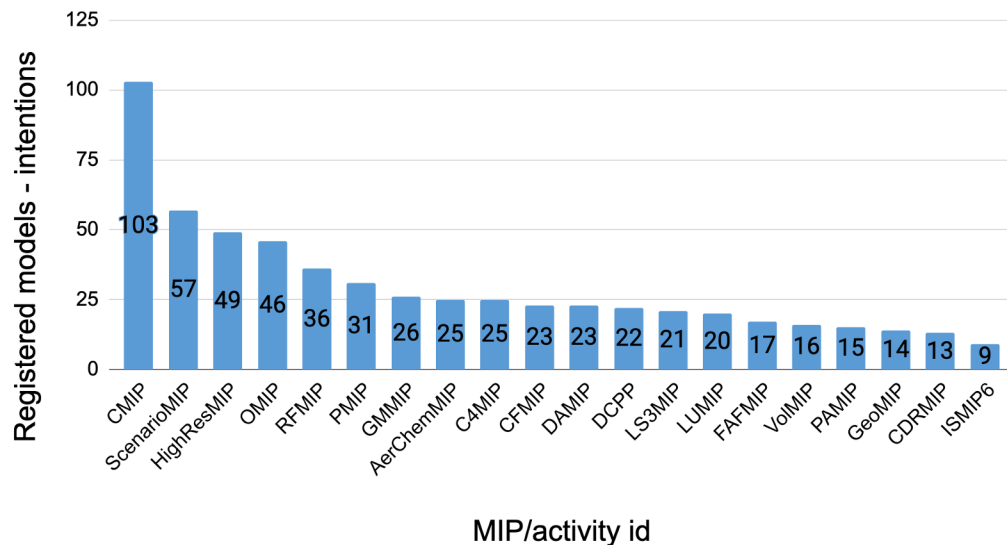
- Over 16 PB of recorded downloads (to November 2020)
- Use of secondary or local “dark” archives (e.g. Pangeo) likely means amount of data globally accessed is significantly larger
- Interesting to see COVID activity impacts



CMIP6 controlled vocabulary

- 137 models registered with CMIP6 CVs
- Each model involved in 6* activities on average
- Experiments grown from ~280 to 322 including six recent “CovidMIP” experiments added to DAMIP
- Added CDRMIP and PAMIP in March 2018

CMIP6_CVs registered models - intentions to contribute



https://wcrp-cmip.github.io/CMIP6_CVs/docs/CMIP6_experiment_id.html
https://wcrp-cmip.github.io/CMIP6_CVs/docs/CMIP6_institution_id.html
https://wcrp-cmip.github.io/CMIP6_CVs/docs/CMIP6_source_id.html

Data Request & CMOR

Data Request

- Documented in GMD
- Many releases, responsive to MIP tweaks
- Latest version 01.00.33 published to include 6 “CovidMIP” experiments sponsored by DAMIP cloned from hist-aer experiment entries

CMOR3

- Updates for python 3 compatibility
- MIP tables generated from data request and CVs when released
- Composite tables require lower investment for modeling group configuration and use

ES-DOC

- Comprehensive documentation that expands upon CMIP6_CVs
- Allows modeling groups to comprehensively document forcing and model configuration per RPPF
- Errata system invaluable for recording dataset issues and extensions

es-doc

Earth System Documentation

Documentation Search v1.0.1 [Support](#)

Project / MIP Era:

CMIP6

 Document Type:

Experiment

 Document Version:

Latest

 MIP:

*

Total Documents = 314. Filtered Documents = 314. << < Page 1 of 13 > >> 25 / page

Name	Alternative Name	Description	Version
1pctCO2	--	1 percent per year increase in CO2	1
1pctCO2-4xext	--	extension from year 140 of 1pctCO2 with 4xCO2	1
1pctCO2-bgc	--	biogeochemically-coupled version of 1 percent per year increasing CO2 experiment	1
1pctCO2-cdr	CDR-reversibility	1 percent per year decrease in CO2 from 4xCO2	1
1pctCO2-rad	--	radiatively-coupled version of 1 percent per year increasing CO2 experiment	1
1pctCO2Ndep	--	1 percent per year increasing CO2 experiment with increasing N-deposition	1
1pctCO2Ndep-bgc	--	biogeochemically-coupled version of 1 percent per year increasing CO2 experiment with increasing N-deposition	1
1pctCO2to4x-withism	--	Experiment with interactive ice sheets forced by 1 percent per year increase in CO2 to 4xCO2 (subsequently held fixed)	1
a4SST	--	control plus warming pattern SSTs	1

es-doc

Earth System Documentation

Dataset Errata - Search v0.0.0 [Support](#) [Docs](#) [Search](#) [Login](#)

Project:

