



Rapid Evaluation Framework for the CMIP7 Assessment Fast Track

Forrest M. Hoffman¹, Birgit Hassler², Ranjini Swaminathan³, Jared Lewis⁴, Bouwe Andela⁵, Nathaniel Collier¹, Dóra Hegedűs^{6,7}, Jiwoo Lee⁸, Charlotte Pascoe⁶, Mika Pflüger⁴, Martina Stockhause⁹, Paul Ullrich⁸, Min Xu¹, Lisa Bock², Felicity Chun¹⁰, Bettina K. Gier^{11,2}, Douglas I. Kelley¹², Axel Lauer², Julien Lenhardt¹³, Manuel Schlund², Mohanan G. Sreeush¹⁴, Katja Weigel^{11,2}, Ed Blockley¹⁵, Rebecca Beadling¹⁶, Romain Beucher¹⁰, Demiso D. Dugassa¹⁷, Valerio Lembo¹⁸, Jianhua Lu¹⁹, Swen Brands²⁰, Jerry Tjiputra²¹, Elizaveta Malinina²², Brian Mederios²³, Enrico Soccimarro²⁴, Jeremy Walton¹⁵, Phil Kershaw⁶, André Lanfer Marquez²⁵, Malcolm J. Roberts¹⁵, Eleanor O'Rourke⁷, Elisabeth Dingley⁷, Briony Turner⁷, Helene Hewitt¹⁵, and John P. Dunne²⁶

¹Oak Ridge National Laboratory, Oak Ridge, Tennessee 37831-6301, United States of America

²Deutsches Zentrum für Luft- und Raumfahrt (DLR), Institut für Physik der Atmosphäre, Oberpfaffenhofen, Germany

³Department of Meteorology and National Centre for Earth Observation, University of Reading, Reading RG6 6UR, United Kingdom

⁴Climate Resource, Northcote, Victoria, Australia

⁵Netherlands eScience Center, Amsterdam, The Netherlands

⁶RAL Space, Science and Technology Facilities Council, Rutherford Appleton Laboratory, Harwell Campus, Didcot OX11 0FD, United Kingdom

⁷CMIP International Project Office, European Space Agency, Harwell Campus, Didcot OX11 0FD, United Kingdom

⁸Lawrence Livermore National Laboratory, Livermore, California 94550, United States of America

⁹German Climate Computing Center (DKRZ), Hamburg 20146, Germany

¹⁰Australia's Climate Simulator (ACCESS-NRI), Australian National University, Canberra, Australia

¹¹University of Bremen, Institute of Environmental Physics (IUP), Bremen, Germany

¹²United Kingdom Centre for Ecology & Hydrology, Wallingford, Oxfordshire OX10 8BB, United Kingdom

¹³Swedish Meteorological and Hydrological Institute, Norrköping, Sweden

¹⁴Alfred-Wegener-Institut, Helmholtz Zentrum für Polar- und Meeresforschung, Bremerhaven, Germany

¹⁵Met Office, FitzRoy Road, Exeter EX1 3PB, United Kingdom

¹⁶Department of Earth and Environmental Science, Temple University, Philadelphia, Pennsylvania 19122, United States of America

¹⁷Water Technology Institute, Faculty of Hydraulic and Water Resources Engineering, Arba Minch University, Arba Minch, Ethiopia

¹⁸Institute of Atmospheric Sciences and Climate, National Research Council of Italy, I-00133 Rome, Italy

¹⁹School of Atmospheric Sciences, Sun Yat-sen University & Southern Marine Science and Engineering Guangdong Laboratory (Zhuhai), Zhuhai, China

²⁰Instituto de Física de Cantabria (CSIC-UC), Santander, Spain

²¹NORCE Norwegian Research Centre, Bjerknes Centre for Climate Research, Bergen, Norway

²²Canadian Centre for Climate Modeling and Analysis, Environment and Climate Change Canada, Victoria, British Columbia, Canada

²³NSF National Center for Atmospheric Research, Boulder, Colorado 80307, United States of America

²⁴CMCC Foundation - Euro-Mediterranean Center on Climate Change, 73100 Lecce LE, Italy

²⁵National Institute for Space Research (INPE), São José dos Campos, Brazil

²⁶NOAA Geophysical Fluid Dynamics Laboratory, Princeton, New Jersey 08540, United States of America

Correspondence: Forrest M. Hoffman (hoffmanfm@ornl.gov), Birgit Hassler (Birgit.Hassler@dlr.de), and Ranjini Swaminathan (R.Swaminathan@reading.ac.uk)



Abstract. As Earth system models (ESMs) grow in complexity and in volumes of output data, there is an increasing need for rapid, comprehensive evaluation of their scientific performance. The upcoming Assessment Fast Track for the Seventh Phase of the Coupled Model Intercomparison Project (CMIP7) will require expeditious response for model analyses designed to inform and drive integrated Earth system assessments. To meet this challenge, the Rapid Evaluation Framework (REF), a community-driven platform for benchmarking and performance assessment of ESMs, was designed and developed. The initial implementation of the REF, constructed to meet the near-term needs of the CMIP7 Assessment Fast Track, builds upon community evaluation and benchmarking tools. The REF runs within a containerized workflow for portability and reproducibility and is aimed at generating and organizing diagnostics covering a variety of model variables. The REF leverages best-available observational datasets to provide assessments of model fidelity across a collection of diagnostics. All diagnostics were identified and finally selected with community involvement and consultation. Operational integration with the Earth System Grid Federation (ESGF) will permit automated execution of the REF for specific diagnostics as soon as model data are published on ESGF by the originating modelling centres. The REF is designed to be portable across a range of current computational platforms to facilitate use by modelling centres for assessing the evolution of model versions or gauging the relative performance of CMIP simulations before being published on ESGF. When integrated into production simulation workflows, results from the REF provide immediate quantitative feedback that allows model developers and scientists to quickly identify model biases and performance issues. After the REF is released to the community, its subsequent development and support will be prioritized by an international consortium of scientists and engineers, enabling a broader impact across Earth science disciplines. For instance, the REF will facilitate improvements to models and reductions in uncertainties for projections since ESMs are the main tool for studying the global Earth system. Production of reproducible diagnostics and community-based assessments are the key features of the REF that help to inform mitigation and adaptation policies.

Copyright statement. This manuscript has been authored by UT-Battelle, LLC, under contract DE-AC05-00OR22725 with the US Department of Energy (DOE). The US government retains and the publisher, by accepting the article for publication, acknowledges that the US government retains a nonexclusive, paid-up, irrevocable, worldwide license to publish or reproduce the published form of this manuscript, or allow others to do so, for US government purposes. DOE will provide public access to these results of federally sponsored research in accordance with the DOE Public Access Plan (<http://energy.gov/downloads/doe-public-access-plan>).

1 Introduction

Earth system models (ESMs) are the primary tools for the science community to study interactions between the atmosphere, land, ocean, cryosphere, and biosphere in the Earth system and how it responds to human-induced and natural forcings (e.g., Mauritsen et al., 2019; Séférian et al., 2019; Yukimoto et al., 2019; Boucher et al., 2020; Danabasoglu et al., 2020; Senior et al., 2020; Döscher et al., 2022). Currently, the next generation of ESMs is being finalized and new simulations will be generated for the Seventh Phase of the Coupled Model Intercomparison Project (CMIP7) and the precursory CMIP7 Assessment Fast Track. As the number of models, ensemble sizes, complexity and output requirements continue to grow, there is an urgent



need to objectively evaluate the fidelity of these models and exploit this wealth of information in order to efficiently advance our understanding of the Earth system and to inform climate mitigation and adaptation policies. This includes, specifically, identification of model uncertainties or systematic biases that may prevent us from objectively constraining model-derived projections of future climate change. To address this need, it is critical to develop efficient evaluation methods that make use of the growing archive of output from these simulations and reduce the time to interpret the output as meaningful scientific insights that can be used by stakeholders and policy makers. The rapid growth of ESM data, driven by model complexity and computational advances, creates both opportunities and challenges. Effective, reproducible, accurate, and unbiased data processing is crucial for translating model outputs into actionable insights for climate policy (IPCC, 2023).

The CMIP Model Benchmarking Task Team (MBTT) was created to address this challenge in preparation for CMIP7 and the Intergovernmental Panel on Climate Change (IPCC) Seventh Assessment Report (AR7) and its deadline for contributions. Following the first phase of the MBTT, which reviewed ESM evaluation and benchmarking approaches and identified the collection and collation of existing community benchmarking software packages (Hassler et al., 2025), and outlining best practices for the use of observational datasets for model evaluation (Beadling et al., in preparation); the MBTT has now expanded its efforts toward developing a community-designed Rapid Evaluation Framework (REF) for routine and rapid benchmarking of CMIP simulations. The conceptual design of the REF was developed at a MBTT workshop in May 2024 and approved by the CMIP Panel in July 2024, with development work commencing in October 2024. The initial design was strongly motivated by the ideas and vision developed for CMIP6 (Eyring et al., 2019); the goal of the REF for CMIP7 is to deliver a complete end-to-end system that will provide a systematic and rapid performance assessment of CMIP models, initially targeting the model experiments contributing to the CMIP7 Assessment Fast Track, which will support IPCC AR7 (Dunne et al., 2024). The vision of the REF is to be a community-owned evaluation framework, leveraging existing community-built model evaluation packages and incorporating an application programming interface (API) that will execute modules for generation of diagnostics and the metrics that underlie them.

Rather than directly ranking models, the REF is primarily concerned with providing objective measures of model performance to allow the wider community to make informed decisions about models that are most appropriate for their specific needs. This requires a standard set of diagnostics and performance metrics to facilitate the comparison of key variables simulated by models with standardized observational reference datasets, and assessment of whether fundamental processes in the Earth system are adequately represented in the models. Once expanded for community use beyond the initial Assessment Fast Track version, the REF will have a wider array of applications and users, including other modelling communities and scientific domains, as well as organisations utilising CMIP models for conducting feedback analysis, impacts assessment, or financial planning. To facilitate understanding of the descriptions of the REF, key terms used throughout the manuscript are defined here. This terminology may not be used consistently across Earth science disciplines.

– **Reference Datasets** – A reference dataset is a collection of observationally constrained or model data used as a standard within a model evaluation diagnostic. Examples may include *in situ* measurements, extrapolated data (from statistical or AI/ML methods), remote sensing data, reanalysis data, or any other dataset that is meant to represent a best estimate of a geophysical quantity or a physical, chemical, biological, or ecological state or process.



- **Model Variables** – A model variable is any quantitative representation or characterization of a physical, chemical, biological, or ecological state or process that changes during execution of the model. Variables are used to represent mass, energy, velocity, momentum, flux rates, and other parameters within models. Model variables may or may not represent observable quantities, and they may be inferred, estimated, or calculated from other related variables or observables.
- **Diagnostics** – A diagnostic is a comparison of a model variable or some combination of model variables with a reference dataset or an intercomparison across models of a model variable or some combination of model variables. A diagnostic may also represent an evaluation of a relationship between multiple model variables and/or multiple reference datasets (i.e., Relationship Diagnostics). Diagnostics have sometimes been called “confrontations” since the objective is to confront models with best-available observations or with best-available model or model ensemble outputs. A diagnostic consists of one or more model performance metrics.
- **Metrics** – A metric is a single statistical evaluation contained within a diagnostic. A diagnostic may consist of more than one metric. Examples include bias, root mean squared error (RMSE), spatial or temporal correlations (Taylor, 2001), Earth Mover’s Distance, Hellinger Distance (Hellinger, 1909), phase/timing of the seasonal cycle, amplitude of the seasonal cycle, inter-annual variability (Giorgi and Francisco, 2000). Not all metrics are useful for all variables or should be used with every observationally constrained dataset. Each metric may be evaluated to produce a metric scalar.
- **Metric Scalars** – A metric scalar is the numerical output resulting from the calculation of a performance metric (e.g., the calculated bias).
- **Scores** – A score is a scalar value (0.0–1.0) transformed from a metric scalar or produced by aggregating multiple metric scalars or multiple scores together.
- **Verification** – Verification is the process of assessing model consistency in terms of correct implementation of the represented processes as expected from the model and simulation experimental design. Sometimes, this is accomplished as the model simulations are being produced (such as monitoring the conservation of total energy, total atmospheric mass, etc.), and the focus is often on the artifacts introduced by the numerical discretization scheme (e.g., Lauritzen et al., 2022) or by changes to software or hardware used for the simulations (e.g., the Ensemble Consistency Test in Baker et al. (2015) and the Time Step Consistency Test in Wan et al. (2017)).
- **Validation** – Validation is the process of determining the degree to which a model accurately represents processes in the real world, particularly for the intended uses of the model. Validation can include a broad range of aspects from ensuring correct units and the sign of the data produced, to the interactions between model components or variables and process representations.
- **Fidelity** – The fidelity of a model is a quantitative assessment of the degree to which model output corresponds to the reference data in aggregate, resulting from a validation exercise. One approach for deriving a fidelity metric is to aggregate relevant scores.



100 2 Conceptual Design of the REF

2.1 Overview of the REF

The REF was designed to be a community-owned evaluation framework that leverages existing open source, community-built model evaluation and benchmarking packages that are integrated together through a standard application programming interface (API) that will execute modules for generation of diagnostics and the metrics that underlie them. By incorporating
105 existing tools and metrics with publicly available reference datasets, the REF standardizes a community workflow for CMIP activities, reduces duplication of efforts for evaluating models, and, through deployment on the Earth System Grid Federation (ESGF), provides rapid feedback to the research community about relative model performance across a wide range of model components and variables. Moreover, the REF is expected to offer a starting point for the research community to expand model evaluation and benchmarking capabilities and applications within their own institutions or modelling centres.

110 The high-level workflow for the REF as it runs on ESGF nodes, as shown in Fig. 1, is triggered by the publication of new simulations requiring evaluation. After the model data have undergone a quality control check to assure the metadata are correct, data storage and compute resources are allocated for a new execution of the diagnostics generation process. As agreed upon by the CMIP Panel, model data that do not conform to the required Controlled Vocabulary (CV) and metadata standards will not be evaluated by the REF. Next, an optimized directed acyclic graph of tasks is produced, processes are initiated to
115 calculate model evaluation metrics and construct diagnostics. The outputs of diagnostics are then staged on public websites for sharing with the research community. The initial implementation of the REF was created to evaluate Assessment Fast Track simulations (Dunne et al., 2024), using five to six diagnostics across five Earth system realms with the expectation that additional diagnostics would be added in the future.

Modelling centres usually evaluate their models during development with the overall performance documented and published
120 once the models are finalised and key simulations are completed. However, this approach has three main limitations. First, the process is slow, and evaluation results often come too late to inform IPCC assessment reports or even later for stakeholders who need timely information. Second, these evaluation results are hard to access; they are scattered across different types of papers (technical or peer-reviewed, open access or not, etc.) across the modelling centres and often only partially included in final IPCC reports, making them difficult to locate and piece together. Third, no consistent approach is applied across centres;
125 each group runs its own evaluations using sometimes different methods, which makes comparisons between models difficult. The REF aims to address these issues by bringing together core evaluation outputs in one place. It helps IPCC experts quickly assess, understand, and synthesise the performance of a new generation of ESMs for Assessment Fast Track simulations. The REF also makes it easier for stakeholders to access the information they need to support regional analysis for adaptation and mitigation, and it supports best practices for evaluation during model development. Releasing the REF before model outputs
130 are submitted to ESGF offers modelling centres the opportunity to systematically assess their models and to make targeted improvements during the development phase.

The REF aims to integrate best-available reference datasets for comparison with model outputs. These reference datasets were collected in obs4MIPs, a project created to distribute data that support evaluation of ESMs via ESGF (nodes listed at



Rapid Evaluation Framework Overview

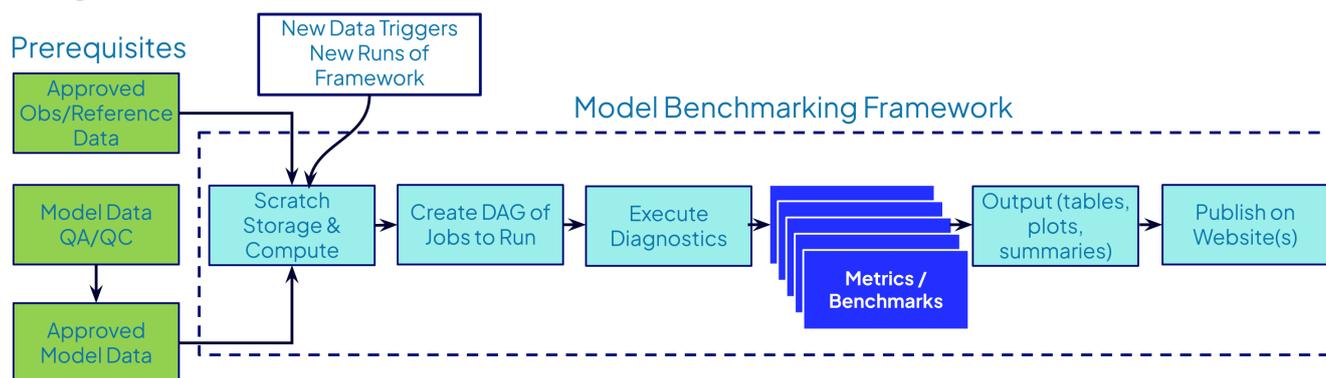


Figure 1. The high-level workflow for the Rapid Evaluation Framework (REF) run on Earth System Grid Federation (ESGF) nodes, shown here, combines quality-assured model output with a collection of observational reference data to initiate execution of relevant diagnostics and generation of tabular and graphical representations of a variety of metrics. The results are then published on websites for community use. Figure available at <https://doi.org/10.5281/zenodo.15594502> (Hoffman et al., 2024).

135 <https://esgf.github.io/nodes.html>). This project compiles a range of observationally constrained data, formatted according to
 140 output.

145 While obs4MIPs serves an important role in distributing datasets the community agrees are essential for model evaluation, the REF is flexible and designed to use additional reference data that meet the required format. This adaptability allows for the incorporation of new reference datasets as they become available, or when specific needs arise, ensuring that the framework can accommodate a broad range of observational data. The REF is not intended to reformat reference or model data, and REF users must ensure the data they wish to use follow the CF conventions and meet the CMIP metadata standards. Thus, data validator software will be used for quality control (QC) of CMIP model output published on ESGF and of reference datasets published in obs4MIPs. This QC is required to assure the data used for benchmarking are consistent and reliable, which enables meaningful model-data comparisons (Hawkins and Sutton, 2016).

2.2 Community Engagement

150 Key stakeholders for the co-development of the REF were identified as modelling centres involved in the Assessment Fast Track, reference dataset providers, tool and package developers, diagnostic developers and climate scientists, including IPCC



155 authors, seeking to analyse the Assessment Fast Track outputs. The MBTT initiated work on the REF based on the results of the CMIP6 Community Survey (O'Rourke, 2025), conducted by the World Climate Research Programme (WCRP) from January to March 2022. Detailed analysis of the survey results by the Task Team (Lee, 2024) revealed calls for open and transparent model benchmarking and evaluation and the ability to run CMIP evaluation tools alongside ESGF, as part of the data publication process, with provision of evaluation results display.

160 Members of Fresh Eyes on CMIP worked with the CMIP-International Program Office and the MBTT to develop a survey that was circulated to modelling centre science and technical leads, as well as more widely in the Fresh Eyes on CMIP and CMIP community mailing lists and via social media channels. The survey was open for the period 7th through 17th May 2024 and generated 152 unique responses, and the responses were analysed by members of the Fresh Eyes on CMIP and reported at a MBTT workshop, which informed design and development of the initial REF proposal. Suggestions included considering multiple software tools, strategies for ensuring accessibility, tools for quality control, and transparency and provenance for reference datasets and outputs. Additional suggestions included acceptance of data not processed by the Computer Model Output Rewriter (CMOR; Taylor et al., 2004), potentially with AI-assisted reformatting, improved accessibility through use of ESGF and cloud-based computing to enhance data inclusivity and flexibility, as well as errata tracking. A detailed analysis of the survey results is available in the survey report by (Wang et al., 2025).

170 A separate community survey was subsequently disseminated to modelling centres in June 2024 that included requests for input on evaluation tools used, interest in use of a benchmarking framework within their own computing environments, and willingness to submit preliminary outputs for quality assurance and control checks (O'Rourke, 2024b). This survey informed refinement of the framework structure and implementation plan. In September 2024, a survey was disseminated to modellers, modelling centres, reference dataset and data infrastructure providers requesting input regarding diagnostics that should be prioritized for inclusion in the REF, receiving 53 responses. A co-development session with members of the Earth observation community was held at the European Space Agency (ESA) Climate Change Initiative (CCI) and Climate Modelling User Group (CMUG) integration meeting in October 2024, which resulted in the inclusion of ozone-related metrics in existing diagnostics where appropriate. The outcomes of the survey, in combination with suggestions from the MBTT, resulted in a table of diagnostics, shown in Section 4 and published separately at CMIP Model Benchmarking Task Team (2024). The REF was thereby launched on 4th November 2024 at a community-engagement event that saw wide participation. Successive months were dedicated to implementation of the diagnostics from four existing evaluation and benchmarking packages within the REF architecture. The packages selected for providing diagnostics for the Assessment Fast Track REF were ESMValTool, PMP, ILAMB, and IOMB, which are described below in Section 4. During this period, engagement efforts focused on technical implementation, involving members of the community responsible for provision of supporting infrastructure and quality control, including the WGCM Infrastructure Panel (WIP), obs4MIPs Steering Panel, the ESGF-WIP Quality Assurance and Control Working Group, the CMIP Data Request and CVs Task Teams, as well as the IPCC Task Group on Data Support for Climate Change Assessments (TG-Data), on their citation and provenance requirements for the REF, see Section 4.6.

185 A preliminary prototype of the working suite was stress-tested for operational usage in time for the REF Hackathon, which was held at the Met Office, in Exeter, UK, from 10th to 13th March 2025 and included dedicated drop-in sessions for modellers



and reference dataset providers. This prototype was able to ingest a sample of CMIP6 model output and obs4MIPs data, running a small subset of ILAMB, PMP, and ESMValTool diagnostics. In May 2025, modelling centre science and technical personnel, as well as previous attendees of the REF launch, were invited to attend the launch of the beta version of the REF. Approximately 107 participants registered for the demo and to attend follow-up feedback sessions in June 2025 regarding the beta release. This beta testing period will be instrumental for building an easy-to-use tutorial and documentation, solving potential issues with REF usage on HPC machines, and receiving feedback from users, with the aim of delivering a public release of the REF for actual usage with Assessment Fast Track outputs in October 2025.

3 Target Applications

Considering the REF scope, stakeholder consultation and initial implementation plan for the Assessment Fast Track simulations, three primary potential applications for the REF were identified. The primary application of the REF is to provide stakeholders from the data analysis and impacts assessment communities with timely information about the scientific performance of ESMs with respect to observationally constrained (reference) datasets across all Earth science realms. Reference data are available primarily but not exclusively from the contemporary observational period, approximately over the last 50 years. Secondly, the REF produces diagnostics focused on key sensitivity indicators that typically do not have corresponding observational constraints, such as equilibrium climate sensitivity (ECS) and transient climate response (TCR). Access to all of this diagnostic information assists stakeholders in selecting simulation results for further analysis, downscaling studies, impacts analysis, or other research. Moreover, the results of the REF can increase equality in climate data access for community members who lack adequate computer resources or Internet access. In some cases, diagnostics from the REF include graphs, charts, or figures directly usable in research publications or assessment reports, offering analysts more time to focus on specific research questions that may use information from the REF.

A key design goal of the REF is that it be usable by modelling centres, research institutions, and individual scientists to enable validation of ESM output prior to publication of simulations on ESGF, intercomparison of model results, and general purpose analysis. Running the REF prior to data submission provides the opportunity for data providers to gauge the performance and sensitivity of their simulations in a standard fashion at any time. Modelling centres typically use their own collections of diagnostics, often developed in-house, for routine model assessment, while other centres or institutions have adopted community-developed evaluation tools or employ a combination of community tools and in-house diagnostics. The REF offers a general purpose framework for model evaluation and is equally useful for tracking the scientific performance of different versions of the same model. The REF diagnostics also provide a convenient means for determining if model changes during development yield improvements. Thus, modelling centres may find that running the REF as a part of their workflow for repeated simulations provides a practical way to track evolving performance of their model. It may also be necessary to reformat the relevant model output variables for each simulation, making them CF-compliant and ensuring CMIP naming standards and metadata are provided, within the workflow for the REF to be able to run correctly. Moreover, analysts from different



research communities may want to add functionality to the initial REF implementation to enable use of other observational
220 datasets, additional diagnostics, and other metrics as discussed in Section 5.4.1.

