

ESMValTool (v1.0) - A community diagnostic and performance metrics tool for routine evaluation of Earth System Models in CMIP

Veronika Eyring¹, Mattia Righi¹, Axel Lauer¹, Martin Evaldsson², Sabrina Wenzel¹, Colin Jones^{3,4}, Alessandro Anav⁵, Oliver Andrews⁶, Irene Cionni⁷, Edouard L. Davin⁸, Clara Deser⁹, Carsten Ehbrecht¹⁰, Pierre Friedlingstein⁵, Peter Gleckler¹¹, Klaus-Dirk Gottschaldt¹, Stefan Hagemann¹², Martin Juckes¹³, Stephan Kindermann¹⁰, John Krasting¹⁴, Dominik Kunert¹, Richard Levine⁴, Alexander Loew^{15,12}, Jarmo Mäkelä¹⁶, Gill Martin⁴, Erik Mason^{14,17}, Adam Phillips⁹, Simon Read¹⁸, Catherine Rio¹⁹, Romain Roehrig²⁰, Daniel Senftleben¹, Andreas Sterl²¹, Lambertus H. van Ulf²¹, Jeremy Walton⁴, Shiyu Wang², and Keith D. Williams⁴

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

[2] Swedish Meteorological and Hydrological Institute (SMHI), 60176 Norrköping, Sweden.

[3] University of Leeds, Leeds, UK

[4] Met Office Hadley Centre, Exeter, UK

[5] University of Exeter, Exeter, UK

[6] Tyndall Centre for Climate Change Research, School of Environmental Sciences, University of East Anglia, Norwich, UK

[7] Agenzia nazionale per le nuove tecnologie, l'energia e lo sviluppo economico sostenibile (ENEA), Rome, Italy

[8] ETH Zurich, Switzerland

[9] National Center for Atmospheric Research (NCAR), Boulder, USA

[10] Deutsches Klimarechenzentrum, Hamburg, Germany

[11] Program for Climate Model Diagnosis and Intercomparison, Lawrence Livermore National Laboratory, Livermore, CA, USA

[12] Max-Planck-Institute for Meteorology, Hamburg, Germany

[13] National Centre for Atmospheric Science, British Atmospheric Data Centre, STFC Rutherford Appleton Laboratory, United Kingdom

[14] Geophysical Fluid Dynamics Laboratory/NOAA, Princeton, NJ, USA

[15] Ludwig-Maximilians-Universität München, Munich, Germany

[16] Finnish Meteorological Institute, Finland

[17] Engility Corporation, Chantilly, VA, USA

[18] University of Reading, Reading, UK

[19] Institut Pierre Simon Laplace, Paris, France

[20] CNRM-GAME, Météo France and CNRS, Toulouse, France

[21] Royal Netherlands Meteorological Institute (KNMI), De Bilt, The Netherlands

Correspondence to: V. Eyring (veronika.eyring@dlr.de)

Abstract

A community diagnostics and performance metrics tool for the evaluation of Earth System Models (ESMs) has been developed that allows for routine comparison of single or multiple models, either against predecessor versions or against observations. The priority of the effort so far has been to target specific scientific themes focusing on selected Essential Climate Variables (ECVs), a range of known systematic biases common to ESMs, such as coupled tropical climate variability, monsoons, Southern Ocean processes, continental dry biases and soil hydrology-climate interactions, as well as atmospheric CO₂ budgets, tropospheric and stratospheric ozone, and tropospheric aerosols. The tool is being developed in such a way that additional analyses can easily be added. A set of standard namelists for each scientific topic reproduces specific sets of diagnostics or performance metrics that have demonstrated their importance in ESM evaluation in the peer-reviewed literature. The Earth System Model Evaluation Tool (ESMValTool) is a community effort open to both users and developers encouraging open exchange of diagnostic source code and evaluation results from the CMIP ensemble. This will facilitate and improve ESM evaluation beyond the state-of-the-art and aims at supporting such activities within the Coupled Model Intercomparison Project (CMIP) and at individual modelling centres. Ultimately, we envisage

running the ESMValTool alongside the Earth System Grid Federation (ESGF) as part of a more routine evaluation of CMIP model simulations while utilizing observations available in standard formats (obs4MIPs) or provided by the user.

1. Introduction

Earth System Model (ESM) evaluation with observations or reanalyses is performed both to understand the performance of a given model and to gauge the quality of a new model, either against predecessor versions or a wider set of models. Over the past decades, the benefits of multi-model intercomparison projects such as the Coupled Model Intercomparison Project (CMIP) have been demonstrated. Since the beginning of CMIP in 1995, participating models have been further developed, with more complex and higher resolution models joining in CMIP5 (Taylor et al., 2012) which supported the Intergovernmental Panel on Climate Change (IPCC) Fifth Assessment Report (AR5) (IPCC, 2013). The main purpose of these internationally coordinated model experiments is to address outstanding scientific questions, to improve the understanding of climate, and to provide estimates of future climate change. Standardization of model output in a format that follows the Network Common Data Format (netCDF) Climate and Forecast (CF) Metadata Convention (<http://cfconventions.org/>) and collection of the model output on the Earth System Grid Federation (ESGF, <http://esgf.llnl.gov/>) facilitated multi-model analyses. However, CMIP has historically lacked a common analysis tool available that could operate directly on submitted model data and deliver a standard evaluation of models against observations.

An important new aspect in the next phase of CMIP (i.e., CMIP6 (Eyring et al., 2015)) is a more distributed organization under the oversight of the CMIP Panel, where a set of standard model experiments, which were common across earlier CMIP cycles, the Diagnostic, Evaluation and Characterization of Klima (DECK) experiments and the CMIP6 historical simulations, will be used to broadly characterize model performance and sensitivity to standard external forcing. Standardization, coordination, common infrastructure, and documentation functions that make the simulation results and their main characteristics available to the broader community are envisaged to be a central part of CMIP6. The Earth System Model Evaluation Tool (ESMValTool) presented here is a community development that can be used as one of the documentation functions in CMIP to help diagnose and understand the origin and consequences of model biases and inter-model spread. Our goal is to develop an evaluation tool that users can run to produce well-established analyses of the CMIP models once the output becomes available on the ESGF. This is realized

1 through text files that we refer to as standard namelists, each calling a certain set of diagnostics and
2 performance metrics to reproduce analyses that have demonstrated to be of importance in ESM
3 evaluation in previous peer-reviewed papers or assessment reports. Through this approach routine
4 and systematic evaluation of model results can be made more efficient. The framework enables
5 scientists to focus on developing more innovative analysis methods rather than constantly having to
6 “re-invent the wheel”. An additional purpose of the ESMValTool is to facilitate model evaluation at
7 individual modelling centres, in particular to rapidly assess the performance of a new model against
8 predecessor versions. Righi et al. (2015) and Jöckel et al. (2015) have applied a subset of the
9 namelists presented here to evaluate a set of simulations using different configurations of the global
10 ECHAM/MESSy Atmospheric Chemistry model (EMAC). In this paper we also highlight the
11 integration of ESMValTool into modelling workflows – including models developed at NOAA’s
12 Geophysical Fluid Dynamics Laboratory (GFDL), the EMAC model, and the NEMO ocean model
13 – through the use of the ESMValTool’s reformatting routine capabilities.