CMIP6

 Experiment ID:

*

 Institution ID:

MOHC

 Source ID:

*

 Variable ID:

*

 Severity:

*

 Status:

*

Total Issues = 278. Filtered Issues = 25. << < Page 1 of 1 > >> 25 / page

#	Institute	Title	Created	Updated	Closed	Severity	Status
1	MOHC	Short datasets submitted for some UKESM1-0-LL ssp245 ...	2020-11-10	--	--	Low	New
2	MOHC	Incorrect source data used for UKESM1-0-LL ssp245 r16 ...	2020-10-28	2020-11-06	--	Critical	Resolved
3	MOHC	Further extension of some UKESM1-0-LL piControl datas ...	2020-10-02	2020-10-02	--	Low	Resolved
4	MOHC	Incorrect processing of some UKESM1-0-LL ssp126 datas ...	2020-09-25	2020-10-02	--	Critical	Resolved
5	NERC	Largely blank files	2020-07-29	--	--	High	New
6	MOHC	Grid point single time step spikes leading to excessi ...	2020-07-27	--	--	High	New
7	MOHC	Extension of a small set of UKESM1-0-LL abrupt-4xCO2 ...	2020-07-03	2020-07-03	--	Low	Resolved
8	MOHC	Extension of early UKESM1-0-LL piControl datasets	2020-06-30	2020-07-10	--	Low	Resolved
9	MOHC	Incorrect experiment id used for UKESM1 piClim-2xNOx ...	2020-06-22	2020-07-10	--	Critical	Resolved
10	MOHC	Incorrect calculation of freshwater fluxes from, and ...	2020-03-02	2020-05-19	--	Critical	On Hold
11	MOHC	Incorrect masking of northern polar row in some UKESM ...	2020-02-20	2020-10-02	--	Critical	Resolved

<https://search.es-doc.org/>
<https://errata.es-doc.org>

Data Citation Service:

- Provide data DOIs on data collection granularities
- Provide information on data usage in papers
- Disseminate data citation information

Data Citation Status:

- 2153 CMIP6 DOIs have been registered, which corresponds to a coverage of 99%.
- 1-3 DOIs are registered per day on average.

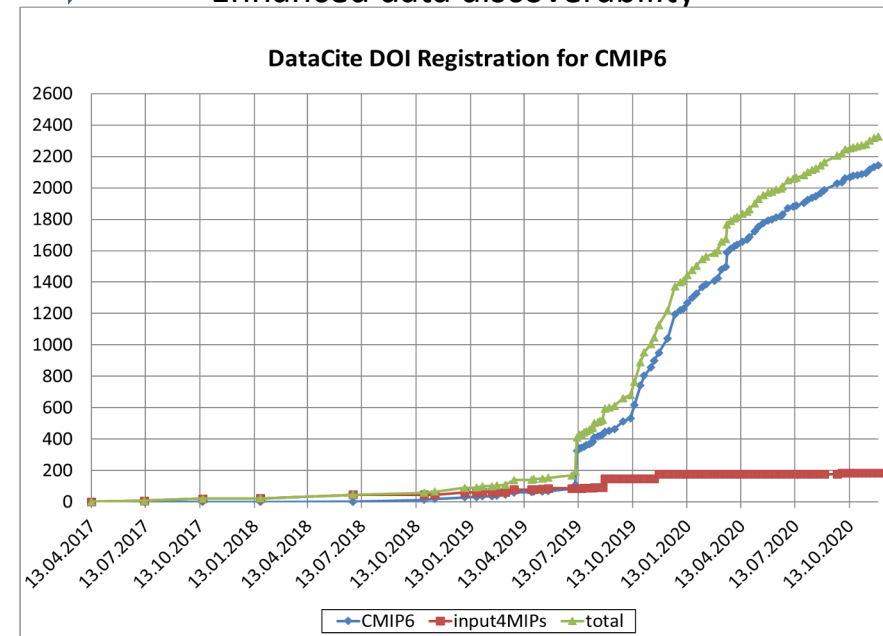
Data Citation Usage Status:

- IPCC WG1 AR6 will include CMIP6 data references.
- 62 papers referencing CMIP6 data have been added.