Furthermore, the REF could serve as an early warning system for modelling centres to identify inconsistencies in the model
variables in a more complex way than was done for the prior CMIP activities (e.g., Taylor et al., 2004) and thus reduce ESGF
data traffic and storage of erroneous data, limiting data use by the wider community. The REF can also be used for a variety
of model intercomparison activities outside of the scope of CMIP, and project leaders can ask simulation contributors to run
225 the REF and provide standard diagnostic results instead of sharing large volumes of model output. Additionally, the REF
enables individual researchers interested in Earth system science to explore model output and apply best-available reference
data to better understand model capabilities and gaps in process representation. The REF provides a standard framework for
integrating additional diagnostics for use by research institutions or local analysts and scientists. New diagnostics integrated
into the REF can be easily shared among modelling centres and researchers, and they become key candidate additions to future
230 public releases of the community version of the REF.

4 System Description of the CMIP7 Assessment Fast Track REF

4.1 Included Evaluation and Benchmarking Packages and the Coupling Strategy

The open source evaluation and benchmarking packages described below—ESMValTool, PMP, and ILAMB & IOMB—were
chosen for inclusion in the first version of the REF. Subsets of diagnostics from each package were selected based on input
235 from the MBTT and the community. For the initial version of the REF, diagnostics were restricted to analyse only monthly
mean output from the expected new simulations. In this section, the packages are briefly described, although a more in-depth
description of them can be found in Hassler et al. (2025), and the CMEC standards are also described, since they offer a strategy
for integrating these disparate evaluation packages.

4.1.1 ESMValTool

240 The Earth System Model Evaluation Tool (ESMValTool) is an open source community-developed software package aimed at
performing many different diagnostics and metrics for the evaluation and benchmarking of ESMs (Righi et al., 2020; Eyring
et al., 2020; Lauer et al., 2020; Weigel et al., 2021; Schlund et al., 2023; Lauer et al., 2025). Many of the diagnostics and
metrics that have been officially released in ESMValTool have then been systematically used for the production of multi-model
intercomparisons embedded in several chapters of Working Group I (WGI) of the IPCC Sixth Assessment Report (AR6) (IPCC,
245 2023). ESMValTool strongly advocates traceability and reproducibility; therefore, all diagnostic results are provided with
metadata documenting the provenance of the model outputs and reference data, the used software packages, and the calculated
metrics and diagnostics. A part of the software package also deals with the adjustment of model or observational datasets that
are not strictly compliant with the CF metadata conventions. While the core capabilities of ESMValTool are fully Python-
based, diagnostics can be based on other open source languages like NCL or R. For all contributions, ESMValTool implements



250 rigorous technical and scientific reviews before new code can be included in the official release, requiring Python Enhancement Proposal (PEP) 8 standards, and testing with pre-commit and Codacy for maturity of the code and standardization.

4.1.2 PMP

The PCMDI Metrics Package (PMP) is an open source Python software package developed for objective and rapid assessments and benchmarking of ESMs, with a focus on atmospheric variables, against the most up-to-date observational datasets (Lee et al., 2025). It has been playing an important role in the systematic evaluation of thousands of simulations from CMIPs, with a strong emphasis on physical climate metrics, particularly atmospheric means and variability. Among its diverse suite of metrics, a subset of metrics that are calculated based on the monthly time series of model variables were chosen for the first implementation in the REF (Table 1). The subset of metrics includes the annual cycle (Gleckler et al., 2008), El Niño Southern Oscillation (ENSO) CLIVAR (Climate and Ocean: Variability, Predictability and Change) metrics (Planton et al., 2021), extra-
260 tropical modes of variability (Lee et al., 2021), and the monsoon (Wang et al., 2011). By offering a database of pre-computed statistics for CMIP6 models, the PMP streamlines the comparison process, making it easier for modelling centres to evaluate their results against established benchmarks.

4.1.3 ILAMB and IOMB

The International Land Model Benchmarking (ILAMB) and International Ocean Model Benchmarking (IOMB) are open
265 source Python software packages that share a large portion of the same codebase. They were developed to provide systematic assessment of land and ocean model performance, primarily for terrestrial and marine biogeochemistry, through comparison with reference datasets (Collier et al., 2018; Luo and Hoffman, 2022; Fu et al., 2022). Diagnostics and metrics within ILAMB and IOMB were developed with engagement of land and ocean modellers and with the *in situ* and remote sensing observational communities (Luo et al., 2012; Hoffman et al., 2017). Both packages were used to evaluate and intercompare historical
270 simulations from CMIP5 and CMIP6 (IPCC, 2023, Chapter 5, Figure 5.7), as well as serving important roles in informing the development of land and ocean models for DOE's Energy Exascale Earth System Model (E3SM; Burrows et al., 2020; Zhu et al., 2019; Yang et al., 2019) and the Community Earth System Model (Lawrence et al., 2019). ILAMB and IOMB both offer a variety of statistical metrics, including bias, RMSE, timing/phase of the seasonal cycle, spatial correlation, and interannual variability. Scores from these metrics are aggregated to provide high-level scores for each model-dataset pairing.
275 Functional relationship metrics within ILAMB and IOMB are used to evaluate the degree to which model variable-to-variable relationships correspond to those of observational data. ILAMB and IOMB produce hierarchical webpages designed to offer users and analysts the ability to view many metrics at once for a given model-dataset pair, as well as to intercompare graphical representations of metrics across all models at once.



4.1.4 CMEC

280 Coordinated Model Evaluation Capabilities (CMEC) is an effort to bring together a diverse set of analysis packages that
have been developed to facilitate the systematic evaluation of ESMs (Ordonez et al., 2025). CMEC provides the strategy for
coupling multiple community benchmarking packages in the REF. CMEC includes capabilities supported by multiple agencies,
and capabilities that have been contributed by community-based experts and international agencies. With widespread and rapid
285 growth in the number of available diagnostic and model evaluation tools, a lack of standards within the evaluation community
have meant that running even a single evaluation tool can require extensive user intervention. Given the significant commonality
in how these evaluation tools operate, interoperability is a natural goal achievable through robust and light-weight standards.
The three goals of the CMEC project include: (1) to develop robust and light-weight standards for operation of evaluation
packages and their output; (2) to develop accompanying tools for installation of evaluation packages, coordinated execution
of evaluation packages, and obtaining data products necessary for operation of these tools; and (3) to build connections across
290 groups, research centres, and individual investigators performing model evaluation. The CMEC standards were adopted for
integrating the output of diagnostics produced by the model benchmarking packages described above.

4.2 Diagnostics

At the heart of the REF are the diagnostics that were selected in an iterative process (see Section 2.2) and that can be calculated
with each new simulation that is presented to the REF. The diagnostics are grouped according to the five different realms
295 (Ocean & Sea Ice, Land & Land Ice, Atmosphere, Earth System, and Impacts & Adaptation), and each diagnostic included in
this first version of the REF is calculated by only one software package. Table 1 provides an overview of all the diagnostics that
were selected for the first version of the REF, based on the diagnostics table published at <https://zenodo.org/records/14284375>,
their realm, the software package with which the diagnostic is calculated, and the reference datasets used in the comparison.
A more detailed table of the diagnostics is presented in Appendix B. For some of the diagnostics, different methodologies
300 adopting different software packages across the diagnostic providers were discovered (e.g., double Inter-tropical Convergence
Zone (ITCZ) biases); in this case, the principle of minimal computational resources and least number of required variables was
adopted in order to choose the appropriate tool to be responsible for the diagnostic in the REF.

During the REF development phase, it became clear that two of the identified diagnostics would not be available from
the software packages implemented in the Assessment Fast Track REF. These diagnostics, with their unique IDs 5.1 (High
305 amplitude Rossby waves) and 5.2 (Internal variability or ensemble spread for individual models), were then removed from the
list of those to be integrated in the initial version of the REF, following consultation with, and agreement from, the Impacts
& Adaptation CMIP Data Request Author team leads and co-chairs of the Vulnerability, Impacts, Adaptation and Climate
Services (VIACS) Advisory Board.

A more detailed description of each diagnostic and the rationale and methodology for its implementation is presented in
310 Appendix C.



Table 1. Based on community recommendations, an initial set of diagnostics was selected, published at <https://zenodo.org/records/14284375>, for incorporation into the initial version of the REF for evaluating relevant CMIP7 Assessment Fast Track simulations. Diagnostics 5.1 and 5.2 are listed in strikethrough style because they will not be implemented for the initial version of the REF.

ID	Diagnostic	Package(s)	Reference Dataset(s)
<i>Ocean & Sea Ice Realm</i>			
1.1	Antarctic annual mean, Arctic September rate of sea ice area (SIA) loss per degree warming (dSIA / dGMST)	ESMValTool	OSI SAF/CCI, HadCRUT
1.2	Atlantic meridional overturning circulation (AMOC)	IOMB	RAPID array
1.3	El Niño Southern Oscillation (ENSO) diagnostics (lifecycle, seasonality, amplitude, teleconnections)	PMP, ESMValTool	TropFlux, GPCP, HadISST, ERA5
1.4	Sea surface temperature (SST) bias, Sea surface salinity (SSS) bias	IOMB	GLODAP2 & WOA (climatology), HadISST (transient)
1.5	Ocean heat content (OHC)	IOMB	IAP v4.2
1.6	Antarctic & Arctic sea ice area seasonal cycle	ESMValTool	OSI SAF/CCI
<i>Land & Land Ice Realm</i>			
2.1	Soil carbon	ILAMB	HWSD2, NCSdV22
2.2	Gross primary production (GPP)	ILAMB	WECANN, FLUXNET2015
2.3	Runoff	ILAMB	Dai, LORA
2.4	Surface soil moisture	ILAMB	Wang & Mao
2.5	Net ecosystem carbon balance	ILAMB	Hoffman & Khatiwala
2.6	Leaf area index (LAI)	ILAMB	AVHRR & VIIRS, GIMMS
2.7	Snow cover	ILAMB	JASMES
<i>Atmosphere Realm</i>			
3.1	Annual cycle and seasonal mean of multiple variables	PMP, ESMValTool	ERA5, ESACCI-OZONE, GPCP, CERES-EBAF
3.2	Radiative and heat fluxes at the surface and top of atmosphere (TOA)	PMP	CERES-EBAF
3.3	Climate variability modes (e.g., ENSO, Madden-Julian Oscillation (MJO), Extratropical modes of variability, monsoon)	PMP	NOAA-20CR, HadISST
3.4	Evaporation minus precipitation ($E - P$)	ILAMB	GPCP (tentative)
3.5	Double inter-tropical convergence zone (ITCZ)	PMP	GPCP-SG
3.6	Cloud radiative effects	ESMValTool	CERES-EBAF, ESACCI-CLOUD
3.7	Scatterplots of two cloud-relevant variables (for specific regions of the globe and specific cloud regimes)	ESMValTool	ESACCI-CLOUD, GPCP-SG, CERES-EBAF, CALIPSO-ICECLOUD, ERA5
<i>Earth System Realm</i>			
4.1	Equilibrium climate sensitivity (ECS)	ESMValTool	N/A
4.2	Transient climate response (TCR)	ESMValTool	N/A
4.3	Transient climate response to cumulative emissions of carbon dioxide (TCRE)	ESMValTool	N/A
4.4	Zero emissions commitment (ZEC)	ESMValTool	N/A
4.5	Historical changes in climate variables (time series, trends)	ESMValTool	HadCRUT5, GCPC, ERA5
<i>Impacts & Adaptation Realm</i>			
5.1	High amplitude Rossby waves		
5.2	Internal variability or ensemble spread for individual models (precipitation and surface temperature)		
5.3	Evaluation of key climate variables at global warming levels	ESMValTool	N/A
5.4	Climate drivers for fire (fire burnt area, fire weather and fuel continuity)	ILAMB, ESMValTool	GFED5, MODIS MOD44B, ESA CCI Biomass, ISIMIP3a, GSWP3, W5E5



4.3 Data Request Opportunity

For CMIP7 the data request to modelling centres has been based on different “opportunities” that the community could submit to the CMIP IPO, following a community-wide call, each focusing on a specific scientific topic. Each opportunity contains a concise description of its scientific goals and a list of all variables that are requested from simulation output to achieve these
315 goals. The different data request opportunities are grouped together according to five realms (Ocean & Sea Ice, Land & Land Ice, Atmosphere, Earth System, and Impacts & Adaptation), which are the same as those that have then been used as a basis for selecting diagnostics for the REF (see Section 4.2).

A REF opportunity was submitted to the Data Request Task Team’s open call to ensure that the required variables for the REF were clearly defined and that modelling centres participating in the REF with their simulations would have a checklist
320 of variables needed for the REF diagnostics. Since the REF diagnostics span all five data request realms, a variable group request was submitted for each of the realms. The variable groups combined, form the REF CMIP7 data request opportunity. The opportunity was published as part of the v1.2.1 release of the data request. It consists of these five variable groups, containing 80 variables in total, many of which are from the recently developed list of baseline climate variables for Earth system modelling (Juckes et al., 2024), which is a subset of variables reflecting the most frequently used elements of CMIP6.
325 To facilitate finding information about the REF opportunity in the different data request documentation papers, it was decided that each paper would contain a very brief description of the opportunity and refer to the documentation paper about the atmosphere realm, where a more detailed REF opportunity description was added (Dingley et al., in preparation).

4.4 Observations

Observations play a key role in the REF and several diagnostics require observed quantities for different climate variables across
330 domains for model evaluation. For each diagnostic, we have identified at least one reference dataset for inclusion in the REF, and the complete list of reference datasets is contained in Appendix A. To ensure compliance with the Findability, Accessibility, Interoperability, and Reusability (FAIR) data standards (Wilkinson et al., 2016), all observational datasets included in the REF are required to have a fully open access license (e.g., CC-BY-4.0, CC0, OGL) and follow the CF Metadata Conventions (Davis et al., 2024) to ensure technical alignment with CMIP standard output. This is enabled by making the datasets available for
335 downloading on ESGF servers (Cinquini et al., 2014) as part of the obs4MIPs project (Gleckler et al., 2011; Teixeira et al., 2025; Ferraro et al., 2015; Waliser et al., 2020). For any datasets not meeting the open access criteria the REF Delivery Team obtained a relaxation of licence constraint by formally requesting and receiving agreement from individual data providers. Observational datasets for the REF were then processed in compliance with the obs4MIPs Data Specifications 2.5 (ODS2.5) (Gleckler et al., 2024) with approval from the obs4MIPs Steering Panel (OSP). In facilitating observational data ingestion for
340 the REF, two additional criteria not fully addressed by the ODS2.5 specifications were identified, for which the REF delivery team has drafted guidelines. These relate to the treatment of uncertainty information in the data and for the formulation of complex citations, which are discussed in Sections 4.5 and 4.6.



Two possible pathways were identified for observational dataset providers to submit their datasets for possible inclusion with the REF in the future. The first step in both cases is to submit a dataset proposal for approval by the OSP. This can
345 either be done directly by the dataset provider or a third party with permission from the dataset provider. Once the dataset has been approved by the OSP, the registered content, including the dataset name, version, data provider details and release date should be submitted to the Program for Climate Model Diagnosis and Intercomparison (PCMDI) obs4MIPs CMOR Tables repository, a process which provides the dataset with a unique `source_id`, following the CMIP conventions. Datasets can then be prepared for compliance in one of two possible ways:

- 350
- using the CMOR software as advised by the OSP to prepare their dataset, following instructions and examples on the PCMDI obs4MIPs CMOR GitHub repository – labelled the “CMOR pathway”.
 - using software packages such as ESMValTool to generate CMOR-like datasets that have additional scripts ensuring CMOR and obs4MIPs compliance – labelled the “CMOR-like pathway”.

The CMOR pathway intrinsically provides a form of validation for datasets before publication through the use of the CMOR
355 software. As the CMOR-like pathway may not provide the same level of validation, the REF Delivery Team has developed a validation script that the prepared CMOR-like datasets need to pass before publication. The prepared datasets are then published to ESGF via one of the two Assessment Fast Track REF nodes, at the Center for Environmental Data Analysis (CEDA, United Kingdom) or at Oak Ridge National Laboratory (ORNL, United States of America).

For certain diagnostics (1.1 and 5.4 in Table 1), the reference datasets listed in Table 1 are pre-processed to produce one
360 or more static files with information needed for the diagnostic. In such instances, these files are stored internally within the corresponding diagnostic package, and the REF includes acknowledgement and references to the input datasets used to produce the files. For such exceptions, the input datasets themselves are not required to be published on ESGF.

4.5 REF requirements in addition to obs4MIPs compliance – Treatment of uncertainties

Currently, most diagnostics within the REF do not incorporate uncertainty information from reference datasets. Where available
365 for a reference dataset and used in the diagnostics, the standard error on the mean or upper and lower bounds around the mean are used. More commonly, multiple datasets are used to account for observational uncertainty. The unavailability of comprehensive and accurately characterised uncertainty information with the observational data was identified as a key barrier to incorporating this information by diagnostic developers. Feedback from observational data providers also indicated that uncertainty information and how to correctly use this information with the data was best provided by the data providers
370 themselves. This prompted the REF Delivery Team to develop guidance for including uncertainty information in observational datasets; this was necessary to allow ingestion of uncertainty information by the REF. The initial proposal, originating from community engagement at the REF Hackathon with observational dataset providers and metrics package developers, was refined following consultation with the obs4MIPs Steering Panel and the CMIP-CVs Task Team. For the CMOR pathway, the following requirements are outlined:



- 375 – All additional uncertainty information (currently accepted information on uncertainty are described in Table 2) is provided in a separate file. Initially, between one and three additional files will be used, depending on the extent of uncertainty information given.
- The netCDF file containing the main geophysical variable has the global attribute `has_auxdata`, and is set to TRUE when additional files containing uncertainty information are provided. It should be set to FALSE when no uncertainty information is provided. If this attribute does not exist, the REF will assume there is no uncertainty information provided to ensure back-compatibility with datasets already published through obs4MIPs.
- 380 – If `has_auxdata` is set to TRUE, the netCDF file containing the main geophysical variable must have the global attribute `aux_variable_id`. This contains all additional `variable_ids` for the uncertainty information provided, in the form of a string with spaces as delimiters.
- 385 – Global attribute `variable_id` of file containing the uncertainty information corresponds to variable name+accepted suffix (no separator between body and suffix). For accepted suffixes, see Table 2.
- Variable name within the netCDF file corresponds to `variable_id`.
- A technical note providing detailed explanation of each additional uncertainty information provided, including type of uncertainty. The note should discuss what is included in the total uncertainty, by providing a breakdown of its components; such as sources of random, structured and systematic uncertainty.
- 390

For the CMOR-like pathway, there is no requirement to add the uncertainty information in separate files. Instead, additional uncertainty fields may be provided following CF conventions, by adding them as ancillary variables in a netCDF-CF file. The additional global attributes `has_auxdata` and `aux_variable_id` are still required, and the `variable_id` should be constructed as described above only using the suffixes from Table 2.

395 This distinction between how uncertainty information may be included in CMOR and CMOR-like datasets is based on the current capabilities of the CMOR software. Currently, CMOR does not accommodate the addition of ancillary variables in its output and the required software update is not feasible within the Assessment Fast Track REF timeline. This guidance was developed as a basic framework to enable the community to utilise a wider range of uncertainty information during the development of REF diagnostics.

400 **4.6 REF support of IPCC AR7 related to the enhanced traceability of its results – Complex citation**

The REF Delivery Team consulted with the IPCC Task Group on Data Support for Climate Change Assessments (TG-Data) on their citation and provenance requirements for the REF. IPCC TG-Data is currently enhancing the traceability of key results of the current Seventh Assessment Report (AR7; Stockhause et al., 2024) cycle using the new Complex Citation standard (Agarwal et al., 2025) for documentation and standardized provenance records for gathering the required information. This simple but flexible Complex Citation approach allows for the traceability of a data product generation and the citation of

405



Table 2. Suffixes for the `variable_id` were proposed to describe the uncertainty information included with the reference datasets.

Suffix	Long Name	Description
<code>utot</code>	Total uncertainty	The total per-datum uncertainty associated with the geophysical variable. If the independent, structured and common uncertainty components are also provided, this would be equal to the sum in quadrature of these components.
<code>nobs</code>	Number of observations	The number of discrete observations or measurements from which a data value has been derived.
<code>stderr</code>	Standard error	The standard error on the mean.
<code>ustr</code>	Structured uncertainty	The per-datum component of uncertainty that is structured and correlated over a defined space/time scale. This correlation length scale in space and time must be provided in the variable metadata.
<code>lbnd</code>	Lower bound	Alternative to <code>stderr</code> and <code>utot</code> for observations with asymmetrical uncertainty distribution, lower bound of the uncertainty.
<code>ubnd</code>	Upper bound	Alternative to <code>stderr</code> and <code>utot</code> for observations with asymmetrical uncertainty distribution, upper bound of the uncertainty.

multiple datasets or data subsets in a single referenced object called Complex Citation Object (CCO). It is planned to include CCO references in every figure caption. Prerequisite is the provision of detailed information on input data usage in form of persistent identifiers (PID) for each file and each citable entity.

The REF has identified the need for the reference datasets published on ESGF with their Handle IDs and data collections with their DOIs. ODS2.5 already requires each file to have a global attribute called `tracking_id` used as PIDs on ESGF, a unique identifier with specific prefixes specified for each ESGF project. In addition, reference datasets used by the REF also require to include the DOI as global attribute to ensure compatibility with CCO. REF leads the task of defining a provenance template for authors of the IPCC AR7 as guidance on providing the CCO-related information.

Through the CMIP7 data citation, the CCO captures provenance information at the granularity of a model's contribution to an experiment. The list of the file handles of the specific data that were used from within the CMIP7 citation resource provides traceability.

4.7 Technical Workflow of Both Versions

Figure 2 shows how the key services for the REF are integrated within the ESGF deployment. At the core of this system is the Compute Engine, which orchestrates the workflow by managing data ingestion and processing tasks. It interacts with the ESGF-Next Generation Event Stream, a Kafka-based service, to trigger data ingestion when new data become available. For the ESGF deployment, each of the services (blue boxes in Fig. 2) is deployed as a Kubernetes service. This decoupling of services enables each service to be scaled according to demand. This might be particularly important if one of the execution services requires more computational resources than the others. Users can access the REF through two primary interfaces: the

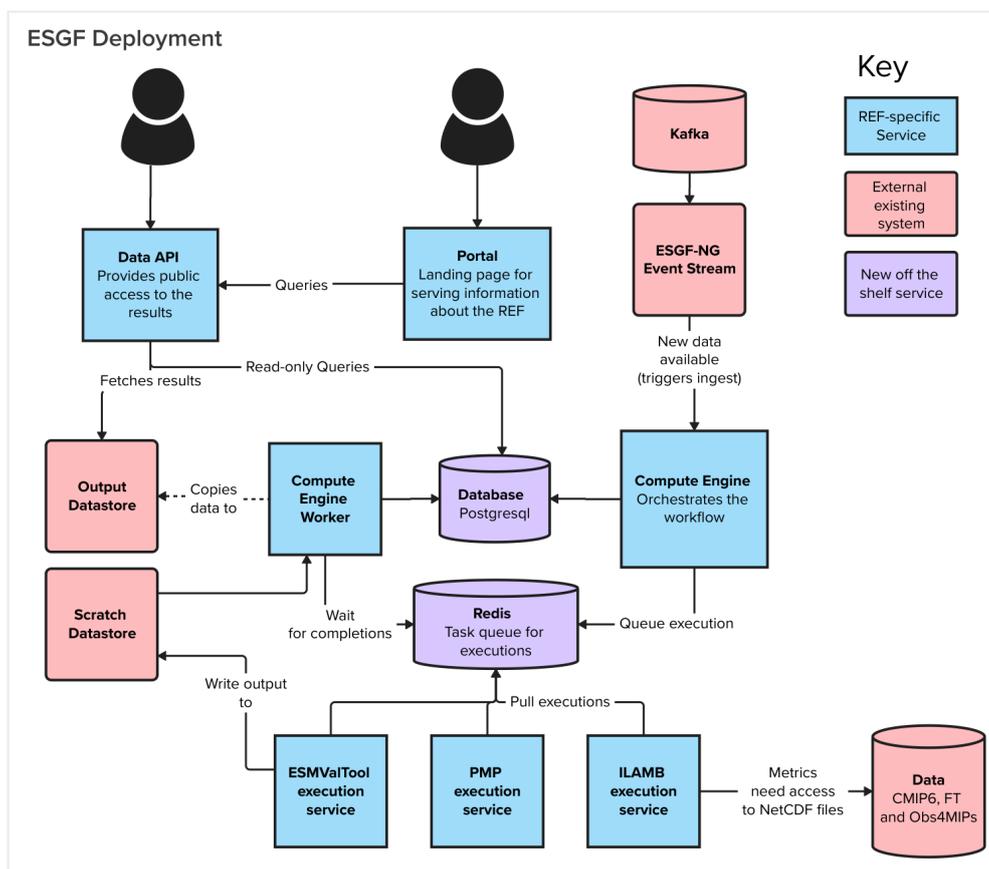


Figure 2. This detailed workflow diagram identifies key services provided by the REF and their relationships with external services and off-the-shelf components for the ESGF deployment. Figure available at <https://doi.org/10.5281/zenodo.15595006> (Lewis et al., 2025b).