14 In addition to standardized model output, the ESGF hosts observations for Model Intercomparison
15 Projects (obs4MIPs (Ferraro et al., 2015; Teixeira et al., 2014)) and reanalyses data (ana4MIPs,
16 <https://www.earthsystemcog.org/projects/ana4mips>). The obs4MIPs and ana4MIPs projects provide
17 the community with access to CMIP-like data sets (in terms of variables, temporal and spatial
18 frequencies, and time periods) of satellite data and reanalyses, together with the corresponding
19 technical documentation. The ESMValTool makes use of these observations as well as observations
20 available from other sources to evaluate the models. In several of the diagnostics and metrics, more
21 than one observational data set or meteorological reanalysis is used to account for uncertainties in
22 observations. This is crucial for assessing model performance in a more robust and scientifically
23 valid way.

24 For the model evaluation we apply diagnostics and in several cases also performance metrics.
25 Diagnostics (e.g., the calculation of zonal means or derived variables in comparison to
26 observations) provide a qualitative comparison of the models with observations. Performance
27 metrics are defined as a quantitative measure of agreement between a simulated and observed
28 quantity which can be used to assess the performance of individual models or generation of models.
29 Quantitative performance metrics are routinely calculated for numerical weather forecast models,
30 but have been increasingly applied to Atmosphere-Ocean General Circulation Models (AOGCMs)
31 or ESMs. Performance metrics used in these studies have mainly focused on climatological mean
32 values of selected ECVs (Connolley and Bracegirdle, 2007; Gleckler et al., 2008; Pincus et al.,

2008; Reichler and Kim, 2008), and only a few studies have developed process-based performance metrics (SPARC-CCMVal, 2010; Waugh and Eyring, 2008; Williams and Webb, 2009). The implementation of performance metrics in the ESMValTool enables a quantitative assessment of model improvements, both for different versions of individual ESMs and for different generations of model ensembles used in international assessments (e.g., CMIP5 versus CMIP6). Application of performance metrics to multiple models helps highlighting when and where one or more models represent a particular process well. While quantitative metrics provide a valuable summary of overall model performance, they usually do not give information on how particular aspects of a model's simulation interact to determine the overall fidelity. For example, a model could simulate a mean state (and trend) in global mean surface temperature that agrees well with observations, but this could be due to compensating errors. To learn more about the sources of errors and uncertainties in models and thereby highlight specific areas requiring improvement, evaluation of the underlying processes and phenomena is necessary. A range of diagnostics and performance metrics focussing on a number of key processes are also included in ESMValTool.

This paper describes ESMValTool version 1.0 (v1.0) which is the first release of the tool to the wider community for application and further development as open source software. It demonstrates the use of the tool by showing example figures for each namelist for either all or a subset of CMIP5 models. Section 2 describes the technical aspects of the tool, and Section 3 the type of modelling and observational data currently supported by ESMValTool (v1.0). In Section 4 an overview of the namelists of ESMValTool (v1.0) is given along with their diagnostics and performance metrics and the variables and observations used. Section 5 describes the use of the ESMValTool in a typical model development cycle and evaluation workflow and Section 6 closes with a summary and an outlook.

2. Brief overview of the ESMValTool

In this section we give a brief overview of ESMValTool (v1.0) which is schematically depicted in Fig. 1. A detailed user's guide is provided in the Supplement.

The ESMValTool consists of a workflow manager and a number of diagnostic and graphical output scripts. It builds on a previously published diagnostic tool for chemistry-climate model evaluation (CCMVal-Diag Tool, Gettelman et al. (2012)), but is different in its focus. In particular, it extends to ESMs by including diagnostics and performance metrics relevant for the coupled Earth system, and also focuses on evaluating models with a common set of diagnostics rather than being mostly flexible as the CCMVal-Diag tool. In addition, several technical and structural changes have been

1 made that facilitate development by multiple users. The workflow manager is written in Python,
2 while a multi-language support is provided in the diagnostic and the graphic routines. The current
3 version supports Python (www.python.org), the NCAR Command Language (NCL, 2016) and R
4 (Ihaka and Gentleman, 1996), but it can be extended to other open-source languages. The
5 ESMValTool is executed by invoking the *main.py* script, which takes a namelist as a single input
6 argument. The namelists are text files written using the XML (eXtensible Markup Language) syntax
7 and define the data to be read (models and observations), the variables to be analysed and the
8 diagnostics to be applied. The XML-syntax has been chosen in order to allow users to express the
9 relationship between these three elements (data, variables and diagnostics) in a structured, easy to
10 use way.

11 Within the workflow, the input data are checked for compliance with the CF and Climate Model
12 Output Rewriter (CMOR, <http://pcmdi.github.io/cmor-site/tables.html>) standards required by the
13 tool (see Section 3) via a set of dedicated reformatting routines, which are also able to fix the most
14 common errors in the input data (e.g., wrong coordinates, undefined or missing values, non-
15 compliant units, etc.). It is additionally possible to define new variables using variable-specific
16 scripts, for example to calculate the total column ozone from a 3D ozone field (tro3), temperature
17 (ta) and surface pressure (ps). The diagnostic and graphic routines are written in a modular and
18 flexible way so that they can be customized by the user via diagnostic-specific settings in the
19 configuration file (cfg-file) and variable-specific settings (in the directory *variable_defs/*) without
20 changing the source code. These routines are complemented by a set of libraries, providing general-
21 purpose code for the most common operations (statistical analyses, regridding tools, graphic styles,
22 etc.). The output of the tool can be both NetCDF and graphics files in various formats. In addition, a
23 log file is written containing all the information of a specific call of the main script: creation date of
24 running the script, version number, analysed data (models and observations), applied diagnostics
25 and variables, and corresponding references. This helps to increase the traceability and
26 reproducibility of the results.

27 To facilitate the development of new namelists and diagnostics by multiple developers from various
28 institutions while preserving code quality and reliability, an automated testing framework is
29 included in the package. This allows the developers to verify that modifications and new code are
30 compatible with the existing code and do not change the results of existing diagnostics. Automated
31 testing within the ESMValTool is implemented on two complementary levels:

- 32 • unittests are used to verify that small code units (e.g., functions/subroutines) provide the

1 expected results.