General Information: cmip6cite.wdc-climate.de
Available Data References: bit.ly/CMIP6_Citation_Search
Data Citation Statistics: bit.ly/CMIP6_DOI_Statistics

Benefit:

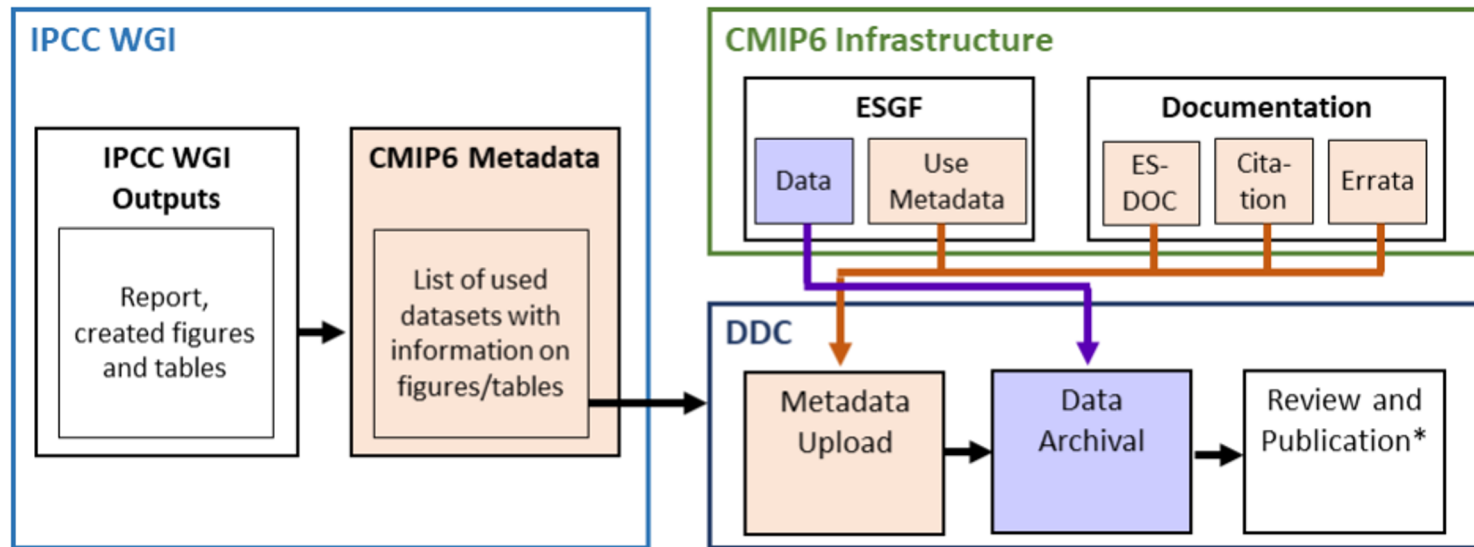
- Data is citable in scientific publications.
- Receive credit for data creation
- Enhanced data discoverability



Long-Term Data Archival in the IPCC DDC

IPCC Data Distribution Centre (DDC):

The DDC at DKRZ is the Reference Data Archive for the GCM data underpinning the IPCC Assessment Reports since the SAR. The DDC has the function to steward the data subset on the long-term complying to international data repository standards like CoreTrustSeal. The significance of data has increased in the AR6.



* DataCite publication and publication on IPCC webpages

Lessons learned

- WIP is important as the operational phase of CMIPx occurs
 - Also important in planning as infrastructure needs (storage, compute, user-base) begin to coalesce
- Funding is not operational - currently “research funding”
 - Many single points of failures
 - Dependence on single contributors with no backup
- Licensing issues have been raised by the IPCC/AR6
 - Inconsistent licensing across institutions causes problems for downstream use
 - IPCC proposes to publish AR6 and associated data (notably, the “Atlas”) under an attribution license *without the ShareAlike restriction*
- Dialogue between panels (CMIP and WIP) essential to ensure clean delivery across complex contributor network

Lessons learned (cnt'd)

- As complexity of CMIP experiments has expanded, maintaining consistency is an ongoing problem
 - E.g. historical experiments can be spawned from piControl AND past1000 experiments, but their lineage (and difference) may not be clear to a user ([CMIP6 CVs#957](#))
 - E.g. past2k and past1000 overlap same experiment duplicated across the temporal extent of the experiment ([CMIP6 CVs#979](#))
- In future phases an iterative loop of propose, review, and revise is required – likely requiring numerous iterations
- This becomes more important as the interlinkages and dependencies between MIPs/experiments continue

Paul J. Durack
durack1@llnl.gov

Matthew Mizielski
matthew.mizielski@metoffice.gov.uk

Karl E. Taylor
taylor13@llnl.gov



Work completed by the PCMDI project
is funded by the U.S. Department of
Energy, Office of Science, Office of
Biological and Environmental Research,
Regional and Global Model Analysis
Program



Disclaimer

This document was prepared as an account of work sponsored by an agency of the United States government. Neither the United States government nor Lawrence Livermore National Security, LLC, nor any of their employees makes any warranty, expressed or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States government or Lawrence Livermore National Security, LLC. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States government or Lawrence Livermore National Security, LLC, and shall not be used for advertising or product endorsement purposes.