Data API, which provides public access to full results, and the Portal, which serves as the landing page for information about the REF and provides synthesised results. Both interfaces allow users to query data, with the Data API fetching results from the Output Dastore.

Data processing is handled by the Compute Engine Worker, which executes tasks and manages data flow between the Database (PostgreSQL) and the Redis task queue. The Redis service queues execution tasks, which are then pulled by various execution services, including ESMValTool, PMP, and ILAMB. These services require access to the netCDF files, sourced from CMIP, and obs4MIPs datasets, to compute the diagnostics and underlying metrics.

Processed data are temporarily stored in the Scratch Dastore before being copied to the Output Dastore, ensuring data integrity and availability. This robust architecture supports efficient data processing and dissemination, enabling comprehensive climate data analysis within the REF framework.



Non-ESGF deployments will look very similar except that the ingestion and solving is performed manually, and there will
435 be an option to install the tool without using Kubernetes.

4.8 Technical Implementation

The REF workflow consists of four steps:

- **Ingestion:** The user registers the source datasets that can be used (reference and model). The metadata from these
440 datasets are extracted and added to a local database along with a path to the file. Only ingested files are used in any
execution calculations.
- **Solve:** For each diagnostic, the possible executions that would be required are determined using the data requirements
of a metric and the catalogue of datasets that have been ingested. A hash of the datasets required for each execution is
stored in the database and subsequently used to determine if an execution needs to be performed or has already been
performed.
- 445 – **Execute:** The executions that require running are then executed. The REF supports three key methods for executing
diagnostics: Local-Serial, Local-Parallel and Celery. The ESGF deployment uses the Celery-based executor, which runs
diagnostics out-of-process and in parallel. Once the execution is complete, the CMEC-based outputs are parsed and any
scalar/timeseries or figures are added to the database for later display and consumption.
- **Visualise:** The results are then made available via a REST API (Representational State Transfer Application Program-
450 ming Interface) and Typescript-based frontend. Python-based tooling may be developed in future to interact with the
results, but that is not planned as of the current release. The Frontend allows users to see the diagnostics that have been
executed and the corresponding results. This includes the ability to track which datasets were used for a given execution,
and the ability to download generated figures, datasets and log output. Box and whisker and time series figures are used
to summarise the metric values across all executions for a diagnostic.

455 4.9 Prerequisites for Use

The REF operates on CF-compliant netCDF model output files that utilize the CMIP CVs and reference datasets. When run
within ESGF deployments, the REF automatically accesses published model output and obs4MIPs reference data; however,
the REF also supports the analysis of datasets that are not published on ESGF, which allows modelling centres and individual
researchers to perform their own model evaluations. A deliberate decision was made to be able to evaluate data that do not con-
460 form with the CMIP6 or CMIP7 controlled vocabularies, enabling assessment of model versions not intended to be published
on ESGF. The local datasets must be available to all of the execution services, and ideally should have undergone a quality
control (QC) process to assure the metadata are correct and the data represent the desired output in the correct units.

The REF is designed to run on a variety of hardware platforms, from desktop computers to high performance computing
(HPC) environments, and can be run either within virtual Conda environments or within Docker containers. Depending on



465 the use case either may be used; however, Docker containers are recommended for a production deployment as they can be
scaled independently. The instructions for installation and getting started are available via the REF documentation site at <https://climate-ref.readthedocs.io/en/latest/> (last access: 19.05.2025) or the GitHub repository at <https://github.com/Climate-REF/climate-ref> (last access: 19.05.2025).

5 Discussion

470 5.1 Community Co-development of a Framework

The REF Delivery Team was composed of an international team that brought together three different evaluation package providers, each of which possessed their own workflows, assumptions, and conventions. An independent delivery team leader was chosen to coordinate the development and to architect the end-to-end workflow. Significant time was invested in identifying commonalities among the providers with the goal of minimising the amount of additional development needed to integrate all of
475 the selected diagnostics into the REF. These tools were never designed to be interoperable and have long-standing communities that use the packages outside of the REF, so a requirement for minimal changes to the underlying packages was also mandated. Developing the appropriate abstractions and common language within the REF depended upon ongoing discussions and often required the prototyping of multiple potential solutions.

The REF system can be quite complex since it is composed of multiple services. A key requirement for the framework was
480 to make it easy for additional packages to be integrated into the REF, both now and in the future. Care was taken to ensure that the amount of context that must be understood for a package developer to contribute diagnostics was kept to a minimum. This has two impacts, it provides a separation of concerns and minimises the complexity required simply to contribute diagnostics, as well as providing the scope to support the different assumptions made by different providers. The ability to quickly refactor and modify the interfaces between the packages was critical. The use of a monolithic repository, where all of the REF-related
485 packages were in a common GitHub repository with issue tracking and a Kanban board, was crucial for development and testing, in addition to enforcing the need for type hints throughout the code base and providing high test coverage. The REF is openly developed and available on GitHub at <https://github.com/Climate-REF/climate-ref> under the Apache License 2.0 open source licence. The version of the REF contemporaneous with the submission of this manuscript is archived on Zenodo at <https://doi.org/10.5281/zenodo.15103441> (Lewis et al., 2025a).

490 The REF Hackathon, hosted by the UK Met Office in Exeter in March 2025, was convened to bring together the REF Delivery Team and to solicit early feedback from potential users. The Hackathon successfully accelerated code development and reference data conversion, obtained feedback and discussions about approaches for computing a few of the diagnostics and supporting traceability for CCO, and tested an alpha version on the UK Met Office HPC facilities. The Hackathon was held relatively early in the development lifecycle of the project, so significant time was spent communicating and getting feedback
495 on the concepts of the REF as the product was still being actively developed. The research community continued to offer ideas for additional capabilities that could be enabled by the REF framework.



5.2 Key Reflections from the CMIP Panel

The CMIP Panel initiated the MBTT, with a call for members in November 2022. The intention of the Task Team was to put together a group of experts who would work together to provide a systematic, open, and rapid performance assessment of the expected large number of models participating in CMIP7, providing a set of informative diagnostics and performance metrics. The Task Team originally focussed on assessing evaluation approaches and software packages. However, it then became apparent that there was a bigger opportunity here to fully integrate evaluation tools into the CMIP publication workflow with diagnostic outputs published alongside model output on ESGF and results displayed on an easily accessible website.

The goal, informed by outcomes of the CMIP6 Community Survey, was to fully integrate the evaluation tools into the CMIP publication workflow, and their diagnostic outputs to be published alongside the model output on ESGF, with a request for results display in an easily accessible website.

The REF became an extremely attractive option to the CMIP panel as it enables 1) consistent assessment of CMIP output by modelling centres, 2) support for author teams in the context of the IPCC and other national climate assessments and 3) input for model selection for downstream applications. The modelling centres are essential to the CMIP endeavour so having a tool that supports their activities is the most important factor. Many of the CMIP panel members have contributed to IPCC assessments or national climate assessments, so it was clear that the work of author teams would be better supported if model diagnostics were easily available to authors (reducing the workload of both authors and chapter scientists). Wider community discussions facilitated by the CMIP panel have also revealed enthusiasm from users of CMIP output for regional downscaling and impacts applications, to have the diagnostics quickly available to support the choice of a limited selection of climate model output to serve their particular application.

The CMIP Panel view the REF as a potential game changer for users of CMIP data. CMIP Panel co-chair Helene Hewitt reflected, during opening remarks at the March 2025 Hackathon onboarding session, that “as someone that was a coordinating lead author of the IPCC AR6 WGI, I saw how much work our chapter scientists did in producing all the metrics and plots, this builds on the then-CMIP6 Panel vision, we are so happy that the task team have taken this on. We hope that being able to integrate these tools into the workflow and ESGF publication process, on a website, will make it much better for users to easily know what they are looking at. Going forward we want to streamline this chain, supporting model selection and downstream use of CMIP.”

5.3 Suggested Uses of the REF

The REF will produce diagnostics spanning a variety of metrics, maps, charts, and figures, as well as calculate scalar scores as a gauge of correspondence between model output and reference datasets. These scores and aggregations of scores across categories of diagnostics are not meant to discriminate “good models” from “bad models.” Instead, the REF is intended to assist analysts in quickly identifying relative differences among models or model versions that researchers must then interpret by viewing the plots and maps that underlie the individual diagnostics. Consumers of REF results must consider that only a limited number of diagnostics can be produced and evaluated, that reference datasets bring with them their own (usually



530 unquantified) uncertainties, and that all models are calibrated by making tradeoffs among parameter values and subjective
choices about the importance of predictions among different variables or model process representations. Most multi-model
assessments indicate that each model has strengths and weaknesses in different areas. For example, one model may perform
well with regard to the distribution of precipitation but may not exhibit the desired distribution of sea surface temperatures,
while other models may show the opposite behaviour. Similarly, some models may capture biogeochemical responses on land
535 well but have a weak representation of the hydrological cycle, while other models may capture runoff well but fail to capture
the seasonal cycle of terrestrial productivity. The REF results are best used to identify which subset of models may be best
for individual studies based on their performance in key science areas or spatial regions of interest. REF results may help
inform the selection of models to include in downscaling studies or the choice of multi-model weights when optimizing for a
particular metric. All of these suggested uses require the researcher to “drill down” into the detailed results produced in the
540 diagnostics, looking beyond performance scores that are merely a high-level indicator of large-scale average correspondence
of model output with a reference dataset.

5.4 The Future of the REF

5.4.1 Potential New Features and Capabilities

A significant advancement for the REF will be its adaptation for daily or sub-daily GCM output in order to capture synoptic-
545 and even mesoscale phenomena such as extra-tropical blocking events (Davini and D’Andrea, 2016; Dorrington et al., 2022;
Dolores-Tesillos et al., 2025; Bacer et al., 2022), storm tracks (Priestley et al., 2020), weather types (Brands, 2022a; Brands
et al., 2023), tropical cyclones (Roberts et al., 2020) and other phenomena. These smaller-scale phenomena are known to
drive extreme events, which have direct implications for societal impacts. Model evaluation across model generations involve
dimensionality reduction through empirical orthogonal function (EOF) analysis (Fasullo et al., 2020; Hannachi et al., 2023),
550 k -means clustering (Hoffman et al., 2005) and/or weather regime detection (Grams et al., 2017), for which a number of
efficient diagnostics tools have already been developed, e.g., Mid-latitude Evaluation System (MiLES; Davini, 2019) or GCM
evaluation with Lamb Weather Types (Brands, 2025).

The REF will likely be utilised by ongoing regional WCRP initiatives such as the Coordinated Regional Downscaling Ex-
periment (CORDEX; Giorgi and Gutowski, 2015; Diez-Sierra et al., 2022), the Inter-Sectoral Impact Model Intercomparison
555 Project (ISIMIP; Warszawski et al., 2014) and the Ice Sheet Model Intercomparison Project (ISMIP; Nowicki et al., 2016), and
significant potential exists to connect with other initiatives of this kind, e.g., Atmospheric River Tracking Method Intercom-
parison Project (ARTMIP; Rutz et al., 2019). Through such collaborations, the REF can serve to foster synergies that will be
mutually beneficial across research initiatives. A new collaboration effort, called the CORDEX Collection of Regional-Scale
Climate Processes and Metrics for Climate Model Evaluation, will develop high-resolution diagnostics that could be integrated
560 into the REF. As the need for higher spatial and temporal resolution grows, increasing temporal resolution is equally important
for understanding extreme weather events, such as the development of tropical cyclones or heatwaves. These high-frequency



variations, such as diurnal cycles or sub-daily phenomena, are key to better predicting impacts on society, and they are likely to be useful additions to the REF in the future.

Internal (or unforced) climate variability is commonly sampled by assessing the output of multiple historical (or scenario) runs of the same GCM, each initialized from different dates of the corresponding pre-industrial control run. In these equally probable surrogates of the real climate system, the main drivers of natural variability, such as ENSO, as well as its associated teleconnections, evolve freely through time. This “initial conditions uncertainty” (Stainforth et al., 2007) produces random noise, which compared to the predictable signal exerted by external forcing agents (e.g., greenhouse gases and aerosols), is particularly large for atmospheric circulation variables such as sea-level pressure and are not directly affected by global warming (Deser et al., 2012, 2020). Consequently, GCM performance estimates for these variables are expected to be likewise affected by internal variability, especially if they are based on short time periods. However, this kind of error uncertainty has seldom been evaluated in past model performance assessments, particularly outside the atmosphere. Considering internal variability in future versions of the REF is a priority, as it will enable assignment of uncertainty ranges to model performance estimates, thereby making them more robust. This step will also improve understanding of the natural drivers of regional climates and provide better insights into how climate extremes evolve under different scenarios.

Another dimension along which the REF is expected to grow is the consideration of other components of the Earth system beyond the physical processes included in the coupled model configurations contributing to CMIP, such as atmospheric chemistry and terrestrial and marine biogeochemistry, which aim to provide a more accurate representation of the global carbon cycle (Séférian et al., 2020) and its feedback on the Earth system (Arora et al., 2020). As the community increasingly focuses on carbon emissions-driven experiments in CMIP7, the REF must be extended to quantify and reconcile carbon cycle biases (Hoffman et al., 2014) and to include the spatial and temporal evaluation of atmospheric CO₂ and CH₄ variability (Keppel-Aleks et al., 2013). This could help improve understanding of the relationships between carbon emissions and Earth system feedbacks, particularly in high-resolution models. Additionally, systematic evaluation of observed emerging signals of subsurface ocean acidification and deoxygenation, alongside warming, is desirable to benchmark transient changes (Tjiputra et al., 2023). Including these aspects will allow the REF to support a more comprehensive evaluation of oceanic changes that have substantial long-term impacts on marine ecosystems and the global carbon cycle.

A key component of model evaluation is understanding model similarity, which can inform how multi-model outputs should be weighted. This can be accomplished *a priori* by analysing similarities in the model architecture or *a posteriori* by analysing similarities in the model output data (Boé, 2018; Brands, 2022b; Merrifield et al., 2023). For the *a priori* approach to work, it is essential to have knowledge of the sub-models used in the coupled model configurations. A metadata archive that collects the pedigree of model components, resolution details, and other relevant information (e.g., Brands et al., 2023) will allow for better tracking of model characteristics and help improve the comparability of model simulations. As the REF evolves, exploring model similarity could become a more prominent feature, helping to refine how ensemble simulations are interpreted and used in Earth system studies and assessments.

The REF will be continuously exercised and updated through the CMIP7 process, and that will undoubtedly reveal opportunities for improvements. Along with such optimizations and incremental improvements, the future evolution of the REF will



emerge as it is increasingly used. Many potential extensions of the REF are already envisioned, as illustrated by the above discussion. The expectation is that additional opportunities will be identified and new features will be desired as the next generation of CMIP models are scrutinized. In designing the REF, such evolution has been anticipated, and the design of the software is modular and flexible. This should facilitate community contributions to the REF. Although all the realms of ESMS are already represented by the REF, we anticipate new diagnostic packages may need to be incorporated into the REF. Packages that focus on high-frequency variability (e.g., diurnal cycle, extreme events at the sub-daily time scale) or specific phenomena (e.g., Madden-Julian Oscillation (MJO), ENSO) would be natural extensions of the current capabilities of the REF. Similarly, observation-based reference datasets will need to be updated as new observations are collected or new products are made available. In all of these future developments, a key aspect for the vitality of the REF will be community engagement. Community contributions and comments are welcome and necessary, as the REF is intended to be an open, community-driven project.

5.4.2 REF Governance

The MBTT and co-sponsors of the Assessment Fast Track version of the REF agreed, based on significant community interest, that continued development of the REF should be coordinated and prioritized by an international consortium of community evaluation package developers, modellers, reference data providers, and Earth system scientists under WCRP. Such a scientific steering panel (SSP) would engage with the community to identify high priority diagnostics and reference datasets for the full set of CMIP7 simulations and other participating WCRP activities; coordinate with ESGF on quality assurance, platform upgrades, and performance enhancements; initiate and coordinate a technical development and support community; and ensure the open extensibility and portability of subsequent REF developments to support use of the evolving REF framework by modelling centres and individual researchers beyond CMIP activities. An interim scientific steering panel (SSP) will be established to carry out these duties and to identify and propose to WCRP preferred long-term governance arrangements. This SSP will remain in place until a formal governance approach is approved by WCRP and handover arrangements can be implemented. An open global recruitment process will be conducted to select interim SSP members. Contributors to sustainment and the future engineering of the REF will likely acquire their own funding or support for collaborative maintenance and development, which will be coordinated through the REF-SSP.

6 Conclusions

The Rapid Evaluation Framework (REF) is an open source Python-based toolkit (compatible with versions ≥ 3.11) under active development, designed to automate and manage computational evaluations of Earth system model (ESM) output. Its primary objective is to enable near-real-time assessment of ESM output through comparison with best-available reference (observationally constrained benchmark) datasets, updating outputs dynamically as new simulation results are published. Functionally analogous to a Continuous Integration/Continuous Deployment (CI/CD) pipeline in software development, the REF streamlines continuous evaluation workflows for Earth system science. A beta version of the REF, targeted for used and testing by modelling centres and interested researchers, was released to the public on 27 May 2025.



Proposed by the MBTT, the REF is initially deployed to support the Assessment Fast Track simulation campaign, providing
630 diagnostics collaboratively selected by the MBTT and the broader science community. While initially tailored for the Assessment Fast Track, the framework is intentionally agnostic to data types and analytical metrics, ensuring adaptability for diverse Earth science applications beyond its initial scope. This extensibility underscores its potential long-term research utility. Key technical features include integration with CI/CD systems, availability on PyPI from version v0.5.0, comprehensive documentation, and community-driven development under an Apache License 2.0 open source license. Emphasizing collaboration, the
635 REF invites contributions from researchers and developers, fostering a shared ecosystem for advancing model-data evaluation tools. The project's design prioritizes scalability and flexibility, aiming to serve as a foundational resource for real-time ESM benchmarking in both current and future research contexts. Interest in the REF across the research community is strong, and future governance of REF development will be coordinated and prioritized by a scientific steering panel that has not yet been formed. A wide variety of new diagnostics are already proposed for integration in subsequent developments of the REF.

640 *Code and data availability.* The REF is openly developed and available on GitHub at <https://github.com/Climate-REF/climate-ref> under the Apache License 2.0 open source licence. The version of the REF contemporaneous with the submission of this manuscript is archived on Zenodo at <https://doi.org/10.5281/zenodo.15103441> (Lewis et al., 2025a). Reference data used by the REF are available from the Earth System Grid Federation (ESGF) nodes listed at <https://esgf.github.io/nodes.html>. Original reference data are listed with appropriate citations in Table A1.



645 Appendix A: Reference Datasets Used by the REF

Table A1. Observationally constrained reference datasets shown here are used in the REF for comparison with model output.

Reference Dataset	Citation	Associated Diagnostic ID(s)
AVHRR & VIIRS	Claverie et al. (2014, 2024)	2.6
CALIPSO-ICECLOUD	Winker (2024); Winker et al. (2024)	3.7
CERES-EBAF	Loeb et al. (2018)	3.1, 3.2, 3.6, 3.7
Dai	Dai (2017); Dai and Trenberth (2002); Dai et al. (2009); Dai (2016, 2021)	2.3
ERA5	Hersbach et al. (2020)	1.3, 3.1, 3.7, 4.5
ESA CCI Biomass	Santoro and Cartus (2024)	5.4
ESACCI-OZONE	Copernicus Climate Data Store (2020); Sofieva et al. (2023); Coldewey-Egbers et al. (2025)	3.1
FLUXNET2015	Pastorello et al. (2020)	2.2
GFED5	Chen et al. (2023a, b)	5.4
GIMMS	Cao et al. (2023a, b)	2.6
GLODAP2	Key et al. (2004); Olsen et al. (2016)	1.4
GPCP & GPCP-SG	Adler et al. (2017, 2003, 2018)	1.3, 3.1, 3.4, 3.5, 3.7
GSWP3	Dirmeyer et al. (2006)	5.4
HadCRUT5	Morice et al. (2021)	1.1, 4.5
HadISST	Rayner et al. (2003)	1.3, 1.4, 3.3
Hoffman	Hoffman et al. (2014)	2.5
HWSD2.0	FAO and IIASA (2023)	2.1
IAP v4.2	Cheng et al. (2024)	1.5
ISIMIP3a	Lange et al. (2022, 2024)	5.4
JASMES	Hori et al. (2017)	2.7
LORA	Hobeichi et al. (2019)	2.3
MODIS MOD44B	DiMiceli et al. (2015)	5.4
NCSDv22	Hugelius et al. (2013)	2.1
NOAA-20CR	Slivinski et al. (2019)	3.3
OSI SAF/CCI	Lavergne et al. (2019)	1.1, 1.6
RAPID array	Moat et al. (2025)	1.2
TropFlux	Kumar et al. (2012)	1.3
W5E5	Lange (2019); Cucchi et al. (2020)	5.4
Wang & Mao	Wang and Mao (2021); Wang et al. (2021)	2.4
WECANN	Alemohammad et al. (2017)	2.2
WOA	Reagan et al. (2023)	1.4

<https://doi.org/10.5194/egusphere-2025-2685>

Preprint. Discussion started: 11 July 2025

© Author(s) 2025. CC BY 4.0 License.



Appendix B: Diagnostics Produced by the REF



Table B1. Based on community recommendations, an initial set of diagnostics was selected, published at <https://zenodo.org/records/14284375>, for incorporation into the initial version of the REF for evaluating relevant CMIP7 Assessment Fast Track simulations. Diagnostics 5.1 and 5.2 are listed in strikethrough style because they will not be implemented for the initial version of the REF for the CMIP7 Assessment Fast Track.