- 2 • integration testing is used to verify that a diagnostic integrates well into the ESMValTool
3 framework and that a diagnostic provides expected results. This is verified by comparison of
4 the results against a set of reference data generated during the implementation of the
5 diagnostic.

6 Each diagnostic is expected to produce a set of well-defined results, i.e. files in a variety of formats
7 and types (e.g., graphics, data files, ASCII files). While testing results of a diagnostic, a special
8 namelist file is executed by ESMValTool which runs a diagnostic on a limited set of test data only
9 minimizing executing time for testing while ensuring that the diagnostic produces the correct
10 results. The tests implemented include:

- 11 • file availability: a check that all required output data have been successfully generated by the
12 diagnostic. A missing file is always an indicator for a failure of the program.
- 13 • file checksum: currently the MD5 checksum is used to verify that contents of a file are the
14 same.
- 15 • graphics check: for graphic files an additional test is implemented which verifies that two
16 graphical outputs are identical. This is in particular useful to verify that outputs of a
17 diagnostic remain the same after code changes.

18 Unittests are implemented for each diagnostic independently using nose
19 (<https://nose.readthedocs.org/en/latest/>). Test files are searched recursively, executed and a statistic
20 on success and failures is provided at the end of the execution. In order to run integration tests for
21 each diagnostic, a small script needs to be written once. As for the unittests, a summary of success
22 and failures is provided as output (see the Supplement for details).

23 For the documentation of the code, Sphinx is used (<http://sphinx-doc.org/>) to organize and format
24 ESMValTool documentation, including text which has been extracted from source code. Sphinx can
25 help to create documentation in a variety of formats, including HTML, LaTeX (and hence printable
26 PDF), manual pages and plain text. Sphinx was originally developed for documenting Python code,
27 and one of its features is that it is able – using the so-called autodoc extension – to extract
28 documentation strings from Python source files and use them in the documentation it generates.
29 This feature apparently does not exist for NCL source files (such as those which are used in
30 ESMValTool), but it has been mimicked here via a Python script, which walks through a subset of
31 the ESMValTool NCL scripts, extracts function names, argument lists and descriptions (from the

comments immediately following the function definition), and assembles them in a subdirectory for usage with Sphinx. The documentation includes a listing of the functions, procedures, and plotting routines in order to encourage the reuse of existing code in multiple namelists.

3. Models and observations

The open-source release of ESMValTool (v1.0) that accompanies this paper is intended to work with CMIP5 model output, but the tool is compatible with any arbitrary model output, provided that it is in CF-compliant netCDF format and that the variables and metadata are following the CMOR tables and definitions. The namelists are designed such that it is straightforward to execute the same diagnostics with either CMIP DECK or CMIP6 model output rather than CMIP5 output, and these will be provided when the new simulations are available. As mentioned in the previous section, routines are provided for checking CF/CMOR compliance and fixing the most common minor flaws in the model output submitted to CMIP5. More substantial deviations from the required standards in the model output may be corrected via project- and model-specific procedures defined by the user and automatically applied within the workflow. The current reformatting routines are, however, not able to convert arbitrary model output to the full CF/CMOR standard. In this case, it is the responsibility of the individual modelling groups to perform that conversion. Currently, model-specific reformatting routines are provided for EMAC (Jöckel et al., 2015; Jöckel et al., 2010), the GFDL CM3 and ESM models (Donner et al., 2011; Dunne et al., 2012; Dunne et al., 2013), and for NEMO (Madec, 2008) which is the ocean model used in for example EC-Earth (Hazeleger et al., 2012). Users can develop similar reformatting routines specific to their model using the template included in the package allowing the tool to run directly on the original model output rather than having to reformat the model output to CF/CMOR beforehand.

The observations are organized in tiers. Where available, observations from the obs4MIPs and reanalysis from the ana4MIPs archives at the ESGF are used in the ESMValTool. These data sets form “Tier 1”. Tier 1 data are freely available for download to be directly used by the tool since they are formatted following the CF/CMOR standard and do not need any additional processing. For other observational data sets, the user has to retrieve the data from their respective source and reformat them into the CF/CMOR standard. To facilitate this task, we provide specific reformatting routines for a large number of such data sets together with detailed information of the data source, as well as download and processing instructions (see Table 1). “Tier 2” includes other freely available data sets and “Tier 3” includes restricted data sets (e.g., requiring the user to accept a

license agreement issued by the data owner). For Tier 2 and 3 data, links and help scripts are provided, so that these observations can be easily retrieved from their respective sources and processed by the user. A collection of all observational data used in ESMValTool (v1.0) is hosted at DLR and the ESGF nodes at BADC and DKRZ, but depending on the license terms of the observations these might not be publicly available.

4. Overview of namelists included in ESMValTool (v1.0)

A number of namelists have been included in ESMValTool (v1.0) that group a set of performance metrics and diagnostics for a given scientific topic. Namelists that focus on the evaluation of physical climate process for respectively, the atmosphere, ocean, and land surface are presented in Sections 4.1, 4.2, and 4.3. These can be applied to simulations with prescribed SSTs (i.e., AMIP runs) or the CMIP5 historical simulations (simulations for 1850 to present-day conducted with the best estimates of natural and anthropogenic climate forcing) that are run by either coupled AOGCMs or ESMs. Another set of namelists has been developed to evaluate biogeochemical biases present in ESMs when additional components of the Earth system such as the carbon cycle, atmospheric chemistry or aerosols are simulated interactively (Sections 4.4 and 4.5 for carbon cycle and aerosols/chemistry, respectively).

In each subsection, we first scientifically motivate the inclusion of the namelist by reviewing the main systematic biases in current ESMs and their importance and implications. We then give an overview of the namelists that can be used to evaluate such biases along with the diagnostics and performance metrics included, and the required variables and corresponding observations that are used in ESMValTool (v1.0). For each namelist we provide 1-2 example figures that are applied to either all or a subset of the CMIP5 models. An assessment of CMIP5 models is however not the focus of this paper. Rather, we attempt to illustrate how the namelists contained within ESMValTool (v1.0) can facilitate the development and evaluation of climate model performance in the targeted areas. Therefore, the results of each figure are only briefly described in each figure caption.

Table 1 provides a summary of all namelists included in ESMValTool (v1.0) along with information on the quantities and ESMValTool variable names for which the namelist is tested, the corresponding observations or reanalyses, the section and example figure in this paper, and references for the namelist. Table 2 then provides an overview of the diagnostics included for each

1 namelist along with specific calculations, the plot type, settings in the configuration file (cfg-file),
2 and comments.

3 **4.1. Detection of systematic biases in the physical climate: atmosphere**

4 **4.1.1. Quantitative performance metrics for atmospheric ECVs**

5 A starting point for the calculation of performance metrics is to assess the representation of
6 simulated climatological mean states and the seasonal cycle for essential climate variables (ECVs,
7 GCOS (2010)). This is supported by a large observational effort to deliver long-term, high quality
8 observations from different platforms and instruments (e.g., obs4MIPs and the ESA Climate
9 Change Initiative (CCI)) and ongoing efforts to improve global reanalysis products (e.g.,
10 ana4MIPs).

11 Following Gleckler et al. (2008) and similar to Fig. 9.7 of Flato et al. (2013), a namelist has been
12 implemented in the ESMValTool that produces a “portrait diagram” by calculating the relative
13 space-time root-mean square error (RMSE) from the climatological mean seasonal cycle of
14 historical simulations for selected variables [*namelist_perfmetrics_CMIP5.xml*]. In Fig. 2 the
15 relative space-time RMSE for the CMIP5 historical simulations (1980-2005) against a reference
16 observation and, where available, an alternative observational data set, is shown. The overall mean
17 bias can additionally be calculated and adding other statistical metrics is straightforward. Different
18 normalizations (mean, median, centered median) can be chosen and the multi model mean/median
19 can also be added. In order to calculate the RMSE, the data is regridded to a common grid using a
20 bilinear interpolation method. The user can select which grid to use as a target grid. The results
21 shown in this section have been obtained after regridding the data to the grid of the reference
22 dataset. With this namelist it is also possible to perform more in-depth analyses of the ECVs, by
23 calculating seasonal cycles, Taylor diagrams (Taylor, 2001), zonally averaged vertical profiles and
24 latitude-longitude maps. In the latter two cases, it is also possible to produce difference plots
25 between a given model and a reference (usually the observational data set) or between two versions
26 of the same model, and to apply a statistical test to highlight significant differences. As an example,
27 Fig. 3 (left panel) shows the zonal profile of seasonal mean temperature differences between the
28 MPI-ESM-LR model (Giorgetta et al., 2013) and ERA-Interim reanalysis (Dee et al., 2011), and
29 Fig. 3 (right panel) a Taylor diagram for temperature at 850 hPa for CMIP5 models compared to
30 ERA-Interim. A similar analysis can be performed with *namelist_righi15gmd_ECVs.xml*, which
31 reproduces the ECV plots of Righi et al. (2015) for a set of EMAC simulations.

1 Tested variables in ESMValTool (v1.0) that are shown in Fig. 2 are selected levels of temperature
2 (ta), eastward (ua) and northward wind (va), geopotential height (zg), and specific humidity (hus),
3 as well as near-surface air temperature (tas), precipitation (pr), all-sky longwave (rlut) and
4 shortwave (rsut) radiation, long-wave (LW_CRE) and shortwave (SW_CRE) cloud radiative effect,
5 and aerosol optical depth (AOD) at 550 nm (od550aer). The models are evaluated against a wide
6 range of observations and reanalysis data: ERA-Interim and NCEP (Kistler et al., 2001) for
7 temperature, winds and geopotential height, AIRS (Aumann et al., 2003) for specific humidity,
8 CERES-EBAF for radiation (Wielicki et al., 1996), Global Precipitation Climatology Project
9 (GPCP, Adler et al. (2003)) for precipitation, Moderate Resolution Imaging Spectrometer (MODIS,
10 Shi et al. (2011)) and the ESA CCI aerosol data (Kinne et al., 2015) for AOD. Additional
11 observations or reanalyses can be provided by the user for these variables and easily added. The
12 tool can also be applied to additional variables if the required observations are made available in an
13 ESMValTool compatible format (see Section 2 and Supplement).

14 **4.1.2. Multi-model mean bias for temperature and precipitation**

15 Near-surface air temperature (tas) and precipitation (pr) are the two variables most commonly
16 requested by users of ESM simulations. Often, diagnostics for tas and pr are shown for the multi-
17 model mean of an ensemble. Both of these variables are the end result of numerous interacting
18 processes in the models, making it challenging to understand and improve biases in these quantities.
19 For example, near surface air temperature biases depend on the models' representation of radiation,
20 convection, clouds, land characteristics, surface fluxes, as well as atmospheric circulation and
21 turbulent transport (Flato et al., 2013), each with their own potential biases that may either augment
22 or oppose one another.

23 The *namelist_flato13ipcc.xml* reproduces a subset of the figures from the climate model evaluation
24 chapter of IPCC AR5 (Chapter 9, Flato et al. (2013)). This namelist will be further developed and a
25 more complete version included in future releases. The diagnostic that calculates the multi-model
26 mean bias compared to a reference data set is part of this namelist and reproduces Figures 9.2 and
27 9.4 of Flato et al. (2013). Figure 4 shows the CMIP5 multi-model average as absolute values and as
28 biases relative to ERA-Interim and the GPCP data for the annual mean surface air temperature and
29 precipitation, respectively. Model output is regridded using bilinear interpolation to the reanalysis
30 or observational grid by default, but alternative options that can be set in the *cfg*-file include
31 regridding of the data to the lowest or highest resolution grid in the entire input data set. Such
32 figures can also be produced for individual seasons as well as for a single model simulation or other

1 2D variables if suitable observations are provided.

2 **4.1.3. Monsoon**

3 Monsoon systems represent the dominant seasonal climate variation in the tropics, with profound
4 socio-economic impacts. Current ESMs still struggle to capture the major features of both the South
5 Asian summer monsoon (SASM, Section 4.1.3.1) and the West African monsoon (WAM, Section
6 4.1.3.2). Sperber et al. (2013) and Roehrig et al. (2013) provide comprehensive assessments of the
7 ability of CMIP5 models to represent these two monsoon systems. By implementing diagnostics
8 from these two studies into ESMValTool (v1.0), we aim to facilitate continuous monitoring of
9 progress in simulating the SASM and WAM systems in ESMs.