ID	Realm	Diagnostic	Package(s)	CMIP Variable(s)	Assessment Fast Track Experiment(s)	Reference Dataset(s)
1.1	Ocean & Sea Ice	Antarctic annual mean, Arctic September rate of sea ice area (SIA) loss per degree warming (dSIA/dGMST)	ESMValTool	siconc, tas, areacello, areacella	historical, esm-historical, hist-nat, hist-aer, hist-GHG	OSI SAF/CCI, HadCRUT
1.2	Ocean & Sea Ice	Atlantic meridional overturning circulation (AMOC)	IOMB	msftmz	historical, esm-historical, hist-nat, hist-aer, hist-GHG	RAPID array
1.3	Ocean & Sea Ice	El Niño Southern Oscillation (ENSO) diagnostics (lifecycle, seasonality, amplitude, teleconnections)	PMP, ESMValTool	pr, tos, areacello, ts, tauu	historical, esm-historical, hist-nat, hist-aer, hist-GHG	TropFlux, GPCP, HadISST, ERA5
1.4	Ocean & Sea Ice	Sea surface temperature (SST) bias, Sea surface salinity (SSS) bias	IOMB	tos	historical, esm-historical, hist-nat, hist-aer, hist-GHG	GLODAP2 & WOA (climatology), HadISST (transient)
1.5	Ocean & Sea Ice	Ocean heat content (OHC)	IOMB		historical, esm-historical, hist-nat, hist-aer, hist-GHG	IAP v4.2
1.6	Ocean & Sea Ice	Antarctic & Arctic sea ice area seasonal cycle	ESMValTool	siconc, areacello	historical, esm-historical, hist-nat, hist-aer, hist-GHG	OSI SAF/CCI
2.1	Land & Land Ice	Soil carbon	ILAMB	cSoil	historical, esm-historical, land-hist, hist-nat, hist-aer, hist-GHG	HWSD2, NCSDv22

continued on next page



continued from previous page

ID	Realm	Diagnostic	Package(s)	CMIP Variable(s)	Assessment Fast Track Experiment(s)	Reference Dataset(s)
2.2	Land & Land Ice	Gross primary production (GPP)	ILAMB	gpp	historical, esm-historical, land-hist, hist-nat, hist-aer, hist-GHG	WECANN, FLUXNET2015
2.3	Land & Land Ice	Runoff	ILAMB	mrro, mrros	historical, esm-historical, land-hist, hist-nat, hist-aer, hist-GHG	Dai, LORA
2.4	Land & Land Ice	Surface soil moisture	ILAMB	mrsos	historical, esm-historical, land-hist, hist-nat, hist-aer, hist-GHG	Wang & Mao
2.5	Land & Land Ice	Net ecosystem carbon balance	ILAMB	nbp, netAtmosLandCO2Flux	historical, esm-historical, land-hist, hist-nat, hist-aer, hist-GHG	Hoffman
2.6	Land & Land Ice	Leaf area index (LAI)	ILAMB	lai	historical, esm-historical, land-hist, hist-nat, hist-aer, hist-GHG	AVHRR & VIIRS, GIMMS
2.7	Land & Land Ice	Snow cover	ILAMB	snc	historical, esm-historical, land-hist, hist-nat, hist-aer, hist-GHG	JASMES
3.1	Atmosphere	Annual cycle and seasonal mean of multiple variables	PMP, ESMValTool	o3, pr, prw, psl, rlds, rlus, rlut, rlutcs, rlds, rdsdcs, rsdt, rsut, rsutes, sfcWind, ta, tas, tauu, toz, ts, ua, va, zg	amip, historical, esm-historical, hist-nat, hist-aer, hist-GHG	ERA5, ESACCI-OZONE, GPCP, CERES-EBAF

continued on next page



continued from previous page

ID	Realm	Diagnostic	Package(s)	CMIP Variable(s)	Assessment Fast Track Experiment(s)	Reference Dataset(s)
3.2	Atmosphere	Radiative and heat fluxes at the surface and top of atmosphere (TOA)	PMP	rlds, rlus, rlut, rlutcs, rsds, rsdscs, rsdt, rsut, rsutcs	amip, historical, esm-historical, hist-nat, hist-aer, hist-GHG	CERES-EBAF
3.3	Atmosphere	Climate variability modes (e.g., ENSO, Madden-Julian Oscillation (MJO), Extratropical modes of variability, monsoon)	PMP	hfsl, hfss, pr, rlds, rlus, rsds, rsus, taux, ts	amip, historical, esm-historical, hist-nat, hist-aer, hist-GHG	NOAA-20CR, HadISST
3.4	Atmosphere	Evaporation minus precipitation ($E - P$)	ILAMB	pr, prsn, evspbl	historical, esm-historical, land-hist, hist-nat, hist-aer, hist-GHG	GPCP (tentative)
3.5	Atmosphere	Double inter-tropical convergence zone (ITCZ)	PMP	pr	amip, historical, esm-historical, hist-nat, hist-aer, hist-GHG	GPCP-SG
3.6	Atmosphere	Cloud radiative effects	ESMValTool	rlut, rlutcs, rsut, rsutcs	amip, historical, esm-historical, hist-nat, hist-aer, hist-GHG	CERES-EBAF, ESACCI-CLOUD
3.7	Atmosphere	Scatterplots of two cloud-relevant variables (for specific regions of the globe and specific cloud regimes)	ESMValTool	clt, cli, clivi, clwvi, pr, rlut, rlutcs, rsut, rsutcs, ta	amip, historical, esm-historical, hist-nat, hist-aer, hist-GHG	ESACCI-CLOUD, GPCP-SG, CERES-EBAF, CALIPSO-ICECLOUD, ERA5
4.1	Earth System	Equilibrium climate sensitivity (ECS)	ESMValTool	tas, rsdt, rsut, rlut, (rtmt)	piControl, abrupt-4xCO2	N/A
4.2	Earth System	Transient climate response (TCR)	ESMValTool	tas	1pctCO2*, piControl	N/A
4.3	Earth System	Transient climate response to cumulative emissions of carbon dioxide (TCRE)	ESMValTool	tas, fco2antt	esm-1pctCO2*, esm-piControl, esm-flat10	N/A

continued on next page



continued from previous page

ID	Realm	Diagnostic	Package(s)	CMIP Variable(s)	Assessment Fast Track Experiment(s)	Reference Dataset(s)
4.4	Earth System	Zero emissions commitment (ZEC)	ESMValTool	tas	esm-flat10, esm-flat10-zec, 1pctCO2*, esm-1pct-brch-1000Pg*	N/A
4.5	Earth System	Historical changes in climate variables (time series, trends)	ESMValTool	tas, pr, psl, ua, hus	amip, historical, esm-historical, hist-nat, hist-aer, hist-GHG	HadCRUT5, GCPC, ERA5
5.1	Impacts & Adaptation	High amplitude Rossby waves				
5.2	Impacts & Adaptation	Internal variability or ensemble spread for individual models (precipitation and surface temperature)				
5.3	Impacts & Adaptation	Evaluation of key climate variables at global warming levels	ESMValTool	tas, pr	historical, esm-historical, scenarios, hist-nat, hist-aer, hist-GHG	N/A
5.4	Impacts & Adaptation	Climate drivers for fire (fire burnt area, fire weather and fuel continuity)	ILAMB, ESMValTool	pr, tasmax, treeFrac, vegFrac, baresoil, ps, huss, sfcWind, lai, mrso, vegFrac, npp	historical, esm-historical, land-hist, hist-nat, hist-aer, hist-GHG	GFED5, MODIS MOD44B, ESA CCI Biomass, ISIMIP3a, GSWP3, W5E5

*CMIP6 variable; not available/necessary for CMIP7 Assessment Fast Track

650



Appendix C: Rationale for Diagnostics Produced by the REF

This Appendix contains the list of diagnostics, v1.0 (CMIP Model Benchmarking Task Team, 2024), originally devised by the CMIP Model Benchmarking Task Team to be implemented in the CMIP7 Assessment Fast Track Rapid Evaluation Framework
655 (REF).

C1 Ocean & Sea Ice Realm

1.1 Antarctic annual mean, Arctic September rate of sea ice area (SIA) loss per degree warming (dSIA / dGMST)

Previous sea ice benchmarking assessments have highlighted the fact that CMIP models systematically underestimate the amount of Arctic sea ice area loss per degree of global warming. Very few CMIP models are able to simulate both a plausible sea
660 ice loss and a plausible change in global mean temperature over the satellite period. This metric evaluates the rate of sea ice loss per degree of global warming, following the approach used for sea ice benchmarking within the Sea Ice Model Intercomparison Project analysis (Notz and SIMIP Community, 2020; Roach et al., 2020). The metric is calculated by regressing the time-series of sea ice area on global mean temperature. Sea ice responds strongly to climate forcing and warming. The rapid decline in Arctic summer sea ice areal extent is a highly visible indicator of climate change.

665 1.2 Atlantic meridional overturning circulation (AMOC)

The Atlantic Meridional Overturning Circulation (AMOC) provides a key indicator of the strength of ocean circulation, which redistributes freshwater, heat and carbon across the Atlantic Basin (Le Bras et al., 2023). A weakening of AMOC, expected to occur as a result of ocean warming, will have global climate consequences. The strength of the AMOC at 26.5°N is commonly used for evaluation of model fidelity since it can be compared with the long-term RAPID-MOCHA (Rapid Climate Change -
670 Meridional Overturning Circulation and Heatflux Array) observational dataset (Moat et al., 2025). RAPID datasets are widely used to validate the AMOC strength in models and are available from 1st April 2005 to 11th February 2023. The AMOC at 26.5°N is calculated as the maximum of the meridional overturning streamfunction, which is provided by CMIP models as the variable “msftmz”. The AMOC is a key component of the global ocean conveyor belt and plays an important role in transporting heat poleward and ocean biogeochemical tracers from the surface into the ocean interior.

675 1.3 El Niño Southern Oscillation (ENSO) diagnostics (lifecycle, seasonality, amplitude, teleconnections)

The El Niño Southern Oscillation (ENSO) is the primary mode of the global interannual climate variability, mainly reflected by the variations in surface wind stress and ocean temperature in the tropical Pacific Ocean. Through teleconnections, the ENSO affects seasonal temperature and precipitation in other parts of the globe (Chen and Wallace, 2015; Vaittinada Ayar et al., 2023). The ENSO variability can be calculated from both sea surface temperature and atmospheric pressure differences
680 between different tropical Pacific areas. The Southern Oscillation Index (SOI) uses pressure differences between the Tahiti and Darwin regions. The Oceanic Niño Index (ONI) summarizes SST anomalies in the Niño 3.4 region. Given its implications for



regional climate variability, capturing the observed ENSO spatial and temporal characteristics would increase the fidelity and robustness in a model's climate projections.

1.4 Sea surface temperature (SST) bias, Sea surface salinity (SSS) bias

685 The SST and SSS distributions provide large scale patterns of surface ocean circulation as well as reflecting dynamical air-sea interactions and ocean-sea ice interactions in the polar regions. SST and SSS biases have a significant impact on the coupling of ESM's two major components, the atmosphere and the ocean. Satellite data products and localized moored sensors are used to produce measurements that are incorporated into reference data to calculate SST and SSS biases in models.

1.5 Ocean heat content (OHC)

690 The majority of the extra energy associated with climate change is stored and circulated in oceanic layers at almost all depths. The OHC may provide one of the most reliable signals about the long-term climate change and decadal to multidecadal variability, including their temporal variation and spatial patterns. It is compared, between models and observations, on a gridded basis ($1^\circ \times 1^\circ$), based on almost all available *in situ* ocean observations (e.g., Argo, conductivity-temperature-depth (CTD) profilers, Mechanical Bathythermographs, bottles, moorings, gliders, and animal-borne ocean sensors; Cheng et al.,
695 2024). Before use, the data are carefully bias corrected, vertically and horizontally interpolated and mapped onto a grid for comparison with models.

1.6 Antarctic & Arctic sea ice area seasonal cycle

The sea ice area, calculated as the sum over the Northern (Arctic) and Southern (Antarctic) Hemisphere grid cell areas multiplied by the sea ice fraction within each cell, exhibits a distinct seasonal cycle. Arctic sea ice area typically has minimum
700 values in September, while Antarctic sea ice area is lowest in February. The seasonal cycle is driven by the seasonal cycle of the insolation, sea ice processes, as well as the exchange with the atmosphere and ocean and can be seen as an overview metric for the general state of the sea ice in a model. Since sea ice has a much higher albedo than the ocean surface, sea ice area plays an important role in the surface energy and radiation budgets. In addition to the multi-year average seasonal cycle of Arctic and Antarctic sea ice area, the diagnostic produces time series of the September (Arctic) and February (Antarctic) sea ice area.

705 C2 Land & Land Ice Realm

2.1 Soil carbon

Soil carbon is the organic matter and inorganic carbon in global soils. It is an important component of the global carbon cycle and affects soil moisture retention and saturation. Soils are the largest stores of carbon on Earth, and warming can impact the ability of soils to store and retain carbon, especially in the Arctic, where permanently frozen soil can release large quantities
710 of carbon into the atmosphere if it thaws due to warming (Trumbore and Czimczik, 2008). Analyzing stored soil carbon helps track quantify the dynamics of the terrestrial carbon cycle within models and the movement of carbon through the Earth system.



2.2 Gross primary production (GPP)

Gross primary production is the process by which plants “fix” atmospheric or aqueous carbon dioxide through photosynthetic reduction into organic compounds, and it is affected by increases in atmospheric carbon dioxide (CO₂) levels and warming (Anav et al., 2015). A fraction of gross primary productivity supports plant respiration and the rest is stored as biomass in stems, leaves, roots, or other plant parts. Land use change, heat and drought stress due to anthropogenic warming, and rising atmospheric CO₂ will differentially influence gross primary production in ecosystems and alter the global carbon cycle. Thus, models must be evaluated to ensure they capture the observed responses to these changes.

2.3 Runoff

Surface water runoff plays an important role in the hydrological cycle by returning excess precipitation to the oceans and controlling how much water flows into water systems (Trenberth et al., 2007; Trenberth and Caron, 2001). Changes in atmospheric circulation and distributions of precipitation have a direct effect on changes in runoff from land. Models must be evaluated to ensure they exhibit the observed responses to precipitation and soil moisture processes that lead to runoff and transport of freshwater into rivers and oceans.

2.4 Surface soil moisture

Surface soil moisture is an important hydrological cycle variable governing interactions between the land surface and atmosphere, and with the oceans through surface water runoff. Soil moisture partitions incoming energy into latent and sensible heat fluxes and also controls how precipitation is partitioned into runoff or for evapotranspiration. It therefore is at the nexus of the carbon, energy, and water cycles and is key to understanding a broad range of processes from drought, floods to agricultural management. Models must be evaluated to ensure they exhibit observed responses in soil moisture to variation in precipitation, temperature, land use change, soil carbon content and biogeochemical cycles (Seneviratne et al., 2010).

2.5 Net ecosystem carbon balance

The net ecosystem carbon balance is a quantitative estimate of the overall terrestrial carbon uptake or loss, sometimes referred to as the land carbon sink Keenan and Williams (2018). The land sink represents the annual carbon uptake on land through gross primary production minus losses through ecosystem respiration, disturbance, wood and agricultural harvest, and land use change. Along with the global marine carbon sink, the global terrestrial carbon sink is important for sequestering anthropogenic carbon from the atmosphere and is influenced by climate change. Models must be evaluated to ensure they exhibit observed responses in the net ecosystem carbon balance since it directly affects the amount anthropogenic carbon retained in the atmosphere (Le Quéré et al., 2018).



740 **2.6 Leaf area index (LAI)**

Leaf area index (LAI) is a measure of the total area of leaves per unit ground area, and it is used to characterize the structure of plant canopies (Fang et al., 2019). Thus, it is an important variable to model mass and energy exchange through leaf surfaces between the biosphere and atmosphere. LAI is influenced gross primary production and plant physiological allocation of carbon to leaves. Models must be evaluated to ensure they exhibit observed responses to seasonal and trend changes in LAI in response to variations in temperature, precipitation, soil moisture, nutrient availability, and atmospheric CO₂.

745 **2.7 Snow cover**

Globally, snow cover is a key determinant of the Earth surface albedo and is also known to affect the large-scale atmospheric circulation, particularly in the Northern Hemisphere (Barnett et al., 1989; Cohen and Entekhabi, 1999). It is also a key element of the Arctic Amplification phenomenon (Cohen et al., 2014) and directly impacts the terrestrial hydrological cycle by providing a meltwater source in mid-latitudes. On the regional scale, snow cover is a paramount climate driver in the mid-to-high latitudes of the Northern Hemisphere, where it determines the length of the winter season and associated significant changes in hydrology, soil properties, and vegetation activity. A persistent snow-cover essentially blocks interactions between the atmosphere and underlying land-surface (Bokhorst et al., 2016). Models must be evaluated to ensure they exhibit observed responses to seasonal and trend changes in precipitation, temperature, and other important driving mechanisms.

755 **C3 Atmosphere Realm**

3.1 Annual cycle and seasonal mean of multiple variables

The annual cycle provides an integrative measure of skill at one of the fundamental forced time scales, yet ESMs often exhibit pathologic biases in the phase or amplitude of key quantities (e.g., Scaife et al., 2010; Hoffman et al., 2014; Miller et al., 2015). Evaluating the seasonal and annual variability in models helps to ensure that key processes are correctly represented in models.

760 **3.2 Radiative and heat fluxes at the surface and top of atmosphere (TOA)**

The radiative and heat fluxes at the surface and top of atmosphere (TOA) represent the fundamental flows of energy through the climate system. The TOA radiative budget indicates the relative disequilibrium of the climate system, and is the primary approach to quantifying the equilibrium climate sensitivity. Similarly, capturing the transient response at TOA is a basic measure of a model's ability to realistically represent the system response to forcing (Loeb et al., 2020). Fluxes at the surface quantify the energy exchange between the atmosphere and surface and are directly related to the redistribution of heat through the system and the ocean heat uptake (Mayer et al., 2024). Imbalances at the TOA in unforced experiments are a typical source of bias, inducing spurious drifts in model behavior (e.g., Mauritsen et al., 2012), inducing an erroneous interpretation of the thermodynamics of the climate system (Lucarini et al., 2017) and in particular of the relation between heat gradients and the general circulation of the coupled atmosphere-ocean system, as well as the patterns of the response of the system to an inho-



770 homogeneous forcing (Irving et al., 2019; Lembo et al., 2019). Calculations of radiative and heat fluxes at the surface and TOA
are usually carried out by the radiative transfer modules in atmospheric models, and parametrization of latent and sensible heat
fluxes at the surface through bulk formulas. The imbalances are obtained combining these fluxes, whereas transports can either
be implied through integration of fluxes (Trenberth and Caron, 2001), or explicitly retrieved through evaluation of the internal
energy within the atmosphere (Moist Static Energy; MSE) or inside the oceans (ocean heat content; Cheng et al., 2024) and
775 other subdomains.

3.3 Climate variability modes (e.g., ENSO, Madden-Julian Oscillation (MJO), Extratropical modes of variability, monsoon)

The main modes of low-frequency variability in the atmosphere, such as the North Atlantic, Arctic and Antarctic Oscilla-
tions, the Pacific North American and Pacific South American patterns, are important drivers of climate variability on the
780 hemispheric to continental-scale because they determine the brought direction and strength of the atmospheric flow (Wallace
and Gutzler, 1981). Since the flow determines the temperature and moisture characteristics of the air masses transported to a
any specific region, these modes also explain a significant fraction of the regional-scale climate variability around the globe
(Hurrell et al., 2001). Originating in the equatorial Pacific, the El Niño-Southern Oscillation (ENSO) is the most important
(ocean-atmosphere) mode of climate variability and is associated with typical climate anomalies in many regions around the
785 globe (Trenberth et al., 1998). Most of these modes are considered internal (or unforced) oscillations, meaning that their
characteristics are largely robust to anthropogenic forcing (Deser et al., 2012). Due to their large scale, they are the primary
diagnostics a global climate model should be able to reproduce, before one would proceed to evaluate model performance on
smaller scales (Fernández-Granja et al., 2024).

3.4 Evaporation minus precipitation ($E - P$)

790 The evaporation minus precipitation is a measure of the hydrologic cycle strength and determines the net water flux at the
surface. The spatial structure of $E - P$ determines the extent of arid regions of the subtropics, and ESMs show diverse changes
in this $E - P$ patterns depending mainly on how their regional circulations vary (Elbaum et al., 2022). Models must be evaluated
to ensure they capture the observed spatial and temporal variability of the global water balance.

3.5 Double inter-tropical convergence zone (ITCZ)

795 The Inter-tropical Convergence Zone (ITCZ) is a circumglobal, narrow low-pressure trough primarily located in the tropics. It
is characterized by converging horizontal winds triggering deep convection and heavy precipitation and its position marks the
onset and cessation of the monsoon. The ITCZ is commonly measured by precipitation, sea-level pressure or outgoing long-
wave radiation anomalies. Over the central and eastern tropical Pacific Ocean, and also over the equatorial Atlantic Ocean,
a single zonal ITCZ structure is located directly North of the Equator. A prominent and long-standing GCM artefact is a
800 second, parallel structure located approximately at 10°S that is most pronounced during austral summer and does not appear



in observations, leading to a “double ITCZ” in the model world (Tian and Dong, 2020; Ma et al., 2023). This error is closely related to a model misinterpretation of the climatological sea-surface temperature patterns the central-to-eastern tropical South Pacific Ocean (20°S–0°, 100°–150°W; Oueslati and Bellon, 2015). A simple metric to quantify the double ITCZ problem is the spatial correlation of the modelled vs. observed grid-box-scale climatological precipitation amounts in the double ITCZ (DI) region (20°S–0°, 100°–150°W; Oueslati and Bellon, 2015) during austral summer (DJF). A second, more sophisticated approach, would imply the construction of a Taylor diagram with the aforementioned fields.

3.6 Cloud radiative effects

Clouds play an important role in climate by reflecting incoming solar radiation (shortwave) and by absorbing and emitting outgoing thermal radiation (longwave). These cloud radiative effects can be quantified by calculating the differences in TOA clear-sky and all-sky radiative fluxes. The diagnostic calculates multi-year annual average shortwave and longwave TOA cloud radiative effects that are then used to create maps and zonal means of the cooling and warming effects of clouds of the models in comparison to observations.

3.7 Scatterplots of two cloud-relevant variables (for specific regions of the globe and specific cloud regimes)

Despite their pivotal role in the radiation budget and the hydrological cycle, clouds have proven notoriously challenging to simulate with global climate models. The diagnostic investigates the relationship between modelled integral cloud properties such as total cloud cover, cloud water path (sum of cloud ice and cloud liquid) and cloud ice water path and climate-relevant quantities such as long- and shortwave TOA cloud radiative effects and precipitation that can then be compared to observations. In addition, the relation between three-dimensional cloud ice water content and three-dimensional air temperature is investigated and compared to reference data.

820 C4 Earth System Realm

4.1 Equilibrium climate sensitivity (ECS)

ECS is defined as the change in global mean near-surface air temperature that results from an instantaneous doubling of the atmospheric CO₂ concentration after the climate system has reached its new equilibrium. To avoid having to simulate thousands of model years until the system has reached equilibrium, ECS is usually approximated with the “effective climate sensitivity” following a method by Gregory et al. (2004), which estimates ECS as the *x*-intercept of a linear regression of the change in global mean TOA net radiation flux against the change in global mean near-surface air temperature (see Schlund et al. (2020) for details). Even though ECS is an idealized metric, it is still relevant for policymakers and scientists since it measures the equilibrium warming response of the climate system to CO₂ forcing.



4.2 Transient climate response (TCR)

830 TCR is defined as the global mean near-surface air temperature change at the time of CO₂ doubling in a simulation where the atmospheric CO₂ concentration is increased by 1% per year from the preindustrial level (Gregory et al., 2009). In practice, it is calculated by averaging the global mean near-surface air temperature change over 20-year period centered around the time of CO₂ doubling (year 70) in the 1pctCO₂ increase simulation (see Meehl et al. (2020) for details). Similar to ECS, it is an idealized metric that describes the warming response of the climate system to CO₂ forcing, but unlike ECS, it does not assume
835 radiative equilibrium of the system.

4.3 Transient climate response to cumulative emissions of carbon dioxide (TCRE)

Unlike ECS and TCR, TCRE describes the warming response of the Earth system to cumulative emissions of CO₂ instead of atmospheric CO₂ concentrations and thus takes carbon cycle feedbacks into account (Gregory et al., 2009). Following Sanderson et al. (2024), TCRE is estimated from an experiment with constant CO₂ emissions of 10 Pg C per year (“esm-
840 flat10”) as the global mean near-surface air temperature change observed after the emission of 1000 Pg C of CO₂ (i.e., at year 100) averaged over a 20-year period (i.e., years 90–110). TCRE is one of the most important policy-relevant climate change metrics since it directly links global warming and CO₂ emissions in a linear relation, which allows a straightforward estimation of carbon budgets that remain to reach specific warming targets.

4.4 Zero Emissions Commitment (ZEC)

845 The Zero Emissions Commitment (ZEC) quantifies the change in global mean temperature expected to occur after net carbon dioxide (CO₂) emissions cease (MacDougall et al., 2020). ZEC is therefore important to consider when estimating the remaining carbon budget. Calculation of ZEC requires dedicated simulations with anthropogenic carbon emissions set to zero, branching off a base simulation. For CMIP7 fast track, the base simulation will be “esm-flat10”, while the dedicated simulation will be called “esm-flat10-zec” (Sanderson et al., 2024).

850 4.5 Historical changes in climate variables (time series, trends)

To assess the ability of climate models to look into the future, it is important to evaluate them on the changes, which have been observed in recent decades, where observations and reanalysis data provide a reference. In addition to the global mean, time series and trends of key climate variable like temperature, pressure, wind and humidity are compared regionally, using, e.g., the IPCC climate reference regions for subcontinental analysis of climate model data (Iturbide et al., 2020). A regional analysis is
855 important to assure climatic consistency and the representation of regional features within global model data as well as climate model data with regional focus, e.g., from the Coordinated Regional climate Downscaling Experiment (CORDEX; Diez-Sierra et al., 2022).



C5 Impacts & Adaptation Realm

5.1 High amplitude Rossby waves

860 High-amplitude Rossby waves are significant meanders in high-altitude winds that are linked to extreme weather events, particularly in the Northern Hemisphere. These waves are characterized by their large size and persistence, often leading to prolonged periods of extreme temperature or precipitation (Fei and White, 2023). Capturing these anomalies correctly in models is challenging, so evaluating Rossby waves is important to ensure that key mechanistic process representations are correct within models. This metric was removed from the list of diagnostics to be evaluated in the initial Assessment Fast
865 Track version of the REF because the initial version will use only monthly model output, while evaluating this diagnostic requires high frequency output. In addition, none of the existing packages evaluates high amplitude Rossby waves and correct implementation of relevant metrics will require more time than was available to produce an operating version of the REF. This diagnostic will be prioritized for implementation in a future version of the REF.

5.2 Internal variability or ensemble spread for individual models (precipitation and surface temperature)

870 The ensemble spread of individual models provides a range of outcomes representing modeled fluctuations in the Earth system due to the chaotic dynamic of individual components. This intrinsic or internal variability is a source of uncertainty in climate model projections especially in the near future and is important for impact risk assessment and adaptation (Mankin et al., 2020). Several of the REF diagnostics have the capacity to capture this spread across domains and climate variables and so for this initial release of the REF we chose to assess internal variability through individual diagnostics rather than a separate
875 one. A separate diagnostic for internal model variability will be evaluated for possible implementation in a future version of the REF. The effect of internal variability on the error metric results for historical climate model simulations (Brands, 2022a; Brands et al., 2023) is another analysis dimension that can be potentially implemented in future versions of the REF.