10 **4.1.3.1. South Asian summer monsoon (SASM)**

11 While individual models vary in their simulations of the SASM, there are known biases in ESMs
12 that span a range of temporal and spatial scales. The namelists in the ESMValTool are targeted
13 toward analysing these biases in a systematic way. Climatological mean biases include excess
14 precipitation over the equatorial Indian Ocean, too little precipitation over the Indian subcontinent
15 and excess precipitation over orography such as the southern slopes of the Himalayas (Annamalai et
16 al., 2007; Bollasina and Nigam, 2009; Sperber et al., 2013), see also Fig. 4. The monsoon onset is
17 typically too late in the models, and the boreal summer intra-seasonal oscillation (BSISO), which
18 has a particularly large socio-economic impact in South Asia, is often weak or not present
19 (Sabeerali et al., 2013). Monsoon low pressure systems, which generate many of the most intense
20 rain events during the monsoon (Krishnamurthy and Misra, 2011) are often too infrequent and weak
21 (Stowasser et al., 2009). In coupled models, biases in SSTs, evaporation, precipitation and air-sea
22 coupling are common (Bollasina and Nigam, 2009) and have been shown to affect both present-day
23 simulations and future projections (Levine et al., 2013). Interannual teleconnections with ENSO
24 (Lin et al., 2008) and the Indian Ocean Dipole (Ashok et al., 2004; Cherchi and Navarra, 2013) are
25 also not well-captured (Turner et al., 2005).

26 Three SASM namelists for the basic climatology, seasonal cycle, intra-seasonal and inter-annual
27 variability and key teleconnections have been implemented into the ESMValTool focusing on
28 SASM rainfall and horizontal winds in June-September (JJAS) [*namelist_SAMonsoon.xml*,
29 *namelist_SAMonsoon_AMIP.xml*, *namelist_SAMonsoon_daily.xml*]. Rainfall and wind
30 climatologies, including their pattern correlations and RMSE against observations, are similar to the
31 metrics proposed by the Climate Variability and Predictability (CLIVAR) Asian–Australian
32 Monsoon Panel (AAMP) Diagnostics Task Team and used by Sperber et al. (2013). Diagnostics for

determining global monsoon domains and intensity follow the definition of Wang et al. (2012) where the global precipitation intensity is calculated from the difference between the hemispheric summer (May-September in the Northern Hemisphere, November-March in the Southern Hemisphere) and winter (vice versa) mean values, and the global monsoon domain is defined by those areas where the precipitation intensity exceeds 2.0 mm/day and the summer precipitation is $> 0.55 \times$ the annual precipitation (Fig. 5). Seasonal cycle diagnostics include monthly rainfall over the Indian region (5° - 30° N, 65° - 95° E) and dynamical indices based on wind-shear (Goswami et al., 1999; Wang and Fan, 1999; Webster and Yang, 1992). Figure 6 shows examples of the seasonal cycle of area-averaged Indian rainfall from selected CMIP5 models and their AMIP counterparts. The namelists include diagnostics to calculate maps of inter-annual standard deviation of JJAS rainfall and horizontal winds at 850 hPa and 200 hPa, and maps of teleconnection diagnostics between Nino3.4 SSTs (defined by the region 190° - 240° E, 5° S to 5° N) and JJAS precipitation across the monsoon region (30° S to 30° N, 40° - 300° E) following (Sperber et al., 2013). To generate difference maps, data are first regridded using an area-conservative binning and using the lowest resolution grid as target. For atmosphere-only models, we also evaluate their ability to represent year to year monsoon variability directly against time-equivalent observations to check whether models, given correct inter-annual SST forcing, can reproduce observed year to year variations and significant events occurring in particular years. This evaluation is done by plotting the time-series across specified years of standardized anomalies (normalized by climatology) of JJAS-averaged dynamical indices and area-averaged JJAS precipitation over the Indian region (defined above) for both the models and observations. Namelists for intra-seasonal variability include maps of standard deviation of 30-50 day filtered daily rainfall, with area-averaged values for key regions including the Bay of Bengal (10° - 20° N, 80° - 100° E) and the Eastern equatorial Indian Ocean (10° S- 10° N, 80° - 100° E) given in the plot titles. To illustrate the northward and eastward propagation of the BSISO, Hovmöller lag-longitude and lag-latitude diagrams show either the latitude-averaged (10° S- 10° N) and plotted for 60° - 160° E, or longitude-averaged (80° E- 100° E) and plotted for 10° S- 30° N, anomalies of 30-80 day filtered daily rainfall correlated against intraseasonal precipitation at the Indian Ocean reference point (75° E- 100° E, 10° S- 5° N). These use a slightly modified (for season, region and filtering band) version of the existing Madden-Julian Oscillation (MJO) NCL scripts, available at <https://www.ncl.ucar.edu/Applications/mjoclivar.shtml>, that are based on the recommendations from the US CLIVAR MJO Working Group (Waliser et al., 2009) and are similar to those shown in Lin et al. (2008) and used in Section 4.1.4.2 for the MJO.

Tested variables in ESMValTool (v1.0), some of which are illustrated in Figs. 5 and 6, include precipitation (pr), eastward (ua) and northward wind (va) at various levels, and skin temperature (ts). The primary reference data sets are ERA-Interim for horizontal winds, Tropical Rainfall Measuring Mission 3B43 version 7 (TRMM-3B43-v7; Huffman et al. (2007) for rainfall and HadISST (Rayner et al., 2003) for SST, although the models are evaluated against a wide range of other observational precipitation data sets (see Table 1) and an alternate reanalysis data set: the Modern-Era Retrospective Analysis for Research and Applications (MERRA; Rienecker et al. (2011)).

4.1.3.2. West African Monsoon Diagnostics

West Africa and the Sahel are highly dependent on seasonal rainfall associated with the WAM. Rainfall in the region exhibits strong inter-decadal variability (Nicholson et al., 2000), with major socio-economic impacts (Held et al., 2005). Projecting the future response of the WAM to increasing concentrations of greenhouse gases (GHG) is therefore of critical importance, as is the ability to make dependable forecasts of the WAM evolution on monthly to seasonal timescales. Current ESMs exhibit biases in their representation of both the mean state (Cook and Vizy, 2006; Roehrig et al., 2013) and temporal variability (Biasutti, 2013) of WAM. Such biases can affect the skill of monthly to seasonal predictions of the WAM as well as long term future projections. CMIP5 coupled models often exhibit warm SST biases in the equatorial Atlantic, which induce a southward shift of the WAM in summer (Richter et al., 2014). Because of the zonal symmetry, the 10°W-10°E meridional transect of any geophysical variable (see below) is particularly informative with respect to the main features of the WAM and their representation in climate models (Redelsperger et al., 2006). For instance, the JJAS-averaged Sahel rainfall has a large inter-model spread with biases ranging from +50% of the observed value (Cook and Vizy, 2006; Roehrig et al., 2013). Differences in simulated surface air temperatures are large over the Sahel and Sahara, with deficiencies in the Saharan heat low inducing feedback errors on the WAM structure. Here, a correct simulation of the surface energy balance is critical, where biases related to the representation of clouds, aerosols and surface albedo (Roehrig et al., 2013). The seasonal cycle also shows large inter-model spread, pointing to deficiencies in the representation of key processes important for the seasonal dynamics of the WAM. Daily precipitation is highly intermittent over the Sahel, mainly caused by a few intense mesoscale convective systems during the monsoon season (Mathon et al., 2002). Intense mesoscale convective systems over Africa as well as the diurnal cycle of the WAM are still a

challenge for most climate models (Roehrig et al., 2013). Improving the quality of the WAM in climate models is therefore urgently needed.