5.3 Evaluation of key climate variables at global warming levels

Global warming level (GWL) exceedance years or the years at which the global mean surface temperature warms over specific
880 values is an important marker of climate change. The exceedance years are calculated as the 21-year mean anomaly of the area averaged global mean surface temperature with respect to the pre-industrial control mean temperature as described in Swaminathan et al. (2022). Assessing global values of key climate variables such as temperature and precipitation at specific GWLs can tell us how different regions are affected by climate change and whether these changes scale linearly with increases in temperature thereby providing important information for adaptation and mitigation efforts. Evaluating climate at GWLs is
885 also relevant for policy, for instance the UNFCCC Paris agreement's central aim is to keep warming below 2°C and if possible below 1.5°C.



5.4 Climate drivers for fire (fire burnt area, fire weather and fuel continuity)

Many climate models do not explicitly simulate fire processes. As a result, assessing fire risk and its impacts often relies on Fire Danger or Fire Weather Indices, which use meteorological information derived from climate models. These indices are traditionally based on daily or sub-daily data. This diagnostic constructs testable fire weather and fire fuel indices specifically designed for evaluating monthly climate model output. To achieve this, we use a Bayesian inference framework, which represents and optimizes different controls of burnt area, to generate fire weather and fire vegetation indicators from monthly observed or reanalysis meteorological and vegetation data. Assessment of meteorological and biological conditions that drive fire would aid the development of these models. Developing fire weather indices that directly target burnt area, rather than relying solely on traditional fire danger formulations, presents a more robust alternative for climate model assessment.

Author contributions. The Model Benchmarking Task Team co-leads conceptualised the structure of the paper and led the writing of the draft. Authors responsible for the Rapid Evaluation Framework conceptualisation: FH and BH. Funding acquisition: FH, BH, BT, and EOR. Methodology: FH and BH. Reference dataset coordination: DH. Data provenance and citation protocol and coding: M. Stockhause, DH, J. Lewis, and BA. Project management: BT, EOR, and ED. Writing and original draft preparation: All authors. Coding and implementation of the REF: J. Lewis, MP, BA, NC, J. Lee, and MX. Development and revision of diagnostics: NC, M. Sreeush. Writing, review, and editing: All authors.

Competing interests. At least one of the (co-)authors is a member of the editorial board of Geoscientific Model Development. The authors have no other competing interests to declare.

Acknowledgements. The work developing the concept of the Rapid Evaluation Framework was coordinated and led by Forrest Hoffman (ORNL) and Birgit Hassler (DLR), co-leads of the Model Benchmarking Task Team. The work overseeing delivery of the Assessment Fast Track Rapid Evaluation Framework is being coordinated by Forrest Hoffman (ORNL), Birgit Hassler (DLR) and Ranjini Swaminathan (NCEO, UoR) as co-leads of the Model Benchmarking Task Team. Birgit Hassler (DLR) led the final selection process of the diagnostics required for the Assessment Fast Track. The REF implementation team for the Assessment Fast Track REF is led by Jared Lewis (Climate Resource) and is supported by expertise from the following software packages and organisations: Climate Resource, eScience Center, ESM-ValTool, ILAMB, IOMB, PMP, and CMEC. We wish to acknowledge and thank all the diagnostic recipe and code developers within those packages. The two ESGF nodes committed to hosting, indexing and replication of the REF are ORNL and CEDA. CEDA is also archiving the Beta reference dataset collection. Quality assurance of reference datasets has been conducted through dataset proposal review by the following members of the obs4MIPs Steering Panel; Peter Gleckler, Simon Pinnock, Greg Elsaesser, Alison Waterfall, Birgit Hassler, and Kate Willett. We also acknowledge the ESGF architecture developers who have directly worked with the REF delivery team. This includes Sasha Ames, LLNL, Lee Liming, Globus, Steve Turoscy, Globus and Dave Poulter CEDA. The project administration is provided by the CMIP IPO, which is hosted by the European Space Agency, with staff provided on contract by HE Space Operations Ltd. All the figures have been



produced by and/or commissioned by the CMIP International Project Office and are under a Creative Commons Attribution 4.0 International licence. We gratefully acknowledge the valuable feedback provided during development of the REF by Maureen Wanzala (WCRP), ESMO Scientific Steering Group and International Project Office, WGCM Infrastructure Panel, CMIP Panel, obs4MIPs Steering Panel, pre-Alpha
920 tester modelling centres; UK Met Office, ISAC-CNR, CCMC, and AS-RCEC as well as all the participants that participated in the stakeholder surveys, drop-ins and the hackathon.

Financial support. The work developing the Assessment Fast Track Rapid Evaluation Framework has been made possible by funding from the European Space Agency and the US Department of Energy. The work of Forrest M. Hoffman, Nathan Collier, and Min Xu was supported by the Reducing Uncertainties in Biogeochemical Interactions through Synthesis and Computation (RUBISCO) Science Focus Area, which is
925 sponsored by the Regional and Global Model Analysis (RGMA) activity of the Earth & Environmental Systems Modeling (EESM) Program in the Earth and Environmental Systems Sciences Division (EESDD) of the Office of Biological and Environmental Research (BER) in the US Department of Energy Office of Science. The Earth System Grid Federation (ESGF) is an international consortium of individually funded data provider institutions; the ESGF2-US Project in the United States of America is sponsored by the Data Management Program in EESDD of BER in the US Department of Energy Office of Science, and the ESGF activity in the United Kingdom is supported by
930 the Centre for Environmental Data Analysis (CEDA), which is sponsored by the Science and Technology Facilities Council (STFC) and the National Environment Research Council (NERC). Oak Ridge National Laboratory (ORNL) is managed by UT-Battelle, LLC, for the US Department of Energy under Contract No. DE-AC05-00OR22725. The work of Birgit Hassler, Lisa Bock, and Manuel Schlund was supported by the European Union's Horizon 2020 research and innovation programme under Grant Agreement No. 101003536 (ESM2025 – Earth System Models for the Future). ESMValTool diagnostic development was performed using resources of the Deutsches Klimarechenzentrum
935 (DKRZ) granted by its Scientific Steering Committee (WLA) under project no. BD0854. The work of Ranjini Swaminathan was supported by UKRI-NERC TerraFIRMA (NE/W004895/1). The work of Jiwoo Lee and Paul Ullrich was supported by the Program for Climate Model Diagnosis and Intercomparison (PCMDI), which is sponsored by the RGMA. Lawrence Livermore National Laboratory (LLNL) is managed by Lawrence Livermore National Security, LLC (LLNS), for the US Department of Energy's National Nuclear Security Administration (NNSA) under contract number DE-AC52-07NA27344. The work of Lisa Bock and Axel Lauer (AL) was supported by the ESA Climate
940 Change Initiative Climate Model User Group (ESA CCI CMUG) under contract 4000125156/18/I-NB. The work of Bettina K. Gier and Katja Weigel was funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) through the Gottfried Wilhelm Leibniz Prize awarded to Veronika Eyring (Reference number EY 22/2-1) and the project S1 ("Diagnosis and Metrics in Climate Models") of the Collaborative Research Centre TRR 181 "Energy Transfers in Atmosphere and Ocean" (Grant No. 274762653). Axel Lauer also received support by the ESA CCI Ozone project (Ozone_cci phase 3) under contract number 4000126562/19/I-NB. The work of Manuel
945 Schlund was also supported by the BMBF under CAP7 project, Grant Agreement No. 01LP2401C. The work of Mohanan G. Sreesh was supported by the European Union's Horizon Europe research and innovation program under Grant 101083922 (OceanICU Improving Carbon Understanding) and Alfred Wegener Institute (AWI) PROCEED short term Research Grant 2025. The work reflects only the authors' views; the European Commission and their executive agency are not responsible for any use that may be made. Ed Blockley was supported by the Met Office Hadley Centre Climate Programme funded by the Department for Science, Innovation and Technology. Valerio Lembo received
950 funding from the European Union's Horizon Europe research and innovation program Grant No. 101081193 (OptimESM) and acknowledges financial support from ICSC – Centro Nazionale di Ricerca in High Performance Computing, Big Data and Quantum Computing, funded by European Union – NextGenerationEU. The work of Jianhua Lu was supported by National Natural Science Foundation of China (Grant

<https://doi.org/10.5194/egusphere-2025-2685>

Preprint. Discussion started: 11 July 2025

© Author(s) 2025. CC BY 4.0 License.



No. 42442507). The work of Brian Medeiros was supported by the RGMA under Award Number DE-SC0022070 and National Science Foundation (NSF) IA 1947282; this work was also supported by the National Center for Atmospheric Research (NCAR), which is a major facility sponsored by the NSF under Cooperative Agreement No. 1852977. The work of Jeremy Walton was supported by the Natural Environment Research Council under the TerraFIRMA project (Grant reference NE/W004895/1) and the Met Office Hadley Centre Climate Programme funded by the Department for Science, Innovation and Technology.



References

- Adler, R., Wang, J.-J., Sapiano, M., Huffman, G., Chiu, L., Xie, P.-P., Ferraro, R., Schneider, U., Becker, A., Bolvin, D., Nelkin, E., Gu, G., and NOAA CDR Program: Global Precipitation Climatology Project (GPCP) Climate Data Record (CDR), Version 2.3 (Monthly), <https://doi.org/10.7289/V56971M6>, 2017.
- Adler, R. F., Huffman, G. J., Chang, A., Ferraro, R., Xie, P.-P., Janowiak, J., Rudolf, B., Schneider, U., Curtis, S., Bolvin, D., Gruber, A., Susskind, J., Arkin, P., and Nelkin, E.: The Version-2 Global Precipitation Climatology Project (GPCP) Monthly Precipitation Analysis (1979–Present), *Journal of Hydrometeorology*, 4, 1147–1167, [https://doi.org/10.1175/1525-7541\(2003\)004<1147:TVGPCP>2.0.CO;2](https://doi.org/10.1175/1525-7541(2003)004<1147:TVGPCP>2.0.CO;2), 2003.
- Adler, R. F., Sapiano, M. R. P., Huffman, G. J., Wang, J.-J., Gu, G., Bolvin, D., Chiu, L., Schneider, U., Becker, A., Nelkin, E., Xie, P., Ferraro, R., and Shin, D.-B.: The Global Precipitation Climatology Project (GPCP) Monthly Analysis (New Version 2.3) and a Review of 2017 Global Precipitation, *Atmosphere*, 9, 138, <https://doi.org/10.3390/atmos9040138>, 2018.
- Agarwal, D., Ayliffe, J., Buck, J. J. H., Damerow, J., Parton, G., Stall, S., Stockhause, M., and Wyborn, L.: Complex Citation Working Group Recommendation, <https://doi.org/10.15497/RDA00130>, report produced by the Research Data Alliance (RDA), 2025.
- Alemohammad, S. H., Fang, B., Konings, A. G., Aires, F., Green, J. K., Kolassa, J., Miralles, D., Prigent, C., and Gentine, P.: Water, Energy, and Carbon with Artificial Neural Networks (WECANN): A Statistically Based Estimate of Global Surface Turbulent Fluxes and Gross Primary Productivity Using Solar-induced Fluorescence, *Biogeosciences*, 14, 4101–4124, <https://doi.org/10.5194/bg-14-4101-2017>, 2017.
- Anav, A., Friedlingstein, P., Beer, C., Ciais, P., Harper, A., Jones, C., Murray-Tortarolo, G., Papale, D., Parazoo, N. C., Peylin, P., Piao, S., Sitch, S., Viovy, N., Wiltshire, A., and Zhao, M.: Spatiotemporal Patterns of Terrestrial Gross Primary Production: A Review, *Reviews of Geophysics*, 53, 785–818, <https://doi.org/https://doi.org/10.1002/2015RG000483>, 2015.
- Arora, V. K., Katavouta, A., Williams, R. G., Jones, C. D., Brovkin, V., Friedlingstein, P., Schwinger, J., Bopp, L., Boucher, O., Cadule, P., Chamberlain, M. A., Christian, J. R., Delire, C., Fisher, R. A., Hajima, T., Ilyina, T., Joetzjer, E., Kawamiya, M., Koven, C. D., Krasting, J. P., Law, R. M., Lawrence, D. M., Lenton, A., Lindsay, K., Pongratz, J., Raddatz, T., Séférian, R., Tachiiri, K., Tjiputra, J. F., Wiltshire, A., Wu, T., and Ziehn, T.: Carbon–concentration and Carbon–climate Feedbacks in CMIP6 Models and Their Comparison to CMIP5 Models, *Biogeosciences*, 17, 4173–4222, <https://doi.org/10.5194/bg-17-4173-2020>, publisher: Copernicus GmbH, 2020.
- Bacer, S., Jomaa, F., Beaumet, J., Gallée, H., Le Bouëdec, E., Ménégoz, M., and Staquet, C.: Impact of Climate Change on Wintertime European Atmospheric Blocking, *Weather and Climate Dynamics*, 3, 377–389, <https://doi.org/10.5194/wcd-3-377-2022>, 2022.
- Baker, A. H., Hammerling, D. M., Levy, M. N., Xu, H., Dennis, J. M., Eaton, B. E., Edwards, J., Hannay, C., Mickelson, S. A., Neale, R. B., Nychka, D., Shollenberger, J., Tribbia, J., Vertenstein, M., and Williamson, D.: A new ensemble-based consistency test for the Community Earth System Model (pyCECT v1.0), *Geoscientific Model Development*, 8, 2829–2840, <https://doi.org/10.5194/gmd-8-2829-2015>, 2015.
- Barnett, T. P., Dümenil, L., Schlese, U., Roeckner, E., and Latif, M.: The Effect of Eurasian Snow Cover on Regional and Global Climate Variations, *Journal of the Atmospheric Sciences*, 46, 661–686, [https://doi.org/10.1175/1520-0469\(1989\)046<0661:TEOESC>2.0.CO;2](https://doi.org/10.1175/1520-0469(1989)046<0661:TEOESC>2.0.CO;2), 1989.
- Beadling, R., Swaminathan, R., Beucher, R., Blockley, E., Brands, S., Hassler, B., Hegedűs, D., Hoffman, F. M., Lee, J., Lewis, J., Lu, J., Malinina, E., Medeiros, B., Scoccimarro, E., Tjiputra, J., Turner, B., and Watson-Parris, D.: Observational Data for Next Generation Climate Model Evaluation: Requirements, Considerations and Best Practices, *Bulletin of the American Meteorological Society*, <https://office.wcrp-cmip.org/sh/29JnZXWyjb>, in preparation.



- Boé, J.: Interdependency in Multimodel Climate Projections: Component Replication and Result Similarity, *Geophysical Research Letters*, 45, 2771–2779, <https://doi.org/10.1002/2017GL076829>, 2018.
- 995
- Bokhorst, S., Pedersen, S. H., Brucker, L., Anisimov, O., Bjerke, J. W., Brown, R. D., Ehrich, D., Essery, R. L. H., Heilig, A., Ingvander, S., Johansson, C., Johansson, M., Jónsdóttir, I. S., Inga, N., Luojus, K., Macelloni, G., Mariash, H., McLennan, D., Rosqvist, G. N., Sato, A., Savela, H., Schneebeli, M., Sokolov, A., Sokratov, S. A., Terzago, S., Vikhamar-Schuler, D., Williamson, S., Qiu, Y., and Callaghan, T. V.: Changing Arctic Snow Cover: A Review of Recent Developments and Assessment of Future Needs for Observations, Modelling, and Impacts, *Ambio*, 45, 516–537, <https://doi.org/10.1007/s13280-016-0770-0>, 2016.
- 1000
- Boucher, O., Servonnat, J., Albright, A. L., Aumont, O., Balkanski, Y., Bastrikov, V., Bekki, S., Bonnet, R., Bony, S., Bopp, L., Braconnot, P., Brockmann, P., Cadule, P., Caubel, A., Cheruy, F., Codron, F., Cozic, A., Cugnet, D., D’Andrea, F., Davini, P., de Lavergne, C., Denvil, S., Deshayes, J., Devilliers, M., Ducharne, A., Dufresne, J.-L., Dupont, E., Éthé, C., Fairhead, L., Falletti, L., Flavoni, S., Foujols, M.-A., Gardoll, S., Gastineau, G., Ghattas, J., Grandpeix, J.-Y., Guenet, B., Guez, Lionel, E., Guilyardi, E., Guimberteau, M., Hauglustaine, D., Hourdin, F., Idelkadi, A., Joussaume, S., Kageyama, M., Khodri, M., Krinner, G., Lebas, N., Levvasseur, G., Lévy, C., Li, L., Lott, F., Lurton, T., Luysaert, S., Madec, G., Madeleine, J.-B., Maignan, F., Marchand, M., Marti, O., Mellul, L., Meurdesoif, Y., Mignot, J., Musat, I., Ottlé, C., Peylin, P., Planton, Y., Polcher, J., Rio, C., Rochetin, N., Rousset, C., Sepulchre, P., Sima, A., Swingedouw, D., Thiéblemont, R., Traore, A. K., Vancoppenolle, M., Vial, J., Vialard, J., Viovy, N., and Vuichard, N.: Presentation and Evaluation of the IPSL-CM6A-LR Climate Model, *Journal of Advances in Modeling Earth Systems*, 12, e2019MS002 010, <https://doi.org/10.1029/2019MS002010>, 2020.
- 1010
- Brands, S.: A Circulation-based Performance Atlas of the CMIP5 and 6 Models for Regional Climate Studies in the Northern Hemisphere Mid-to-high Latitudes, *Geoscientific Model Development*, 15, 1375–1411, <https://doi.org/10.5194/gmd-15-1375-2022>, 2022a.
- Brands, S.: Common Error Patterns in the Regional Atmospheric Circulation Simulated by the CMIP Multi-Model Ensemble, *Geophysical Research Letters*, 49, e2022GL101 446, <https://doi.org/10.1029/2022GL101446>, 2022b.
- Brands, S.: pyLamb - A Climate Model Verification Tool Based on Lamb Weather Types, <https://doi.org/10.5281/zenodo.15346363>, 2025.
- 1015
- Brands, S., Fernández-Granja, J. A., Bedia, J., Casanueva, A., and Fernández, J.: A Global Climate Model Performance Atlas for the Southern Hemisphere Extratropics Based on Regional Atmospheric Circulation Patterns, *Geophysical Research Letters*, 50, e2023GL103 531, <https://doi.org/10.1029/2023GL103531>, 2023.
- Burrows, S. M., Maltrud, M., Yang, X., Zhu, Q., Jeffery, N., Shi, X., Ricciuto, D. M., Wang, S., Bisht, G., Tang, J., Wolfe, J., Harrop, B. E., Singh, B., Brent, L., Baldwin, S., Zhou, T., Cameron-Smith, P., Keen, N., Collier, N., Xu, M., Hunke, E. C., Elliott, S. M., Turner, A. K., Li, H.-Y., Wang, H., Golaz, J.-C., Bond-Lamberty, B., Hoffman, F. M., Riley, W. J., Thornton, P. E., Calvin, K., and Leung, L. R.: The DOE E3SM v1.1 Biogeochemistry Configuration: Description and Simulated Ecosystem-Climate Responses to Historical Changes in Forcing, *Journal of Advances in Modeling Earth Systems*, 12, e2019MS001 766, <https://doi.org/10.1029/2019MS001766>, 2020.
- 1020
- Cao, S., Li, M., Zhu, Z., Wang, Z., Zha, J., Zhao, W., Duanmu, Z., Chen, J., Zheng, Y., Chen, Y., Myneni, R. B., and Piao, S.: Spatiotemporally Consistent Global Dataset of the GIMMS Leaf Area Index (GIMMS LAI4g) from 1982 to 2020 (V1.2), <https://doi.org/10.5281/zenodo.8281930>, data set, 2023a.
- 1025
- Cao, S., Li, M., Zhu, Z., Wang, Z., Zha, J., Zhao, W., Duanmu, Z., Chen, J., Zheng, Y., Chen, Y., Myneni, R. B., and Piao, S.: Spatiotemporally Consistent Global Dataset of the GIMMS Leaf Area Index (GIMMS LAI4g) from 1982 to 2020, *Earth System Science Data*, 15, 4877–4899, <https://doi.org/10.5194/essd-15-4877-2023>, 2023b.
- Chen, X. and Wallace, J. M.: ENSO-Like Variability: 1900–2013, *Journal of Climate*, 28, 9623–9641, [https://doi.org/10.1175/JCLI-D-15-](https://doi.org/10.1175/JCLI-D-15-1030)
- 1030 0322.1, 2015.

Chen, Y., Hall, J., Van Wees, D., Andela, N., Hantson, S., Giglio, L., Van Der Werf, G. R., Morton, D. C., and Randerson, J. T.: Global Fire Emissions Database (GFED5) Burned Area, <https://doi.org/10.5281/ZENODO.7668424>, data set, 2023a.

1035 Chen, Y., Hall, J., van Wees, D., Andela, N., Hantson, S., Giglio, L., van der Werf, G. R., Morton, D. C., and Randerson, J. T.: Multi-decadal Trends and Variability in Burned Area from the Fifth Version of the Global Fire Emissions Database (GFED5), *Earth System Science Data*, 15, 5227–5259, <https://doi.org/10.5194/essd-15-5227-2023>, 2023b.

Cheng, L., Pan, Y., Tan, Z., Zheng, H., Zhu, Y., Wei, W., Du, J., Yuan, H., Li, G., Ye, H., Gouretski, V., Li, Y., Trenberth, K. E., Abraham, J., Jin, Y., Reseghetti, F., Lin, X., Zhang, B., Chen, G., Mann, M. E., and Zhu, J.: IAPv4 Ocean Temperature and Ocean Heat Content Gridded Dataset, *Earth System Science Data*, 16, 3517–3546, <https://doi.org/10.5194/essd-16-3517-2024>, 2024.

1040 Cinquini, L., Crichton, D., Mattmann, C., Harney, J., Shipman, G., Wang, F., Ananthakrishnan, R., Miller, N., Denvil, S., Morgan, M., Pobre, Z., Bell, G. M., Doutriaux, C., Drach, R., Williams, D., Kershaw, P., Pascoe, S., Gonzalez, E., Fiore, S., and Schweitzer, R.: The Earth System Grid Federation: An Open Infrastructure for Access to Distributed Geospatial Data, *Future Generation Computer Systems*, 36, 400–417, <https://doi.org/10.1016/j.future.2013.07.002>, 2014.

1045 Claverie, M., Vermote, E., and NOAA CDR Program: NOAA Climate Data Record (CDR) of Leaf Area Index (LAI) and Fraction of Absorbed Photosynthetically Active Radiation (FAPAR), Version 4 (Version Superseded), <https://doi.org/10.7289/V5M043BX>, published by NOAA National Centers for Environmental Information, 2014.

Claverie, M., Vermote, E., Justice, C., Csizsar, I., Myneni, R., Baret, F., Masuoka, E., Wolfe, R., Ray, J. P., and NOAA CDR Program: NOAA Climate Data Record (CDR) of VIIRS Leaf Area Index (LAI) and Fraction of Absorbed Photosynthetically Active Radiation (FAPAR), Version 1, <https://doi.org/10.25921/9X3M-0E02>, published by NOAA National Centers for Environmental Information, 2024.

1050 CMIP Model Benchmarking Task Team: CMIP7 Assessment Fast Track Diagnostics List for the Rapid Evaluation Framework, <https://doi.org/10.5281/ZENODO.14284375>, 2024.

Cohen, J. and Entekhabi, D.: Eurasian Snow Cover Variability and Northern Hemisphere Climate Predictability, *Geophysical Research Letters*, 26, 345–348, <https://doi.org/10.1029/1998GL900321>, 1999.

1055 Cohen, J., Screen, J. A., Furtado, J. C., Barlow, M., Whittleston, D., Coumou, D., Francis, J., Dethloff, K., Entekhabi, D., Overland, J., and Jones, J.: Recent Arctic Amplification and Extreme Mid-latitude Weather, *Nature Geoscience*, 7, 627–637, <https://doi.org/10.1038/ngeo2234>, 2014.

Coldewey-Egbers, M., Loyola, D. G., Latter, B., Siddans, R., Kerridge, B., Hubert, D., van Roozendaal, M., and Eisinger, M.: The novel GOME-type Ozone Profile Essential Climate Variable (GOP-ECV) Data Record Covering the Past 26 Years, *Atmospheric Measurement Techniques Discussions*, 2025, 1–31, <https://doi.org/10.5194/amt-2024-196>, 2025.