To evaluate key aspects of the WAM, two namelists have been implemented into ESMValTool (v1.0) [*namelist_WAMonsoon.xml*, *namelist_WAMonsoon_daily.xml*]. These include maps and meridional transects (averages over 10°W to 10°E) that provide a climatological picture of the summer (JJAS) WAM structure: (i) precipitation (pr) for the mean position of the WAM, (ii) near-surface air temperature (tas) for biases in the Atlantic cold tongue and the Saharan heat low, (iii) horizontal winds (ua, va) for the mean position and intensity of the monsoon flow at 925 hPa and of the mid- (700 hPa) and upper-level (200 hPa) jets. The surface and top of the atmosphere (TOA) radiation budgets provide a picture of the radiative fluxes associated with the WAM. Figure 7 shows the meridional transect of summer-averaged precipitation over West Africa for a range of CMIP5 models as an example for this namelist. Diagnostic for the mean seasonal cycle of precipitation is also provided to evaluate the WAM onset and withdrawal. Finally, a set of diagnostics for the WAM intra-seasonal variability evaluates the ability of models to capture variability of precipitation on timescales associated with African easterly waves (3-10 day), the MJO (25-90 days) and more broadly the WAM intra-seasonal variability (1-90 days). The strong day-to-day intermittency of precipitation is also diagnosed using maps of 1-day autocorrelation of intra-seasonal precipitation anomalies (Roehrig et al., 2013). To perform the autocorrelation analysis, data is first regridded to a common 1°×1° map using a bilinear interpolation method, whereas for generating difference maps the same regridding method as for the SASM diagnostics is used (see Section 4.1.3.1). Observations for evaluation are based on the following data sets: GPCP version 2.2 and Tropical Rainfall Measuring Mission 3B43 version 7 (TRMM-3B43-v7, Huffman et al. (2007)) precipitation retrievals, Clouds and Earth's Radiant Energy Systems (CERES) Energy Balanced and Filled (EBAF) edition 2.6 radiation estimates (Loeb et al., 2009), NOAA daily TOA outgoing longwave radiation (Liebmann and Smith, 1996), ERA-Interim reanalysis for the dynamics.

4.1.4. Natural modes of climate variability

4.1.4.1. NCAR Climate Variability Diagnostics Package

Modes of natural climate variability from interannual to multi-decadal time scales are important as they have large impacts on regional and even global climate with attendant socio-economic impacts. Characterization of internal (i.e., unforced) climate variability is also important for the detection and attribution of externally-forced climate change signals (Deser et al., 2012; Deser et al., 2014).

Internally-generated modes of variability also complicate model evaluation and intercomparison. As these modes are spontaneously generated, they do not need to exhibit the same chronological sequence in models as in nature. However, their statistical properties (e.g., time scale, autocorrelation, spectral characteristics, and spatial patterns) are captured to varying degrees of skill among climate models. Despite their importance, systematic evaluation of these modes remains a daunting task given the wide range to consider, the length of the data record needed to adequately characterize them, the importance of sub-surface oceanic processes and uncertainties in the observational records (Deser et al., 2010).

In order to assess natural modes of climate variability in models, the NCAR Climate Variability Diagnostics Package (CVDP) (Phillips et al., 2014) has been implemented into the ESMValTool. The CVDP has been developed as a standalone tool. To allow for easy updating of the CVDP once a new version is released, the structure of the CVDP is kept in its original form and a single namelist [*namelist_CVDP.xml*] has been written to enable the CVDP to be run directly within ESMValTool. The CVDP facilitates evaluation of the major modes of climate variability, including ENSO (Deser et al., 2010), PDO (Deser et al., 2010; Mantua et al., 1997), the Atlantic Multi-decadal Oscillation (AMO, Trenberth and Shea (2006)), the Atlantic Meridional Overturning Circulation (AMOC, Danabasoglu et al. (2012)), and atmospheric teleconnection patterns such as the Northern and Southern Annular Modes (NAM (Hurrell and Deser, 2009; Thompson and Wallace, 2000) and SAM (Thompson and Wallace, 2000), respectively), North Atlantic Oscillation (NAO, Hurrell and Deser (2009)), and Pacific North and South American (PNA and PSA, respectively (Thompson and Wallace, 2000)) patterns. For details on the actual calculation of these modes in CVDP we refer to the original CVDP package and explanations available at <http://www2.cesm.ucar.edu/working-groups/cvewg/cvdp>.

Depending on the climate mode analyzed, the CVDP package uses the following variables: precipitation (pr), sea level pressure (psl), near-surface air temperature (tas), skin temperature (ts), snow depth (snd), and basin-average ocean meridional overturning mass stream function (msftmyz). The models are evaluated against a wide range of observations and reanalysis data, for example NCEP for near-surface air temperature, HadISST for skin temperature, and the NOAA-CIRES Twentieth Century Reanalysis Project (Compo et al., 2011) for sea level pressure. Additional observations or reanalysis can be added by the user for these variables. The ESMValTool (v1.0) namelist runs on all CMIP5 models. As an example, Fig. 8 shows the representation of the PDO as

1 simulated by 41 CMIP5 models and observations (HadISST) and Fig. 9 the mean AMOC from 13
2 CMIP5 models.

3 **4.1.4.2. Madden-Julian oscillation (MJO)**

4 The MJO is the dominant mode of tropical intraseasonal variability (30-80 day) and has wide
5 impacts on numerous regional climate and weather phenomena (Madden and Julian, 1971).
6 Associated with enhanced convection in the tropics, the MJO exerts a significant influence on
7 monsoon precipitation, e.g. on the South Asian Monsoon (Pai et al., 2011) and on the west African
8 monsoon (Alaka and Maloney, 2012). The eastward propagation of the MJO into the West Pacific
9 can trigger the onset of some El Nino events (Feng et al., 2015; Hoell et al., 2014). The MJO also
10 influences tropical cyclogenesis in various ocean basins (Klotzbach, 2014). Increased vertical
11 resolution in the atmosphere and better and representation of stratospheric processes have led to an
12 improvement in MJO fidelity in CMIP5 compared with CMIP3 (Lin et al., 2006). However, current
13 generation models still struggle to adequately capture the eastward propagation of the MJO (Hung
14 et al., 2013) and the variance intensity is typically too weak. Identifying and reducing such biases
15 will be important for ESMs to accurately represent important climate phenomena, such as regional
16 precipitation variability in the tropics arising through the differing impact of MJO phases on ENSO
17 and ENSO forced regional climate anomalies (Hoell et al., 2014).