1060 Collier, N., Hoffman, F. M., Lawrence, D. M., Keppel-Aleks, G., Koven, C. D., Riley, W. J., Mu, M., and Randerson, J. T.: The International Land Model Benchmarking (ILAMB) System: Design, Theory, and Implementation, *Journal of Advances in Modeling Earth Systems*, 10, 2731–2754, <https://doi.org/10.1029/2018MS001354>, 2018.

Copernicus Climate Data Store: Ozone Monthly Gridded Data from 1970 to Present Derived from Satellite Observations, <https://doi.org/10.24381/CDS.4EBFE4EB>, published by European Center for Medium-range Weather Forecast, 2020.

1065 Cucchi, M., Weedon, G. P., Amici, A., Bellouin, N., Lange, S., Müller Schmied, H., Hersbach, H., and Buontempo, C.: WFDE5: Bias-adjusted ERA5 Reanalysis Data for Impact Studies, *Earth System Science Data*, 12, 2097–2120, <https://doi.org/10.5194/essd-12-2097-2020>, 2020.

Dai, A.: Historical and Future Changes in Streamflow and Continental Runoff, chap. 2, pp. 17–37, American Geophysical Union (AGU), ISBN 9781118971772, <https://doi.org/10.1002/9781118971772.ch2>, 2016.

Dai, A.: Dai and Trenberth Global River Flow and Continental Discharge Dataset, <https://doi.org/10.5065/D6V69H1T>, 2017.



- 1070 Dai, A.: Hydroclimatic Trends During 1950–2018 Over Global Land, *Climate Dynamics*, 56, 4027–4049, <https://doi.org/10.1007/s00382-021-05684-1>, 2021.
- Dai, A. and Trenberth, K. E.: Estimates of Freshwater Discharge from Continents: Latitudinal and Seasonal Variations, *Journal of Hydrometeorology*, 3, 660–687, [https://doi.org/10.1175/1525-7541\(2002\)003<0660:EOFDFC>2.0.CO;2](https://doi.org/10.1175/1525-7541(2002)003<0660:EOFDFC>2.0.CO;2), 2002.
- Dai, A., Qian, T., Trenberth, K. E., and Milliman, J. D.: Changes in Continental Freshwater Discharge from 1948 to 2004, *Journal of Climate*, 22, 2773–2792, <https://doi.org/10.1175/2008JCLI2592.1>, 2009.
- 1075 Danabasoglu, G., Lamarque, J.-F., Bacmeister, J., Bailey, D. A., DuVivier, A. K., Edwards, J., Emmons, L. K., Fasullo, J., Garcia, R., Gettelman, A., Hannay, C., Holland, M. M., Large, W. G., Lauritzen, P. H., Lawrence, D. M., Lenaerts, J. T. M., Lindsay, K., Lipscomb, W. H., Mills, M. J., Neale, R., Oleson, K. W., Otto-Bliessner, B., Phillips, A. S., Sacks, W., Tilmes, S., van Kampenhout, L., Versteinen, M., Bertini, A., Dennis, J., Deser, C., Fischer, C., Fox-Kemper, B., Kay, J. E., Kinnison, D., Kushner, P. J., Larson, V. E., Long, M. C., Mickelson, S., Moore, J. K., Nienhouse, E., Polvani, L., Rasch, P. J., and Strand, W. G.: The Community Earth System Model Version 2 (CESM2), *Journal of Advances in Modeling Earth Systems*, 12, e2019MS001916, <https://doi.org/10.1029/2019MS001916>, 2020.
- 1080 Davini, P.: MiLES - Mid Latitude Evaluation System, <https://doi.org/10.5281/zenodo.2578139>, 2019.
- Davini, P. and D’Andrea, F.: Northern Hemisphere Atmospheric Blocking Representation in Global Climate Models: Twenty Years of Improvements?, *Journal of Climate*, 29, 8823–8840, <https://doi.org/10.1175/JCLI-D-16-0242.1>, 2016.
- Davis, E., Taylor, K. E., Adloff, F., Gregory, J., Lawrence, B., and Lee, D.: Supporting Open Science with the CF Metadata Conventions for NetCDF, <https://doi.org/10.5281/zenodo.15015065>, slides for the presentation given on 11 December 2024 at the AGU Annual Meeting to the session for the AGU Open Science Recognition Prize, 2024.
- 1085 Deser, C., Phillips, A., Bourdette, V., and Teng, H.: Uncertainty in Climate Change Projections: The Role of Internal Variability, *Climate Dynamics*, 38, 527–546, <https://doi.org/10.1007/s00382-010-0977-x>, 2012.
- Deser, C., Lehner, F., Rodgers, K. B., Ault, T., Delworth, T. L., DiNezio, P. N., Fiore, A., Frankignoul, C., Fyfe, J. C., Horton, D. E., Kay, J. E., Knutti, R., Lovenduski, N. S., Marotzke, J., McKinnon, K. A., Minobe, S., Randerson, J., Screen, J. A., Simpson, I. R., and Ting, M.: Insights from Earth System Model Initial-condition Large Ensembles and Future Prospects, *Nature Climate Change*, 10, 277–286, <https://doi.org/10.1038/s41558-020-0731-2>, 2020.
- 1090 Diez-Sierra, J., Iturbide, M., Gutiérrez, J. M., Fernández, J., Milovac, J., Cofiño, A. S., Cimadevilla, E., Nikulin, G., Levvasseur, G., Kjellström, E., Bülow, K., Horányi, A., Brookshaw, A., García-Díez, M., Pérez, A., Baño-Medina, J., Ahrens, B., Alias, A., Ashfaq, M., Bukovsky, M., Buonomo, E., Caluwaerts, S., Chou, S. C., Christensen, O. B., Ciarlò, J. M., Coppola, E., Corre, L., Demory, M.-E., Djurdjevic, V., Evans, J. P., Fealy, R., Feldmann, H., Jacob, D., Jayanarayanan, S., Katzfey, J., Keuler, K., Kittel, C., Kurnaz, M. L., Laprise, R., Lionello, P., McGinnis, S., Mercogliano, P., Nabat, P., Öno, B., Ozturk, T., Panitz, H.-J., Paquin, D., Pieczka, I., Raffaele, F., Remedio, A. R., Scinocca, J., Sevault, F., Somot, S., Steger, C., Tangang, F., Teichmann, C., Termonia, P., Thatcher, M., Torma, C., Meijgaard, E. v., Vautard, R., Warrach-Sagi, K., Winger, K., and Zittis, G.: The Worldwide C3S CORDEX Grand Ensemble: A Major Contribution to Assess Regional Climate Change in the IPCC AR6 Atlas, *Bulletin of the American Meteorological Society*, 103, E2804–E2826, <https://doi.org/10.1175/BAMS-D-22-0111.1>, 2022.
- 1100 DiMiceli, C., Carroll, M., Sohlberg, R., Kim, D.-H., Kelly, M., and Townshend, J.: MOD44B MODIS/Terra Vegetation Continuous Fields Yearly L3 Global 250m SIN Grid V006, <https://doi.org/10.5067/MODIS/MOD44B.006>, 2015.
- 1105 Dingley, B., Anstey, J. A., Abalos, M., Abraham, C., Bergman, T., Bock, L., Hassler, B., Kramer, R. J., Luo, F., O’Connor, F. M., Šácha, P., Simpson, I. R., Wilcox, L. J., and Zelinka, M. D.: Atmosphere Theme Data Request for CMIP7, Geoscientific Model Development, <https://office.wcrp-cmip.org/sh/29JnZXWyjb>, in preparation.



- Dirmeyer, P. A., Gao, X., Zhao, M., Guo, Z., Oki, T., and Hanasaki, N.: GSWP-2: Multimodel Analysis and Implications for Our Perception of the Land Surface, *Bulletin of the American Meteorological Society*, 87, 1381–1398, <https://doi.org/10.1175/BAMS-87-10-1381>, 2006.
- 1110 Dolores-Tesillos, E., Martius, O., and Quinting, J.: On the Role of Moist and Dry Processes in Atmospheric Blocking Biases in the Euro-Atlantic Region in CMIP6, *Weather and Climate Dynamics*, 6, 471–487, <https://doi.org/10.5194/wcd-6-471-2025>, 2025.
- Dorrington, J., Strommen, K., and Fabiano, F.: Quantifying Climate Model Representation of the Wintertime Euro-Atlantic Circulation Using Geopotential-Jet Regimes, *Weather and Climate Dynamics*, 3, 505–533, <https://doi.org/10.5194/wcd-3-505-2022>, 2022.
- Dunne, J. P., Hewitt, H. T., Arblaster, J., Bonou, F., Boucher, O., Cavazos, T., Durack, P. J., Hassler, B., Juckes, M., Miyakawa, T., Mizielinski, M., Naik, V., Nicholls, Z., O'Rourke, E., Pincus, R., Sanderson, B. M., Simpson, I. R., and Taylor, K. E.: An Evolving Coupled
1115 Model Intercomparison Project Phase 7 (CMIP7) and Fast Track in Support of Future Climate Assessment, *EGUsphere*, pp. 1–51, <https://doi.org/10.5194/egusphere-2024-3874>, 2024.
- Döscher, R., Acosta, M., Alessandri, A., Anthoni, P., Arsouze, T., Bergman, T., Bernardello, R., Boussetta, S., Caron, L.-P., Carver, G., Castrillo, M., Catalano, F., Cvijanovic, I., Davini, P., Dekker, E., Doblas-Reyes, F. J., Docquier, D., Echevarria, P., Fladrich, U., Fuentes-Franco, R., Gröger, M., v. Hardenberg, J., Hieronymus, J., Karami, M. P., Keskinen, J.-P., Koenigk, T., Makkonen, R., Massonnet, F.,
1120 Ménégos, M., Miller, P. A., Moreno-Chamarro, E., Nieradzic, L., van Noije, T., Nolan, P., O'Donnell, D., Ollinaho, P., van den Oord, G., Ortega, P., Prims, O. T., Ramos, A., Reerink, T., Rousset, C., Ruprich-Robert, Y., Le Sager, P., Schmith, T., Schrödner, R., Serva, F., Sicardi, V., Sloth Madsen, M., Smith, B., Tian, T., Tourigny, E., Uotila, P., Vancoppenolle, M., Wang, S., Wärlind, D., Willén, U., Wyser, K., Yang, S., Yepes-Arbós, X., and Zhang, Q.: The EC-Earth3 Earth system model for the Coupled Model Intercomparison Project 6, *Geoscientific Model Development*, 15, 2973–3020, <https://doi.org/10.5194/gmd-15-2973-2022>, 2022.
- 1125 Eaton, B., Gregory, J., Drach, B., Taylor, K., Hankin, S., Caron, J., Signell, R., Bentley, P., Rappa, G., Höck, H., Pamment, A., Juckes, M., Raspaud, M., Blower, J., Horne, R., Whiteaker, T., Blodgett, D., Zender, C., Lee, D., Hassell, D., Snow, A. D., Kölling, T., Allured, D., Jelenak, A., Soerensen, A. M., Gaultier, L., Herlédan, S., Manzano, F., Barring, L., Barker, C., and Bartholomew, S. L.: NetCDF Climate and Forecast (CF) Metadata Conventions, <https://doi.org/10.5281/zenodo.14275599>, report produced by the CF Community, 2024.
- Elbaum, E., Garfinkel, C. I., Adam, O., Morin, E., Rostkier-Edelstein, D., and Dayan, U.: Uncertainty in Projected Changes in Precipitation
1130 Minus Evaporation: Dominant Role of Dynamic Circulation Changes and Weak Role for Thermodynamic Changes, *Geophysical Research Letters*, 49, e2022GL097725, <https://doi.org/https://doi.org/10.1029/2022GL097725>, 2022.
- Eyring, V., Bony, S., Meehl, G. A., Senior, C. A., Stevens, B., Stouffer, R. J., and Taylor, K. E.: Overview of the Coupled Model Intercomparison Project Phase 6 (CMIP6) Experimental Design and Organization, *Geoscientific Model Development*, 9, 1937–1958, <https://doi.org/10.5194/gmd-9-1937-2016>, 2016.
- 1135 Eyring, V., Cox, P. M., Flato, G. M., Gleckler, P. J., Abramowitz, G., Caldwell, P., Collins, W. D., Gier, B. K., Hall, A. D., Hoffman, F. M., Hurtt, G. C., Jahn, A., Jones, C. D., Klein, S. A., Krasting, J. P., Kwiatkowski, L., Lorenz, R., Maloney, E., Meehl, G. A., Pendergrass, A. G., Pincus, R., Ruane, A. C., Russell, J. L., Sanderson, B. M., Santer, B. D., Sherwood, S. C., Simpson, I. R., Stouffer, R. J., and Williamson, M. S.: Taking Climate Model Evaluation to the Next Level, *Nature Climate Change*, 9, 102–110, <https://doi.org/10.1038/s41558-018-0355-y>, 2019.
- 1140 Eyring, V., Bock, L., Lauer, A., Righi, M., Schlund, M., Andela, B., Arnone, E., Bellprat, O., Brötz, B., Caron, L.-P., Carvalho, N., Cionni, I., Cortesi, N., Crezee, B., Davin, E. L., Davini, P., Debeire, K., de Mora, L., Deser, C., Docquier, D., Earnshaw, P., Ehbrecht, C., Gier, B. K., Gonzalez-Reviriego, N., Goodman, P., Hagemann, S., Hardiman, S., Hassler, B., Hunter, A., Kadow, C., Kindermann, S., Koirala, S., Koldunov, N., Lejeune, Q., Lembo, V., Lovato, T., Lucarini, V., Massonnet, F., Müller, B., Pandde, A., Pérez-Zanón, N., Phillips, A., Predoi, V., Russell, J., Sellar, A., Serva, F., Stacke, T., Swaminathan, R., Torralba, V., Vegas-Regidor, J., von Hardenberg, J., Weigel,



- 1145 K., and Zimmermann, K.: Earth System Model Evaluation Tool (ESMValTool) v2.0 – An Extended Set of Large-scale Diagnostics for Quasi-operational and Comprehensive Evaluation of Earth System Models in CMIP, *Geoscientific Model Development*, 13, 3383–3438, <https://doi.org/10.5194/gmd-13-3383-2020>, 2020.
- Fang, H., Baret, F., Plummer, S., and Schaepman-Strub, G.: An Overview of Global Leaf Area Index (LAI): Methods, Products, Validation, and Applications, *Reviews of Geophysics*, 57, 739–799, <https://doi.org/https://doi.org/10.1029/2018RG000608>, 2019.
- 1150 FAO and IIASA: Harmonized World Soil Database Version 2.0, Food and Agriculture Organization (FAO) and International Institute for Applied Systems Analysis (IIASA), ISBN 978-92-5-137499-3, <https://doi.org/10.4060/cc3823en>, 2023.
- Fasullo, J. T., Phillips, A. S., and Deser, C.: Evaluation of Leading Modes of Climate Variability in the CMIP Archives, *Journal of Climate*, 33, 5527–5545, <https://doi.org/10.1175/JCLI-D-19-1024.1>, 2020.
- Fei, C. and White, R. H.: Large-Amplitude Quasi-Stationary Rossby Wave Events in ERA5 and the CESM2: Features, Precursors, and Model Biases in Northern Hemisphere Winter, *Journal of the Atmospheric Sciences*, 80, 2075–2090, <https://doi.org/10.1175/JAS-D-22-0042.1>, 2023.
- 1155 Fernández-Granja, J. A., Bedia, J., Casanueva, A., Brands, S., and Fernández, J.: The Signature of the Main Modes of Climatic Variability as Revealed by the Jenkinson-Collison Classification over Europe, *International Journal of Climatology*, 44, 4076–4088, <https://doi.org/10.1002/joc.8569>, 2024.
- 1160 Ferraro, R., Waliser, D. E., Gleckler, P., Taylor, K. E., and Eyring, V.: Evolving Obs4MIPs to Support Phase 6 of the Coupled Model Intercomparison Project (CMIP6), *Bulletin of the American Meteorological Society*, 96, ES131–ES133, <https://doi.org/10.1175/BAMS-D-14-00216.1>, 2015.
- Fu, W., Moore, J. K., Primeau, F., Collier, N., Ogunro, O. O., Hoffman, F. M., and Randerson, J. T.: Evaluation of Ocean Biogeochemistry and Carbon Cycling in CMIP Earth System Models With the International Ocean Model Benchmarking (IOMB) Software System, *Journal of Geophysical Research: Oceans*, 127, e2022JC018965, <https://doi.org/10.1029/2022JC018965>, 2022.
- 1165 Giorgi, F. and Francisco, R.: Uncertainties in Regional Climate Change Prediction: A Regional Analysis of Ensemble Simulations with the HADCM2 Coupled AOGCM, *Climate Dynamics*, 16, 169–182, <https://doi.org/10.1007/PL00013733>, 2000.
- Giorgi, F. and Gutowski, W. J.: Regional Dynamical Downscaling and the CORDEX Initiative, *Annual Review of Environment and Resources*, 40, 467–490, <https://doi.org/10.1146/annurev-environ-102014-021217>, 2015.
- 1170 Gleckler, P., Ferraro, R., and Waliser, D.: Improving Use of Satellite Data in Evaluating Climate Models, *Eos Transactions American Geophysical Union*, 92, 172–172, <https://doi.org/10.1029/2011EO200005>, 2011.
- Gleckler, P., Taylor, K. E., Durack, P. J., Nadeau, D., Biard, J. C., Elsaesser, G., Ferraro, R., Finkensieper, S., Hassler, B., Manaster, A., Mears, C., Pinnock, S., Stevens, S., Tuma, M., Turner, B., Waterfall, A., and Willet, K. M.: Obs4MIPs Data Specifications Version 2.5 (ODS2.5), <https://doi.org/10.5281/zenodo.11500474>, report produced by the Earth System Model and Observations (ESMO), 2024.
- 1175 Gleckler, P. J., Taylor, K. E., and Doutriaux, C.: Performance Metrics for Climate Models, *Journal of Geophysical Research: Atmospheres*, 113, <https://doi.org/10.1029/2007JD008972>, 2008.
- Grams, C. M., Beerli, R., Pfenninger, S., Staffell, I., and Wernli, H.: Balancing Europe’s Wind-power Output Through Spatial Deployment Informed by Weather Regimes, *Nature Climate Change*, 7, 557–562, <https://doi.org/10.1038/nclimate3338>, 2017.
- Gregory, J. M., Ingram, W. J., Palmer, M. A., Jones, G. S., Stott, P. A., Thorpe, R. B., Lowe, J. A., Johns, T. C., and Williams, K. D.: A New Method for Diagnosing Radiative Forcing and Climate Sensitivity, *Geophysical Research Letters*, 31, <https://doi.org/10.1029/2003GL018747>, 2004.

Gregory, J. M., Jones, C. D., Cadule, P., and Friedlingstein, P.: Quantifying Carbon Cycle Feedbacks, *Journal of Climate*, 22, 5232–5250, <https://doi.org/10.1175/2009JCLI2949.1>, 2009.

1185 Hannachi, A., Finke, K., and Trendafilov, N.: Common EOFs: A Tool for Multi-model Comparison and Evaluation, *Climate Dynamics*, 60, 1689–1703, <https://doi.org/10.1007/s00382-022-06409-8>, 2023.

Hassler, B., Hoffman, F. M., Beadling, R., Blockley, E., Huang, B., Lee, J., Lembo, V., Lewis, J., Lu, J., Madaus, L., Malinina, E., Medeiros, B., Pokam, W., Scoccimarro, E., and Swaminathan, R.: Systematic Benchmarking of Climate Models: Methodologies, Applications, and New Directions, *ESS Open Archive*, <https://doi.org/10.22541/essoar.174196646.65056548/v1>, 2025.

1190 Hawkins, E. and Sutton, R.: Connecting Climate Model Projections of Global Temperature Change with the Real World, *Bulletin of the American Meteorological Society*, 97, 963–980, <https://doi.org/10.1175/BAMS-D-14-00154.1>, 2016.

Hellinger, E.: Neue Begründung der Theorie quadratischer Formen von unendlichvielen Veränderlichen, *Journal für die reine und angewandte Mathematik*, 1909, 210–271, <https://doi.org/10.1515/crll.1909.136.210>, 1909.

1195 Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz-Sabater, J., Nicolas, J., Peubey, C., Radu, R., Schepers, D., Simmons, A., Soci, C., Abdalla, S., Abellan, X., Balsamo, G., Bechtold, P., Biavati, G., Bidlot, J., Bonavita, M., De Chiara, G., Dahlgren, P., Dee, D., Diamantakis, M., Dragani, R., Flemming, J., Forbes, R., Fuentes, M., Geer, A., Haimberger, L., Healy, S., Hogan, R. J., Hólm, E., Janisková, M., Keeley, S., Laloyaux, P., Lopez, P., Lupu, C., Radnoti, G., de Rosnay, P., Rozum, I., Vamborg, F., Villaume, S., and Thépaut, J.-N.: The ERA5 Global Reanalysis, *Quarterly Journal of the Royal Meteorological Society*, 146, 1999–2049, <https://doi.org/10.1002/qj.3803>, 2020.

1200 Hobeichi, S., Abramowitz, G., Evans, J., and Beck, H. E.: Linear Optimal Runoff Aggregate (LORA): A Global Gridded Synthesis Runoff Product, *Hydrology and Earth System Sciences*, 23, 851–870, <https://doi.org/10.5194/hess-23-851-2019>, 2019.

Hoffman, F., Hassler, B., and Model Benchmarking Task Team: Rapid Evaluation Framework Overview, <https://doi.org/10.5281/zenodo.15594502>, 2024.

Hoffman, F. M., Hargrove, W. W., Erickson, D. J., and Oglesby, R. J.: Using Clustered Climate Regimes to Analyze and Compare Predictions from Fully Coupled General Circulation Models, *Earth Interactions*, 9, 1–27, <https://doi.org/10.1175/EI110.1>, 2005.

1205 Hoffman, F. M., Randerson, J. T., Arora, V. K., Bao, Q., Cadule, P., Ji, D., Jones, C. D., Kawamiya, M., Khatiwala, S., Lindsay, K., Obata, A., Shevliakova, E., Six, K. D., Tjiputra, J. F., Volodin, E. M., and Wu, T.: Causes and Implications of Persistent Atmospheric Carbon Dioxide Biases in Earth System Models, *Journal of Geophysical Research: Biogeosciences*, 119, 141–162, <https://doi.org/10.1002/2013JG002381>, 2014.

1210 Hoffman, F. M., Koven, C. D., Keppel-Aleks, G., Lawrence, D. M., Riley, W. J., Randerson, J. T., Ahlström, A., Abramowitz, G., Baldocchi, D. D., Best, M. J., Bond-Lamberty, B., De Kauwe, M. G., Denning, A. S., Desai, A. R., Eyring, V., Fisher, J. B., Fisher, R. A., Gleckler, P. J., Huang, M., Hugelius, G., Jain, A. K., Kiang, N. Y., Kim, H., Koster, R. D., Kumar, S. V., Li, H., Luo, Y., Mao, J., McDowell, N. G., Mishra, U., Moorcroft, P. R., Pau, G. S. H., Ricciuto, D. M., Schaefer, K., Schwalm, C. R., Serbin, S. P., Shevliakova, E., Slater, A. G., Tang, J., Williams, M., Xia, J., Xu, C., Joseph, R., and Koch, D.: International Land Model Benchmarking (ILAMB) 2016 Workshop Report, Tech. Rep. DOE/SC-0186, U.S. Department of Energy, Office of Science, Germantown, Maryland, USA, <https://doi.org/10.2172/1330803>, 2017.

Hori, M., Sugiura, K., Kobayashi, K., Aoki, T., Tanikawa, T., Kuchiki, K., Niwano, M., and Enomoto, H.: A 38-year (1978–2015) Northern Hemisphere Daily Snow Cover Extent Product Derived Using Consistent Objective Criteria from Satellite-borne Optical Sensors, *Remote Sensing of Environment*, 191, 402–418, <https://doi.org/10.1016/j.rse.2017.01.023>, 2017.