18 To assess the main MJO features in ESMs, a namelist with a number of diagnostics developed by
19 the US CLIVAR MJO Working Group (Kim et al., 2009; Waliser et al., 2009) has been
20 implemented in the ESMValTool (v1.0) [*namelist_mjo_mean_state.xml*, *namelist_mjo_daily.xml*].
21 These diagnostics are calculated using precipitation (pr), outgoing longwave radiation (OLR) (rlut),
22 eastward (ua) and northward wind (va) at 850 hPa (u850) and 200 hPa (u200) against various
23 observations and reanalysis data sets for boreal summer (May-October) and winter (November-
24 April).

25 Observation and reanalysis data sets include GPCP-1DD for precipitation, ERA-Interim and NCEP-
26 DOE reanalysis 2 for wind components (Kanamitsu et al., 2002) and NOAA polar-orbiting satellite
27 data for OLR (Liebmann and Smith, 1996). The majority of the scripts are based on example scripts
28 at <http://ncl.ucar.edu/Applications/mjoclivar.shtml>. Daily data is required for most of the scripts.
29 The basic diagnostics include mean seasonal state and 20-100 day bandpass filtered variance for
30 precipitation and u850 in summer and winter. To better assess and understand model biases in the
31 MJO, a number of more sophisticated diagnostics have also been implemented. These include;
32 univariate empirical orthogonal function (EOF) analysis for 20-100 day bandpass filtered daily

1 anomalies of precipitation, OLR, u850 and u200. To illustrate the northward and eastward
2 propagation of the MJO, lag-longitude and lag-latitude diagrams show either the equatorial
3 (latitude) averaged (10°S-10°N) or zonal (longitude) averaged (80°E-100°E) intraseasonal
4 precipitation anomalies and u850 anomalies correlated against intraseasonal precipitation at the
5 Indian Ocean reference point (75°E-100°E, 10°S-5°N). Similar figures can also be produced for
6 other key variables and regions following the definitions of Waliser et al. (2009). To further explore
7 the MJO intraseasonal variability, the wavenumber-frequency spectra for each season is calculated
8 for individual variables. In addition, we also produce cross-spectral plots to quantify the coherence
9 and phase relationships between precipitation and u850. Figure 10 shows examples of boreal
10 summer (May-October) wavenumber-frequency spectra of 10°S-10°N averaged daily precipitation
11 from GPCP-1DD, HadGEM2-ES, MPI-ESM-LR and EC-Earth. Finally, we also calculate the
12 multivariate combined EOF (CEOF) modes using equatorial averaged (15°S-15°N) daily anomalies
13 of U850, U200 and OLR. This analysis demonstrates the relationship between lower- and upper-
14 tropospheric wind anomalies and convection. To further illustrate the spatial-temporal structure of
15 the MJO, the first two leading CEOFs are used to derive a composite MJO life cycle which
16 highlights intraseasonal variability and northward/eastward propagation of the MJO. The data used
17 in these diagnostics are regridded to a common 0.5°×0.5° grid using an area-conservative method.

18 **4.1.5. Diurnal cycle**

19 In addition to the previously discussed biases in precipitation, many ESMs that rely on
20 parameterized convection exhibit biases related to the diurnal cycle and timing of precipitation.
21 Over land, ESMs tend to simulate a diurnal cycle of continental convective precipitation in phase
22 with insolation, while observed precipitation peaks in the early evening. This constitutes one of the
23 endemic biases of ESMs, in which convective precipitation intensity is often related to atmospheric
24 instability. This bias can have important implications for the simulated climate, as the timing of
25 precipitation influences subsequent surface evaporation, and convective clouds affect radiation
26 differently around noon or in late afternoon. The biases in the diurnal cycle are most pronounced
27 over land areas and the diurnal cycles of convection and clouds during the day contribute to the
28 continental warm bias (Cheruy et al., 2014). Similarly, biases in the diurnal cycle also exist over the
29 ocean (Jiang et al., 2015). Another motivation for looking at the diurnal cycle in models is that its
30 representation is more closely linked to the parameterizations of surface fluxes, boundary-layer,
31 convection and cloud processes than any other diagnostics. The phase of precipitation and radiative
32 fluxes during the day is the consequence of surface warming, boundary-layer turbulence mixing and

cumulus clouds moistening, as well as of the triggering criteria used to activate deep convection, and the closure used to compute convective intensity. The evaluation of the diurnal cycle thus provides a direct insight into the representation of physical processes in a model. Recent efforts to improve the representation of the diurnal cycle of precipitation models include modifying the convective entrainment rate, revisiting the quasi-equilibrium hypothesis for shallow and deep convection, and adding a representation of key missing processes such as boundary-layer thermals or cold pools. We envisage that ESMValTool will help to quantify the impact of those improvements in the next generation of ESMs.

To help document progress made in the representation of the diurnal cycle of precipitation (pr) in models, a set of diagnostics has been implemented in ESMValTool. After regridding all data on a common $2.5^{\circ} \times 2.5^{\circ}$ grid using bilinear interpolation, the mean diurnal cycle computed every 3 hours is approximated at each grid-point by a sum of sine and cosine functions (first harmonic analysis) allowing to derive global maps of the amplitude and phase of maximum rainfall over the day. Mean diurnal cycle of precipitation is also provided over specific regions in the tropics. Over land, we contrast semi-arid (Sahel) and humid (Amazonia) regions as well as West-Africa and India. Over the ocean, we focus on the Gulf of Guinea, the Indian Ocean and the East and West Equatorial Pacific. We use TRMM 3B42 V7, as a reference (http://mirador.gsfc.nasa.gov/collections/TRMM_3B42_daily_007.shtml). The ESMValTool also includes diagnostics for the evaluation of the diurnal cycle of radiative fluxes at the top of the atmosphere and at the surface, and their decomposition into LW and SW, total and clear sky components, however not all are available for all models from the CMIP5 archive. As a reference, we use 3-hourly SYN1deg CERES products (Wielicki et al., 1996), derived from measurements at top of the atmosphere and computed using a radiative transfer model at the surface (<http://ceres.larc.nasa.gov/products.php?product=SYN1deg>). These diagnostics provide a first insight into the representation of the diurnal cycle, but further analysis is required to understand the links between the model's parameterizations and the representation of the diurnal cycle, as well as the impact of errors in the diurnal cycle on other, slower timescale climate processes. Figure 11 shows the evaluation against TRMM observations of the mean diurnal cycle averaged over specific regions in the tropics for five summers (2004-2008) simulated by four CMIP5 ESMs.