- Hugelius, G., Bockheim, J. G., Camill, P., Elberling, B., Grosse, G., Harden, J. W., Johnson, K., Jorgenson, T., Koven, C. D., Kuhry, P.,
1220 Michaelson, G., Mishra, U., Palmtag, J., Ping, C.-L., O'Donnell, J., Schirmer, L., Schuur, E. a. G., Sheng, Y., Smith, L. C., Strauss,
J., and Yu, Z.: A New Data Set for Estimating Organic Carbon Storage to 3 m Depth in Soils of the Northern Circumpolar Permafrost
Region, *Earth System Science Data*, 5, 393–402, <https://doi.org/10.5194/essd-5-393-2013>, 2013.
- Hurrell, J. W., Kushnir, Y., and Visbeck, M.: The North Atlantic Oscillation, *Science*, 291, 603–605, <https://doi.org/10.1126/science.1058761>,
2001.
- 1225 IPCC: Climate Change 2021 – The Physical Science Basis: Working Group I Contribution to the Sixth Assessment Re-
port of the Intergovernmental Panel on Climate Change, Cambridge University Press, 1st edn., ISBN 978-1-009-15789-6,
<https://doi.org/10.1017/9781009157896>, 2023.
- Irving, D. B., Wijffels, S., and Church, J. A.: Anthropogenic Aerosols, Greenhouse Gases, and the Uptake, Transport, and Storage of Excess
Heat in the Climate System, *Geophysical Research Letters*, 46, 4894–4903, <https://doi.org/https://doi.org/10.1029/2019GL082015>, 2019.
- 1230 Iturbide, M., Gutiérrez, J. M., Alves, L. M., Bedia, J., Cerezo-Mota, R., Gimeno, E., Góñi, A. S., Di Luca, A., Faria, S. H., Gorodet-
skaya, I. V., Hauser, M., Herrera, S., Hennessy, K., Hewitt, H. T., Jones, R. G., Krakovska, S., Manzananas, R., Martínez-Castro, D., Nar-
isma, G. T., Nurhati, I. S., Pinto, I., Seneviratne, S. I., van den Hurk, B., and Vera, C. S.: An Update of IPCC Climate Reference Regions
for Subcontinental Analysis of Climate Model Data: Definition and Aggregated Datasets, *Earth System Science Data*, 12, 2959–2970,
<https://doi.org/10.5194/essd-12-2959-2020>, 2020.
- 1235 Jukes, M., Taylor, K. E., Antonio, F., Brayshaw, D., Buontempo, C., Cao, J., Durack, P. J., Kawamiya, M., Kim, H., Lovato, T., Mackallah,
C., Mizielinski, M., Nuzzo, A., Stockhause, M., Visoni, D., Walton, J., Turner, B., O'Rourke, E., and Dingley, B.: Baseline Climate
Variables for Earth System Modelling, *EGUsphere*, pp. 1–37, <https://doi.org/10.5194/egusphere-2024-2363>, 2024.
- Keenan, T. F. and Williams, C. A.: The Terrestrial Carbon Sink, *Annual Review of Environment and Resources*, 43, 219–243,
<https://doi.org/10.1146/annurev-environ-102017-030204>, 2018.
- 1240 Keppel-Aleks, G., Randerson, J. T., Lindsay, K., Stephens, B. B., Moore, J. K., Doney, S. C., Thornton, P. E., Mahowald, N. M., Hoffman,
F. M., Sweeney, C., Tans, P. P., Wennberg, P. O., and Wofsy, S. C.: Atmospheric Carbon Dioxide Variability in the Community Earth
System Model: Evaluation and Transient Dynamics during the Twentieth and Twenty-First Centuries, *Journal of Climate*, 26, 4447–4475,
<https://doi.org/10.1175/JCLI-D-12-00589.1>, 2013.
- Key, R. M., Kozyr, A., Sabine, C. L., Lee, K., Wanninkhof, R., Bullister, J. L., Feely, R. A., Millero, F. J., Mordy, C., and Peng, T.-
1245 H.: A Global Ocean Carbon Climatology: Results from Global Data Analysis Project (GLODAP), *Global Biogeochemical Cycles*, 18,
<https://doi.org/10.1029/2004GB002247>, 2004.
- Kumar, P. B., Vialard, J., Lengaigne, M., Murty, V. S. N., and McPhaden, M. J.: TropFlux: Air-sea Fluxes for the Global Tropical Ocean -
Description and Evaluation, *Climate Dynamics*, 38, 1521–1543, <https://doi.org/10.1007/s00382-011-1115-0>, 2012.
- Lange, S.: WFDE5 Over Land Merged with ERA5 Over the Ocean (W5E5), <https://doi.org/10.5880/pik.2019.023>, 2019.
- 1250 Lange, S., Mengel, M., Treu, S., and Büchner, M.: ISIMIP3a Atmospheric Climate Input Data, <https://doi.org/10.48364/ISIMIP.982724.2>,
2022.
- Lange, S., Quesada-Chacón, D., and Büchner, M.: Secondary ISIMIP3b Bias-adjusted Atmospheric Climate Input Data,
<https://doi.org/10.48364/ISIMIP.581124.5>, 2024.
- Lauer, A., Eyring, V., Bellprat, O., Bock, L., Gier, B. K., Hunter, A., Lorenz, R., Pérez-Zanón, N., Righi, M., Schlund, M., Senftleben,
1255 D., Weigel, K., and Zechlau, S.: Earth System Model Evaluation Tool (ESMValTool) v2.0 – Diagnostics for Emergent Constraints and



- Future Projections from Earth System Models in CMIP, *Geoscientific Model Development*, 13, 4205–4228, <https://doi.org/10.5194/gmd-13-4205-2020>, 2020.
- Lauer, A., Bock, L., Hassler, B., Jöckel, P., Ruhe, L., and Schlund, M.: Monitoring and Benchmarking Earth System Model Simulations with ESMValTool v2.12.0, *Geoscientific Model Development*, 18, 1169–1188, <https://doi.org/10.5194/gmd-18-1169-2025>, 2025.
- 1260 Lauritzen, P. H., Kevlahan, N. K.-R., Toniazzo, T., Eldred, C., Dubos, T., Gassmann, A., Larson, V. E., Jablonowski, C., Guba, O., Shipway, B., Harrop, B. E., Lemarié, F., Tailleux, R., Herrington, A. R., Large, W., Rasch, P. J., Donahue, A. S., Wan, H., Conley, A., and Bacmeister, J. T.: Reconciling and Improving Formulations for Thermodynamics and Conservation Principles in Earth System Models (ESMs), *Journal of Advances in Modeling Earth Systems*, 14, e2022MS003117, <https://doi.org/10.1029/2022MS003117>, 2022.
- Lavergne, T., Sørensen, A. M., Kern, S., Tonboe, R., Notz, D., Aaboe, S., Bell, L., Dybkjær, G., Eastwood, S., Gabarro, C., Heygster, G., 1265 Killie, M. A., Brandt Kreiner, M., Lavelle, J., Saldo, R., Sandven, S., and Pedersen, L. T.: Version 2 of the EUMETSAT OSI SAF and ESA CCI Sea-Ice Concentration Climate Data Records, *The Cryosphere*, 13, 49–78, <https://doi.org/10.5194/tc-13-49-2019>, 2019.
- Lawrence, D. M., Fisher, R. A., Koven, C. D., Oleson, K. W., Swenson, S. C., Bonan, G., Collier, N., Ghimire, B., van Kampenhout, L., Kennedy, D., Kluzek, E., Lawrence, P. J., Li, F., Li, H., Lombardozi, D., Riley, W. J., Sacks, W. J., Shi, M., Vertenstein, M., Wieder, W. R., Xu, C., Ali, A. A., Badger, A. M., Bisht, G., van den Broeke, M., Brunke, M. A., Burns, S. P., Buzan, J., Clark, M., Craig, A., 1270 Dahlin, K., Drewniak, B., Fisher, J. B., Flanner, M., Fox, A. M., Gentine, P., Hoffman, F. M., Keppel-Aleks, G., Knox, R., Kumar, S., Lenaerts, J., Leung, L. R., Lipscomb, W. H., Lu, Y., Pandey, A., Pelletier, J. D., Perket, J., Randerson, J. T., Ricciuto, D. M., Sanderson, B. M., Slater, A., Subin, Z. M., Tang, J., Thomas, R. Q., Val Martin, M., and Zeng, X.: The Community Land Model Version 5: Description of New Features, Benchmarking, and Impact of Forcing Uncertainty, *Journal of Advances in Modeling Earth Systems*, 11, 4245–4287, <https://doi.org/10.1029/2018MS001583>, 2019.
- 1275 Le Bras, I. A.-A., Willis, J., and Fenty, I.: The Atlantic Meridional Overturning Circulation at 35°N From Deep Moorings, Floats, and Satellite Altimeter, *Geophysical Research Letters*, 50, e2022GL101931, <https://doi.org/https://doi.org/10.1029/2022GL101931>, 2023.
- Le Quéré, C., Andrew, R. M., Friedlingstein, P., Sitch, S., Hauck, J., Pongratz, J., Pickers, P. A., Korsbakken, J. I., Peters, G. P., Canadell, J. G., Arneeth, A., Arora, V. K., Barbero, L., Bastos, A., Bopp, L., Chevallier, F., Chini, L. P., Ciais, P., Doney, S. C., Gkritzalis, T., Goll, D. S., Harris, I., Haverd, V., Hoffman, F. M., Hoppema, M., Houghton, R. A., Hurtt, G., Ilyina, T., Jain, A. K., Johannessen, T., Jones, C. D., Kato, E., Keeling, R. F., Goldewijk, K. K., Landschützer, P., Lefèvre, N., Lienert, S., Liu, Z., Lombardozi, D., Metzl, N., Munro, D. R., Nabel, J. E. M. S., Nakaoka, S.-I., Neill, C., Olsen, A., Ono, T., Patra, P., Peregón, A., Peters, W., Peylin, P., Pfeil, B., Pierrot, D., Poulter, B., Rehder, G., Resplandy, L., Robertson, E., Rocher, M., Rödenbeck, C., Schuster, U., Schwinger, J., Séférian, R., Skjelvan, I., Steinhoff, T., Sutton, A., Tans, P. P., Tian, H., Tilbrook, B., Tubiello, F. N., van der Laan-Luijkx, I. T., van der Werf, G. R., Viovy, N., Walker, A. P., Wiltshire, A. J., Wright, R., Zaehle, S., and Zheng, B.: Global Carbon Budget 2018, *Earth System Science Data*, 10, 1285 2141–2194, <https://doi.org/10.5194/essd-10-2141-2018>, 2018.
- Lee, J.: CMIP6 Community Survey: Model Benchmarking and Evaluation Results, <https://doi.org/10.5281/zenodo.15478482>, 2024.
- Lee, J., Gleckler, P. J., Ahn, M.-S., Ordóñez, A., Ullrich, P. A., Sperber, K. R., Taylor, K. E., Planton, Y. Y., Guilyardi, E., Durack, P., Bonfils, C., Zelinka, M. D., Chao, L.-W., Dong, B., Doutriaux, C., Zhang, C., Vo, T., Boutte, J., Wehner, M. F., Pendergrass, A. G., Kim, D., Xue, Z., Wittenberg, A. T., and Krasting, J.: Systematic and Objective Evaluation of Earth System Models: PCMDI Metrics Package (PMP) 1290 Version 3, *Geoscientific Model Development*, 17, 3919–3948, <https://doi.org/10.5194/gmd-17-3919-2024>, 2025.
- Lembo, V., Folini, D., Wild, M., and Lionello, P.: Inter-hemispheric Differences in Energy Budgets and Cross-equatorial Transport Anomalies During the 20th Century, *Climate Dynamics*, 53, 115–135, <https://doi.org/10.1007/s00382-018-4572-x>, 2019.



- Lewis, J., Andela, B., Collier, N., Lee, J., Hegedűs, D., Pflüger, M., and Xu, M.: Rapid Evaluation Framework, v0.6.1, <https://doi.org/10.5281/zenodo.15534355>, 2025a.
- 1295 Lewis, J., Hoffman, F., and REF Delivery Team: REF Key Services and Interactions with External Services Including ESGF, <https://doi.org/10.5281/zenodo.15595006>, 2025b.
- Loeb, N. G., Doelling, D. R., Wang, H., Su, W., Nguyen, C., Corbett, J. G., Liang, L., Mitrescu, C., Rose, F. G., and Kato, S.: Clouds and the Earth's Radiant Energy System (CERES) Energy Balanced and Filled (EBAF) Top-of-Atmosphere (TOA) Edition-4.0 Data Product, *Journal of Climate*, 31, 895–918, <https://doi.org/10.1175/JCLI-D-17-0208.1>, 2018.
- 1300 Loeb, N. G., Wang, H., Allan, R. P., Andrews, T., Armour, K., Cole, J. N. S., Dufresne, J.-L., Forster, P., Gettelman, A., Guo, H., Mauritson, T., Ming, Y., Paynter, D., Proistosescu, C., Stuecker, M. F., Willén, U., and Wyser, K.: New Generation of Climate Models Track Recent Unprecedented Changes in Earth's Radiation Budget Observed by CERES, *Geophysical Research Letters*, 47, e2019GL086705, <https://doi.org/https://doi.org/10.1029/2019GL086705>, 2020.
- Lucarini, V., Ragone, F., and Lunkeit, F.: Predicting Climate Change Using Response Theory: Global Averages and Spatial Patterns, *Journal of Statistical Physics*, 166, 1036–1064, <https://doi.org/10.1007/s10955-016-1506-z>, 2017.
- 1305 Luo, Y. and Hoffman, F. M.: Benchmark Analysis, in: *Land Carbon Cycle Modeling: Matrix Approach, Data Assimilation, & Ecological Forecasting*, edited by Luo, Y. and Smith, B., pp. 157–162, CRC Press, Boca Raton, FL, USA, ISBN 978-1-4987-3701-2, <https://doi.org/10.1201/9780429155659-24>, 2022.
- Luo, Y. Q., Randerson, J. T., Abramowitz, G., Bacour, C., Blyth, E., Carvalhais, N., Ciais, P., Dalmonech, D., Fisher, J. B., Fisher, R., Friedlingstein, P., Hibbard, K., Hoffman, F. M., Huntzinger, D., Jones, C. D., Koven, C., Lawrence, D., Li, D. J., Mahecha, M., Niu, S. L., Norby, R., Piao, S. L., Qi, X., Peylin, P., Prentice, I. C., Riley, W., Reichstein, M., Schwalm, C., Wang, Y. P., Xia, J. Y., Zaehle, S., and Zhou, X. H.: A Framework for Benchmarking Land Models, *Biogeosciences*, 9, 3857–3874, <https://doi.org/10.5194/bg-9-3857-2012>, 2012.
- 1310 Ma, X., Zhao, S., Zhang, H., and Wang, W.: The Double-ITCZ Problem in CMIP6 and the Influences of Deep Convection and Model Resolution, *International Journal of Climatology*, 43, 2369–2390, <https://doi.org/https://doi.org/10.1002/joc.7980>, 2023.
- MacDougall, A. H., Frölicher, T. L., Jones, C. D., Rogelj, J., Matthews, H. D., Zickfeld, K., Arora, V. K., Barrett, N. J., Brovkin, V., Burger, F. A., Eby, M., Eliseev, A. V., Hajima, T., Holden, P. B., Jeltsch-Thömmes, A., Koven, C., Mengis, N., Menviel, L., Michou, M., Mokhov, I. I., Oka, A., Schwinger, J., Séférian, R., Shaffer, G., Sokolov, A., Tachiiri, K., Tjiputra, J., Wiltshire, A., and Ziehn, T.: Is There Warming in the Pipeline? A Multi-model Analysis of the Zero Emissions Commitment from CO₂, *Biogeosciences*, 17, 2987–3016, <https://doi.org/10.5194/bg-17-2987-2020>, 2020.
- 1320 Mankin, J. S., Lehner, F., Coats, S., and McKinnon, K. A.: The Value of Initial Condition Large Ensembles to Robust Adaptation Decision-Making, *Earth's Future*, 8, e2012EF001610, <https://doi.org/10.1029/2020EF001610>, 2020.
- Mauritsen, T., Stevens, B., Roeckner, E., Crueger, T., Esch, M., Giorgetta, M., Haak, H., Jungclaus, J., Klocke, D., Matei, D., Mikolajewicz, U., Notz, D., Pincus, R., Schmidt, H., and Tomassini, L.: Tuning the Climate of a Global Model, *Journal of Advances in Modeling Earth Systems*, 4, <https://doi.org/https://doi.org/10.1029/2012MS000154>, 2012.
- 1325 Mauritsen, T., Bader, J., Becker, T., Behrens, J., Bittner, M., Brokopf, R., Brovkin, V., Claussen, M., Crueger, T., Esch, M., Fast, I., Fiedler, S., Fläschner, D., Gayler, V., Giorgetta, M., Goll, D. S., Haak, H., Hagemann, S., Hedemann, C., Hohenegger, C., Ilyina, T., Jahns, T., Jimenez-de-la Cuesta, D., Jungclaus, J., Kleinen, T., Kloster, S., Kracher, D., Kinne, S., Kleberg, D., Lasslop, G., Kornbluh, L., Marotzke, J., Matei, D., Meraner, K., Mikolajewicz, U., Modali, K., Möbis, B., Müller, W. A., Nabel, J. E. M. S., Nam, C. C. W., Notz, D., Nyawira, S.-S., Paulsen, H., Peters, K., Pincus, R., Pohlmann, H., Pongratz, J., Popp, M., Raddatz, T. J., Rast, S., Redler, R., Reick,
- 1330



- C. H., Rohrschneider, T., Schemann, V., Schmidt, H., Schnur, R., Schulzweida, U., Six, K. D., Stein, L., Stemmler, I., Stevens, B., von Storch, J.-S., Tian, F., Voigt, A., Vrese, P., Wieners, K.-H., Wilkenskjeld, S., Winkler, A., and Roeckner, E.: Developments in the MPI-M Earth System Model version 1.2 (MPI-ESM1.2) and Its Response to Increasing CO₂, *Journal of Advances in Modeling Earth Systems*, 11, 998–1038, <https://doi.org/10.1029/2018MS001400>, 2019.
- 1335 Mayer, M., Kato, S., Bosilovich, M., Bechtold, P., Mayer, J., Schröder, M., Behrangi, A., Wild, M., Kobayashi, S., Li, Z., and L’Ecuyer, T.: Assessment of Atmospheric and Surface Energy Budgets Using Observation-Based Data Products, *Surveys in Geophysics*, 45, 1827–1854, <https://doi.org/10.1007/s10712-024-09827-x>, 2024.
- Meehl, G. A., Senior, C. A., Eyring, V., Flato, G., Lamarque, J.-F., Stouffer, R. J., Taylor, K. E., and Schlund, M.: Context for Interpreting Equilibrium Climate Sensitivity and Transient Climate Response from the CMIP6 Earth System Models, *Science Advances*, 6, eaba1981, <https://doi.org/10.1126/sciadv.aba1981>, 2020.
- 1340 Merrifield, A. L., Brunner, L., Lorenz, R., Humphrey, V., and Knutti, R.: Climate Model Selection by Independence, Performance, and Spread (ClimSIPS v1.0.1) for Regional Applications, *Geoscientific Model Development*, 16, 4715–4747, <https://doi.org/10.5194/gmd-16-4715-2023>, 2023.
- Miller, S. M., Hayek, M. N., Andrews, A. E., Fung, I., and Liu, J.: Biases in Atmospheric CO₂ Estimates from Correlated Meteorology Modeling Errors, *Atmospheric Chemistry and Physics*, 15, 2903–2914, <https://doi.org/10.5194/acp-15-2903-2015>, 2015.
- 1345 Moat, B. I., Smeed, D., Rayner, D., Johns, W. E., Smith, R. H., Volkov, D. L., Elipot, S., Petit, T., Kajtar, J. B., Baringer, M. O., and Collins, J.: Atlantic Meridional Overturning Circulation Observed by the RAPID-MOCHA-WBTS Array at 26N from 2004 to 2023 (v2023.1a), <https://doi.org/10.5285/33826d6e-801c-b0a7-e063-7086abc0b9db>, Archive Location: North Atlantic Ocean, Straits of Florida, Atlantic Ocean, Northeast Atlantic Ocean (40°W), Northwest Atlantic Ocean (40°W), published by NERC EDS British Oceanographic Data Centre NOC, 2025.
- 1350 Morice, C. P., Kennedy, J. J., Rayner, N. A., Winn, J. P., Hogan, E., Killick, R. E., Dunn, R. J. H., Osborn, T. J., Jones, P. D., and Simpson, I. R.: An Updated Assessment of Near-Surface Temperature Change From 1850: The HadCRUT5 Data Set, *Journal of Geophysical Research: Atmospheres*, 126, e2019JD032361, <https://doi.org/10.1029/2019JD032361>, 2021.
- Notz, D. and SIMIP Community: Arctic Sea Ice in CMIP6, *Geophysical Research Letters*, 47, e2019GL086749, <https://doi.org/https://doi.org/10.1029/2019GL086749>, 2020.
- 1355 Nowicki, S. M. J., Payne, A., Larour, E., Seroussi, H., Goelzer, H., Lipscomb, W., Gregory, J., Abe-Ouchi, A., and Shepherd, A.: Ice Sheet Model Intercomparison Project (ISMIP6) contribution to CMIP6, *Geoscientific Model Development*, 9, 4521–4545, <https://doi.org/10.5194/gmd-9-4521-2016>, 2016.
- Olsen, A., Key, R. M., van Heuven, S., Lauvset, S. K., Velo, A., Lin, X., Schirnack, C., Kozyr, A., Tanhua, T., Hoppema, M., Jutterström, S., 1360 Steinfeldt, R., Jeansson, E., Ishii, M., Pérez, F. F., and Suzuki, T.: The Global Ocean Data Analysis Project Version 2 (GLODAPv2) – An Internally Consistent Data Product for the World Ocean, *Earth System Science Data*, 8, 297–323, <https://doi.org/10.5194/essd-8-297-2016>, 2016.
- Ordóñez, A., Ullrich, P., and Lee, J.: cmecmetrics/cmec-driver: v1.1.9-release, <https://doi.org/10.5281/zenodo.15548748>, 2025.
- O’Rourke, E.: CMIP6 Community Survey Results, <https://doi.org/10.5281/zenodo.11654909>, 2025.
- 1365 Oueslati, B. and Bellon, G.: The Double ITCZ Bias in CMIP5 Models: Interaction Between SST, Large-scale Circulation and Precipitation, *Climate Dynamics*, 44, 585–607, <https://doi.org/10.1007/s00382-015-2468-6>, 2015.
- Pastorello, G., Trotta, C., Canfora, E., Chu, H., Christianson, D., Cheah, Y.-W., Poindexter, C., Chen, J., Elbashandy, A., Humphrey, M., Isaac, P., Polidori, D., Reichstein, M., Ribeca, A., van Ingen, C., Vuichard, N., Zhang, L., Amiro, B., Ammann, C., Arain, M. A., Ardö, J.,



- 1370 Arkebauer, T., Arndt, S. K., Arriga, N., Aubinet, M., Aurela, M., Baldocchi, D., Barr, A., Beamesderfer, E., Marchesini, L. B., Bergeron, O., Beringer, J., Bernhofer, C., Berveiller, D., Billesbach, D., Black, T. A., Blanken, P. D., Bohrer, G., Boike, J., Bolstad, P. V., Bonal, D., Bonnefond, J.-M., Bowling, D. R., Bracho, R., Brodeur, J., Brümmer, C., Buchmann, N., Burban, B., Burns, S. P., Buysse, P., Cale, P., Cavagna, M., Cellier, P., Chen, S., Chini, I., Christensen, T. R., Cleverly, J., Collalti, A., Consalvo, C., Cook, B. D., Cook, D., Coursolle, C., Cremonese, E., Curtis, P. S., D'Andrea, E., da Rocha, H., Dai, X., Davis, K. J., Cinti, B. D., Grandcourt, A. d., Ligne, A. D., De Oliveira, R. C., Delpierre, N., Desai, A. R., Di Bella, C. M., Tommasi, P. d., Dolman, H., Domingo, F., Dong, G., Dore, S., Duce, P., Dufrêne, E., Dunn, A., Dušek, J., Eamus, D., Eichelmann, U., ElKhidir, H. A. M., Eugster, W., Ewenz, C. M., Ewers, B., Famulari, D., Fares, S., Feigenwinter, I., Feitz, A., Fensholt, R., Filippa, G., Fischer, M., Frank, J., Galvagno, M., Gharun, M., Gianelle, D., Gielen, B., Gioli, B., Gitelson, A., Goded, I., Goeckede, M., Goldstein, A. H., Gough, C. M., Goulden, M. L., Graf, A., Griebel, A., Gruening, C., Grünwald, T., Hammerle, A., Han, S., Han, X., Hansen, B. U., Hanson, C., Hatakka, J., He, Y., Hehn, M., Heinesch, B., Hinko-Najera, N., Hörtnagl, L., Hutley, L., Ibrom, A., Ikawa, H., Jackowicz-Korczynski, M., Janouš, D., Jans, W., Jassal, R., Jiang, S., Kato, T., Khomik, M., Klatt, J., Knohl, A., Knox, S., Kobayashi, H., Koerber, G., Kolle, O., Kosugi, Y., Kotani, A., Kowalski, A., Kruijt, B., Kurbatova, J., Kutsch, W. L., Kwon, H., Launiainen, S., Laurila, T., Law, B., Leuning, R., Li, Y., Liddell, M., Limousin, J.-M., Lion, M., Liska, A. J., Lohila, A., López-Ballesteros, A., López-Blanco, E., Loubet, B., Loustau, D., Lucas-Moffat, A., Lüers, J., Ma, S., Macfarlane, C., Magliulo, V., Maier, R., Mammarella, I., Manca, G., Marcolla, B., Margolis, H. A., Marras, S., Massman, W., Mastepanov, M., Matamala, R., Matthes, J. H., Mazzenga, F., McCaughey, H., McHugh, I., McMillan, A. M. S., Merbold, L., Meyer, W., Meyers, T., Miller, S. D., Minerbi, S., Moderow, U., Monson, R. K., Montagnani, L., Moore, C. E., Moors, E., Moreaux, V., Moureaux, C., Munger, J. W., Nakai, T., Neiryneck, J., Nesic, Z., Nicolini, G., Noormets, A., Northwood, M., Noretto, M., Nouvellon, Y., Novick, K., Oechel, W., Olesen, J. E., Ourcival, J.-M., Papuga, S. A., Parmentier, F.-J., Paul-Limoges, E., Pavelka, M., Peichl, M., Pendall, E., Phillips, R. P., Pilegaard, K., Pirk, N., Posse, G., Powell, T., Prasse, H., Prober, S. M., Rambal, S., Rannik, U., Raz-Yaseef, N., Rebmann, C., Reed, D., Dios, V. R. d., Restrepo-Coupe, N., Reverter, B. R., Roland, M., Sabbatini, S., Sachs, T., Saleska, S. R., Sánchez-Cañete, E. P., Sanchez-Mejia, Z. M., Schmid, H. P., Schmidt, M., Schneider, K., Schrader, F., Schroder, I., Scott, R. L., Sedlák, P., Serrano-Ortíz, P., Shao, C., Shi, P., Shironya, I., Siebicke, L., Šigut, L., Silberstein, R., Sirca, C., Spano, D., Steinbrecher, R., Stevens, R. M., Sturtevant, C., Suyker, A., Tagesson, T., Takahashi, S., Tang, Y., Tapper, N., Thom, J., Tomassucci, M., Tuovinen, J.-P., Urbanski, S., Valentini, R., van der Molen, M., van Gorsel, E., van Huissteden, K., Varlagin, A., Verfaillie, J., Vesala, T., Vincke, C., Vitale, D., Vygodskaya, N., Walker, J. P., Walter-Shea, E., Wang, H., Weber, R., Westermann, S., Wille, C., Wofsy, S., Wohlfahrt, G., Wolf, S., Woodgate, W., Li, Y., Zampedri, R., Zhang, J., Zhou, G., Zona, D., Agarwal, D., Biraud, S., Torn, M., and Papale, D.: The FLUXNET2015 Dataset and the ONEFlux Processing Pipeline for Eddy Covariance Data, *Scientific Data*, 7, 225, <https://doi.org/10.1038/s41597-020-0534-3>, 2020.
- 1385 Planton, Y. Y., Guilyardi, E., Wittenberg, A. T., Lee, J., Gleckler, P. J., Bayr, T., McGregor, S., McPhaden, M. J., Power, S., Roehrig, R., Vialard, J., and Voltaire, A.: Evaluating Climate Models with the CLIVAR 2020 ENSO Metrics Package, *Bulletin of the American Meteorological Society*, 102, E193–E217, <https://doi.org/10.1175/BAMS-D-19-0337.1>, 2021.
- 1400 Priestley, M. D. K., Ackerley, D., Catto, J. L., Hodges, K. I., McDonald, R. E., and Lee, R. W.: An Overview of the Extratropical Storm Tracks in CMIP6 Historical Simulations, *Journal of Climate*, 33, 6315–6343, <https://doi.org/10.1175/JCLI-D-19-0928.1>, 2020.
- Rayner, N. A., Parker, D. E., Horton, E. B., Folland, C. K., Alexander, L. V., Rowell, D. P., Kent, E. C., and Kaplan, A.: Global Analyses of Sea Surface Temperature, Sea Ice, and Night Marine Air Temperature Since the Late Nineteenth Century, *Journal of Geophysical Research: Atmospheres*, 108, <https://doi.org/10.1029/2002JD002670>, 2003.
- 1405 Reagan, J. R., Boyer, T. P., García, H. E., Locarnini, R. A., Baranova, O. K., Bouchard, C., Cross, S. L., Mishonov, A. V., Paver, C. R., Seidov, D., Wang, Z., and Dukhovskoy, D.: World Ocean Atlas 2023, <https://doi.org/10.25921/VA26-HV25>, 2023.



- 1410 Righi, M., Andela, B., Eyring, V., Lauer, A., Predoi, V., Schlund, M., Vegas-Regidor, J., Bock, L., Brötz, B., de Mora, L., Diblen, F., Dreyer, L., Drost, N., Earnshaw, P., Hassler, B., Koldunov, N., Little, B., Loosveldt Tomas, S., and Zimmermann, K.: Earth System Model Evaluation Tool (ESMValTool) v2.0 – Technical Overview, *Geoscientific Model Development*, 13, 1179–1199, <https://doi.org/10.5194/gmd-13-1179-2020>, 2020.
- Roach, L. A., Dörr, J., Holmes, C. R., Massonnet, F., Blockley, E. W., Notz, D., Rackow, T., Raphael, M. N., O’Farrell, S. P., Bailey, D. A., and Bitz, C. M.: Antarctic Sea Ice Area in CMIP6, *Geophysical Research Letters*, 47, e2019GL086729, <https://doi.org/10.1029/2019GL086729>, 2020.
- 1415 Roberts, M. J., Camp, J., Seddon, J., Vidale, P. L., Hodges, K., Vanni ere, B., Mecking, J., Haarsma, R., Bellucci, A., Scoccimarro, E., Caron, L.-P., Chauvin, F., Terray, L., Valcke, S., Moine, M.-P., Putrasahan, D., Roberts, C. D., Senan, R., Zarzycki, C., Ullrich, P., Yamada, Y., Mizuta, R., Kodama, C., Fu, D., Zhang, Q., Danabasoglu, G., Rosenbloom, N., Wang, H., and Wu, L.: Projected Future Changes in Tropical Cyclones Using the CMIP6 HighResMIP Multimodel Ensemble, *Geophysical Research Letters*, 47, e2020GL088662, <https://doi.org/10.1029/2020GL088662>, 2020.
- 1420 Rutz, J. J., Shields, C. A., Lora, J. M., Payne, A. E., Guan, B., Ullrich, P., O’Brien, T., Leung, L. R., Ralph, F. M., Wehner, M., Brands, S., Collow, A., Goldenson, N., Gorodetskaya, I., Griffith, H., Kashinath, K., Kawzenuk, B., Krishnan, H., Kurlin, V., Lavers, D., Magnusdottir, G., Mahoney, K., McClenny, E., Muszynski, G., Nguyen, P. D., Prabhat, M., Qian, Y., Ramos, A. M., Sarangi, C., Sellars, S., Shulgina, T., Tome, R., Waliser, D., Walton, D., Wick, G., Wilson, A. M., and Viale, M.: The Atmospheric River Tracking Method Intercomparison Project (ARTMIP): Quantifying Uncertainties in Atmospheric River Climatology, *Journal of Geophysical Research: Atmospheres*, 124, 13 777–13 802, <https://doi.org/10.1029/2019JD030936>, 2019.
- 1425 Sanderson, B. M., Brovkin, V., Fisher, R., Hohn, D., Ilyina, T., Jones, C., Koenigk, T., Koven, C., Li, H., Lawrence, D., Lawrence, P., Liddicoat, S., Macdougall, A., Mengis, N., Nicholls, Z., O’Rourke, E., Romanou, A., Sandstad, M., Schwinger, J., Seferian, R., Sentman, L., Simpson, I., Smith, C., Steinert, N., Swann, A., Tjiputra, J., and Ziehn, T.: flat10MIP: An Emissions-driven Experiment to Diagnose the Climate Response to Positive, Zero, and Negative CO₂ Emissions, *EGUsphere*, pp. 1–39, <https://doi.org/10.5194/egusphere-2024-3356>, 2024.
- 1430 Santoro, M. and Cartus, O.: ESA Biomass Climate Change Initiative (Biomass_cci): Global datasets of Forest Above-ground Biomass for the Years 2010, 2015, 2016, 2017, 2018, 2019, 2020 and 2021, v5.01, <https://doi.org/10.5285/BF535053562141C6BB7AD831F5998D77>, published by NERC EDS Centre for Environmental Data Analysis, 2024.
- Scaife, A. A., Woollings, T., Knight, J., Martin, G., and Hinton, T.: Atmospheric Blocking and Mean Biases in Climate Models, *Journal of Climate*, 23, 6143–6152, <https://doi.org/10.1175/2010JCLI3728.1>, 2010.
- 1435 Schlund, M., Lauer, A., Gentine, P., Sherwood, S. C., and Eyring, V.: Emergent Constraints on Equilibrium Climate Sensitivity in CMIP5: Do They Hold for CMIP6?, *Earth System Dynamics*, 11, 1233–1258, <https://doi.org/10.5194/esd-11-1233-2020>, 2020.
- Schlund, M., Hassler, B., Lauer, A., Andela, B., J ockel, P., Kazeroni, R., Loosveldt Tomas, S., Medeiros, B., Predoi, V., S en esi, S., Servonnat, J., Stacke, T., Vegas-Regidor, J., Zimmermann, K., and Eyring, V.: Evaluation of Native Earth System Model Output with ESMValTool v2.6.0, *Geoscientific Model Development*, 16, 315–333, <https://doi.org/10.5194/gmd-16-315-2023>, 2023.
- 1440 Seneviratne, S. I., Corti, T., Davin, E. L., Hirschi, M., Jaeger, E. B., Lehner, I., Orlowsky, B., and Teuling, A. J.: Investigating Soil Moisture–Climate Interactions in a Changing Climate: A Review, *Earth-Science Reviews*, 99, 125–161, <https://doi.org/https://doi.org/10.1016/j.earscirev.2010.02.004>, 2010.



- Senior, C. A., Jones, C. G., Wood, R. A., Sellar, A., Belcher, S., Klein-Tank, A., Sutton, R., Walton, J., Lawrence, B., Andrews, T., and Mul-
ahy, J. P.: U.K. Community Earth System Modeling for CMIP6, *Journal of Advances in Modeling Earth Systems*, 12, e2019MS002004,
1445 <https://doi.org/10.1029/2019MS002004>, 2020.
- Slivinski, L. C., Compo, G. P., Whitaker, J. S., Sardeshmukh, P. D., Giese, B. S., McColl, C., Allan, R., Yin, X., Vose, R., Titchner, H.,
Kennedy, J., Spencer, L. J., Ashcroft, L., Brönnimann, S., Brunet, M., Camuffo, D., Cornes, R., Cram, T. A., Crouthamel, R., Domínguez-
Castro, F., Freeman, J. E., Gergis, J., Hawkins, E., Jones, P. D., Jourdain, S., Kaplan, A., Kubota, H., Blancq, F. L., Lee, T.-C., Lorrey, A.,
Luterbacher, J., Maugeri, M., Mock, C. J., Moore, G. K., Przybylak, R., Pudmenzky, C., Reason, C., Slonosky, V. C., Smith, C. A., Tinz,
1450 B., Trewin, B., Valente, M. A., Wang, X. L., Wilkinson, C., Wood, K., and Wyszyński, P.: Towards a More Reliable Historical Reanalysis:
Improvements for Version 3 of the Twentieth Century Reanalysis System, *Quarterly Journal of the Royal Meteorological Society*, 145,
2876–2908, <https://doi.org/10.1002/qj.3598>, 2019.
- Sofieva, V. F., Szlag, M., Tamminen, J., Arosio, C., Rozanov, A., Weber, M., Degenstein, D., Bourassa, A., Zawada, D., Kiefer, M., Laeng,
A., Walker, K. A., Sheese, P., Hubert, D., van Roozendaal, M., Retscher, C., Damadeo, R., and Lumpe, J. D.: Updated Merged SAGE-
1455 CCI-OMPS+ Dataset for the Evaluation of Ozone Trends in the Stratosphere, *Atmospheric Measurement Techniques*, 16, 1881–1899,
<https://doi.org/10.5194/amt-16-1881-2023>, 2023.
- Stainforth, D. A., Allen, M. R., Tredger, E. R., and Smith, L. A.: Confidence, Uncertainty and Decision-support Relevance
in Climate Predictions, *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*,
<https://doi.org/10.1098/rsta.2007.2074>, 2007.
- 1460 Stockhause, M., Huard, D., Al Khourdajie, A., Gutiérrez, J. M., Kawamiya, M., Klutse, N. A. B., Krey, V., Milward, D., Okem, A. E., Pirani,
A., Sitz, L. E., Solman, S. A., Spinuso, A., and Xing, X.: Implementing FAIR Data Principles in the IPCC Seventh Assessment Cycle:
Lessons Learned and Future Prospects, *PLOS Climate*, 3, e0000533, <https://doi.org/10.1371/journal.pclm.0000533>, 2024.
- Swaminathan, R., Parker, R. J., Jones, C. G., Allan, R. P., Quaife, T., Kelley, D. I., Mora, L. d., and Walton, J.: The Physical Climate at Global
Warming Thresholds as Seen in the U.K. Earth System Model, *Journal of Climate*, 35, 29–48, <https://doi.org/10.1175/JCLI-D-21-0234.1>,
1465 2022.
- Séférian, R., Nabat, P., Michou, M., Saint-Martin, D., Voldoire, A., Colin, J., Decharme, B., Delire, C., Berthet, S., Chevallier, M., Sénési,
S., Franchisteguy, L., Vial, J., Mallet, M., Joetzjer, E., Geoffroy, O., Guérémy, J.-F., Moine, M.-P., Msadek, R., Ribes, A., Rocher, M.,
Roehrig, R., Salas-y Mélia, D., Sanchez, E., Terray, L., Valcke, S., Waldman, R., Aumont, O., Bopp, L., Deshayes, J., Éthé, C., and
Madec, G.: Evaluation of CNRM Earth System Model, CNRM-ESM2-1: Role of Earth System Processes in Present-Day and Future
1470 Climate, *Journal of Advances in Modeling Earth Systems*, 11, 4182–4227, <https://doi.org/10.1029/2019MS001791>, 2019.
- Séférian, R., Berthet, S., Yool, A., Palmiéri, J., Bopp, L., Tagliabue, A., Kwiatkowski, L., Aumont, O., Christian, J., Dunne, J., Gehlen, M.,
Ilyina, T., John, J. G., Li, H., Long, M. C., Luo, J. Y., Nakano, H., Romanou, A., Schwinger, J., Stock, C., Santana-Falcón, Y., Takano, Y.,
Tjiputra, J., Tsujino, H., Watanabe, M., Wu, T., Wu, F., and Yamamoto, A.: Tracking Improvement in Simulated Marine Biogeochemistry
Between CMIP5 and CMIP6, *Current Climate Change Reports*, 6, 95–119, <https://doi.org/10.1007/s40641-020-00160-0>, 2020.
- 1475 Taylor, K., Doutriaux, C., and Peterschmitt, J.: Climate Model Output Rewriter (CMOR), Tech. Rep. UCRL-TR-204637, 15014202,
Lawrence Livermore National Laboratory, <https://doi.org/10.2172/15014202>, 2004.
- Taylor, K. E.: Summarizing Multiple Aspects of Model Performance in a Single Diagram, *Journal of Geophysical Research: Atmospheres*,
106, 7183–7192, <https://doi.org/10.1029/2000JD900719>, 2001.
- Teixeira, J., Waliser, D., Ferraro, R., Gleckler, P., Lee, T., and Potter, G.: Satellite Observations for CMIP5: The Genesis of Obs4MIPs,
1480 *Bulletin of the American Meteorological Society*, 95, 1329–1334, <https://doi.org/10.1175/BAMS-D-12-00204.1>, 2025.



- Tian, B. and Dong, X.: The Double-ITCZ Bias in CMIP3, CMIP5, and CMIP6 Models Based on Annual Mean Precipitation, *Geophysical Research Letters*, 47, e2020GL087232, <https://doi.org/10.1029/2020GL087232>, 2020.
- Tjiputra, J. F., Negrel, J., and Olsen, A.: Early Detection of Anthropogenic Climate Change Signals in the Ocean Interior, *Scientific Reports*, 13, 3006, <https://doi.org/10.1038/s41598-023-30159-0>, 2023.
- 1485 Trenberth, K. E. and Caron, J. M.: Estimates of Meridional Atmosphere and Ocean Heat Transports, *Journal of Climate*, 14, 3433–3443, [https://doi.org/10.1175/1520-0442\(2001\)014<3433:EOMAAO>2.0.CO;2](https://doi.org/10.1175/1520-0442(2001)014<3433:EOMAAO>2.0.CO;2), 2001.
- Trenberth, K. E., Branstator, G. W., Karoly, D., Kumar, A., Lau, N.-C., and Ropelewski, C.: Progress During TOGA in Understanding and Modeling Global Teleconnections Associated with Tropical Sea Surface Temperatures, *Journal of Geophysical Research: Oceans*, 103, 14291–14324, <https://doi.org/10.1029/97JC01444>, 1998.
- 1490 Trenberth, K. E., Smith, L., Qian, T., Dai, A., and Fasullo, J.: Estimates of the Global Water Budget and Its Annual Cycle Using Observational and Model Data, *Journal of Hydrometeorology*, 8, 758–769, <https://doi.org/10.1175/JHM600.1>, 2007.
- Trumbore, S. E. and Czimczik, C. I.: An Uncertain Future for Soil Carbon, *Science*, 321, 1455–1456, <https://doi.org/10.1126/science.1160232>, 2008.
- Vaittinada Ayar, P., Battisti, D. S., Li, C., King, M., Vrac, M., and Tjiputra, J.: A Regime View of ENSO Flavors Through Clustering in CMIP6 Models, *Earth's Future*, 11, e2022EF003460, <https://doi.org/10.1029/2022EF003460>, 2023.
- 1495 Waliser, D., Gleckler, P. J., Ferraro, R., Taylor, K. E., Ames, S., Biard, J., Bosilovich, M. G., Brown, O., Chepfer, H., Cinquini, L., Durack, P. J., Eyring, V., Mathieu, P.-P., Lee, T., Pinnock, S., Potter, G. L., Rixen, M., Saunders, R., Schulz, J., Thépaut, J.-N., and Tuma, M.: Observations for Model Intercomparison Project (Obs4MIPs): Status for CMIP6, *Geoscientific Model Development*, 13, 2945–2958, <https://doi.org/10.5194/gmd-13-2945-2020>, 2020.
- 1500 Wallace, J. M. and Gutzler, D. S.: Teleconnections in the Geopotential Height Field during the Northern Hemisphere Winter, *Monthly Weather Review*, 109, 784–812, [https://doi.org/10.1175/1520-0493\(1981\)109<0784:TITGHF>2.0.CO;2](https://doi.org/10.1175/1520-0493(1981)109<0784:TITGHF>2.0.CO;2), 1981.
- Wan, H., Zhang, K., Rasch, P. J., Singh, B., Chen, X., and Edwards, J.: A New and Inexpensive Non-bit-for-bit Solution Reproducibility Test Based on Time Step Convergence (TSC1.0), *Geoscientific Model Development*, 10, 537–552, <https://doi.org/10.5194/gmd-10-537-2017>, 2017.
- 1505 Wang, B., Kim, H.-J., Kikuchi, K., and Kitoh, A.: Diagnostic Metrics for Evaluation of Annual and Diurnal Cycles, *Climate Dynamics*, 37, 941–955, <https://doi.org/10.1007/s00382-010-0877-0>, 2011.
- Wang, H., Pearson, B., Hou, A., Lanfer Marquez, A., and Bonnet, P.: Model Evaluation and Benchmarking: Community Survey Summary Report, <https://doi.org/10.5281/zenodo.15212597>, 2025.
- Wang, Y. and Mao, J.: Global Multi-layer Soil Moisture Products, <https://doi.org/10.6084/M9.FIGSHARE.13661312.V1>, data set, 2021.
- 1510 Wang, Y., Mao, J., Jin, M., Hoffman, F. M., Shi, X., Wullschleger, S. D., and Dai, Y.: Development of Observation-based Global Multilayer Soil Moisture Products for 1970 to 2016, *Earth System Science Data*, 13, 4385–4405, <https://doi.org/10.5194/essd-13-4385-2021>, 2021.
- Warszawski, L., Frieler, K., Huber, V., Piontek, F., Serdeczny, O., and Schewe, J.: The Inter-Sectoral Impact Model Intercomparison Project (ISI-MIP): Project Framework, *Proceedings of the National Academy of Sciences*, 111, 3228–3232, <https://doi.org/10.1073/pnas.1312330110>, 2014.
- 1515 Weigel, K., Bock, L., Gier, B. K., Lauer, A., Righi, M., Schlund, M., Adeniyi, K., Andela, B., Arnone, E., Berg, P., Caron, L.-P., Cionni, I., Corti, S., Drost, N., Hunter, A., Lledó, L., Mohr, C. W., Paçal, A., Pérez-Zanón, N., Predoi, V., Sandstad, M., Sillmann, J., Sterl, A., Vegas-Regidor, J., von Hardenberg, J., and Eyring, V.: Earth System Model Evaluation Tool (ESMValTool) v2.0 – Diagnostics for



- Extreme Events, Regional and Impact Evaluation, and Analysis of Earth System Models in CMIP, Geoscientific Model Development, 14, 3159–3184, <https://doi.org/10.5194/gmd-14-3159-2021>, 2021.
- 1520 Wilkinson, M. D., Dumontier, M., Aalbersberg, I. J., Appleton, G., Axton, M., Baak, A., Blomberg, N., Boiten, J.-W., Da Silva Santos, L. B., Bourne, P. E., Bouwman, J., Brookes, A. J., Clark, T., Crosas, M., Dillo, I., Dumon, O., Edmunds, S., Evelo, C. T., Finkers, R., Gonzalez-Beltran, A., Gray, A. J., Groth, P., Goble, C., Grethe, J. S., Heringa, J., 'T Hoen, P. A., Hooft, R., Kuhn, T., Kok, R., Kok, J., Lusher, S. J., Martone, M. E., Mons, A., Packer, A. L., Persson, B., Rocca-Serra, P., Roos, M., Van Schaik, R., Sansone, S.-A., Schultes, E., Sengstag, T., Slater, T., Strawn, G., Swertz, M. A., Thompson, M., Van Der Lei, J., Van Mulligen, E., Velterop, J., Waagmeester, A.,
- 1525 Wittenburg, P., Wolstencroft, K., Zhao, J., and Mons, B.: The FAIR Guiding Principles for Scientific Data Management and Stewardship, *Scientific Data*, 3, 160018, <https://doi.org/10.1038/sdata.2016.18>, 2016.
- Winker, D.: CALIPSO Lidar Level 3 Ice Cloud Data, Standard V1-00, https://doi.org/10.5067/CALIOP/CALIPSO/L3_ICE_CLOUD-STANDARD-V1-00, published by NASA Langley Atmospheric Science Data Center Distributed Active Archive Center, 2024.
- Winker, D., Cai, X., Vaughan, M., Garnier, A., Magill, B., Avery, M., and Getzewich, B.: A Level 3 Monthly Gridded Ice Cloud Dataset
- 1530 Derived from 12 Years of CALIOP Measurements, *Earth System Science Data*, 16, 2831–2855, <https://doi.org/10.5194/essd-16-2831-2024>, 2024.
- Yang, X., Ricciuto, D. M., Thornton, P. E., Shi, X., Xu, M., Hoffman, F. M., and Norby, R. J.: The Effects of Phosphorus Cycle Dynamics on Carbon Sources and Sinks in the Amazon Region: A Modeling Study Using ELM v1, *Journal of Geophysical Research: Biogeosciences*, 124, 3686–3698, <https://doi.org/10.1029/2019JG005082>, 2019.
- 1535 Yukimoto, S., Kawai, H., Koshiro, T., Oshima, N., Yoshida, K., Urakawa, S., Tsujino, H., Deushi, M., Tanaka, T., Hosaka, M., Yabu, S., Yoshimura, H., Shindo, E., Mizuta, R., Obata, A., Adachi, Y., and Ishii, M.: The Meteorological Research Institute Earth System Model Version 2.0, MRI-ESM2.0: Description and Basic Evaluation of the Physical Component, *Journal of the Meteorological Society of Japan. Ser. II*, 97, 931–965, <https://doi.org/10.2151/jmsj.2019-051>, 2019.
- Zhu, Q., Riley, W. J., Tang, J., Collier, N., Hoffman, F. M., Yang, X., and Bisht, G.: Representing Nitrogen, Phosphorus, and Carbon
- 1540 Interactions in the E3SM Land Model: Development and Global Benchmarking, *Journal of Advances in Modeling Earth Systems*, 11, 2238–2258, <https://doi.org/10.1029/2018MS001571>, 2019.