1 **4.1.6. Clouds**

2 **4.1.6.1. Clouds and radiation**

3 Clouds are a key component of the climate system because of their large impact on the radiation
4 budget as well as their crucial role in the hydrological cycle. The simulation of clouds in climate
5 models has been challenging because of the many nonlinear processes involved (Boucher et al.,
6 2013). Simulations of long-term mean cloud properties from CMIP3 and CMIP5 models show large
7 biases compared to observations (Chen et al., 2011; Klein et al., 2013; Lauer and Hamilton, 2013).
8 Such biases have a range of implications as they affect application of these models to investigate
9 chemistry-climate interactions and aerosol-cloud interactions, while also having an impact on the
10 climate sensitivity of the model.

11 The namelist *namelist_lauer13jclim.xml* computes the climatology and interannual variability of
12 climate relevant cloud variables such as cloud radiative forcing, liquid and ice water path, and cloud
13 cover and reproduces the evaluation results of Lauer and Hamilton (2013). The standard namelist
14 includes a comparison of the geographical distribution of multi-year average cloud parameters from
15 individual models and the multi-model mean with satellite observations. Taylor diagrams are
16 generated that show the multi-year annual or seasonal average performance of individual models
17 and the multi-model mean in reproducing satellite observations. The diagnostic routine also
18 facilitates the assessment of the bias of the multi-model mean and zonal averages of individual
19 models compared with satellite observations. Interannual variability is estimated as the relative
20 temporal standard deviation from multi-year timeseries of data with the temporal standard
21 deviations calculated from monthly anomalies after subtracting the climatological mean seasonal
22 cycle. Data regridding is applied using a bilinear interpolation method and choosing the grid of the
23 reference dataset as target. As an example, Fig. 12 shows the bias of the 20-year average (1985-
24 2005) annual mean cloud radiative effects from CMIP5 models (multi-model mean) against the
25 CERES EBAF satellite climatology (2001-2012) (Loeb et al., 2012; Loeb et al., 2009), similar to
26 Flato et al. (2013) their Figure 9.5.

27 The cloud namelist focuses on precipitation (pr) and four cloud parameters that largely determine
28 the impact of clouds on the radiation budget and thus climate in the model simulations: total cloud
29 amount (clt), liquid water path (lwp), ice water path (iwp), and TOA cloud radiative effect (CRE)
30 consisting of the longwave CRE and shortwave CRE that can also separately be evaluated with the
31 performance metrics namelist (see Section 4.1.1). Precipitation is evaluated with GPCP data, total
32 cloud amount with MODIS, liquid water path with passive-microwave satellite observations from

the University of Wisconsin (O'Dell et al., 2008), and the ice water path with MODIS Cloud Model Intercomparison Project (MODIS-CFMIP, Pincus et al. (2012), King et al. (2003)) data.

4.1.6.2. Quantitative performance assessment of cloud regimes

The cloud-climate radiative feedback process remains one of the largest sources of uncertainty in determining the climate sensitivity of models (Boucher et al., 2013). Traditionally, clouds have been evaluated in terms of their impact on the mean top of atmosphere fluxes. However, it is possible to achieve good performance on these quantities through compensating errors, for example boundary layer clouds may be too reflective but have insufficient horizontal coverage (Nam et al., 2012). Williams and Webb (2009) proposed a Cloud Regime Error Metric (CREM) which critically tests the ability of a model to simulate both the relative frequency of occurrence and the radiative properties correctly for a set of cloud regimes determined by the daily mean cloud top pressure, in-cloud albedo and fractional coverage at each grid-box. Having previously identified the regimes by clustering joint cloud-top pressure-optical depth histograms from the International Satellite Cloud Climatology Project (ISCCP, Rossow and Schiffer (1999)) as per Williams and Webb (2009), each daily model grid box is assigned to the regime cluster centroid with the closest cloud top pressure, in-cloud albedo and fractional coverage as determined by the 3-element Euclidean distance. The fraction of grid points assigned to each of the regimes and the mean radiative properties of those grid points are then compared to the observed values. This routine also uses a bilinear regridding method with a $2.5^{\circ} \times 2.5^{\circ}$ target grid.

This metric is now implemented in ESMValTool (v1.0), with references in the code to tables in the Williams and Webb (2009) study defining the cluster centroids [*namelist_williams09climdyn_CREM.xml*]. Required are daily data from ISCCP mean cloud albedo (albiscpp), ISCCP Mean Cloud Top Pressure (pctisccp), ISCCP Total Total Cloud Fraction (cltisccp), TOA outgoing short- and long-wave radiation (rsut, rlut), TOA outgoing shortwave radiation (rlutes), surface snow area fraction (snc) or surface snow amount (snw), and sea ice area fraction (sic). The metric has been applied over the period January 1985 to December 1987 to those CMIP5 models with the required diagnostics (daily data) available for their AMIP simulation (see caption of Fig. 13). A perfect score with respect to ISCCP would be zero. Williams and Webb (2009) also compared data from the MODIS and the Earth Radiation Budget Experiment (ERBE, Barkstrom (1984)) to ISCCP in order to provide an estimate of observational uncertainty. This observational regime characteristic was found to be 0.96 as marked on Fig. 13 when calculated over the period March 1985 to February 1990. Hence a model with a score that is similar to this value

can be considered to be within observational uncertainty, although it should be noted that this does not necessarily mean that the model lies within the observations for each regime. Error bars are not plotted since experience has shown that the metric has little sensitivity to interannual variability and models that are visibly different on Fig. 13 are likely to be significantly so. A minimum of two years, and ideally five years or more, of daily data are required for the scientific analysis.

4.2. Detection of systematic biases in the physical climate: ocean

4.2.1. Handling of ocean grids

Analysis of ocean model data from ESMs poses several unique challenges for analysis. First, in order to avoid numerical singularities in their calculations, ocean models often use irregular grids where the poles have been rotated or moved to be located over land areas. For example, the global configuration of the Nucleus for European Modelling of the Ocean (NEMO) framework uses a tripolar grid (Madec, 2008), with the three poles located over Siberia, Canada and Antarctica. Second, transports of scalar quantities (e.g., overturning stream functions and heat transports) can only be calculated accurately on the original model grids as interpolation to other grids introduces errors. This means that, e.g. for the calculation of water transport through a strait, both the horizontal and vertical extent of the grids on which the u and v currents are defined is required. Therefore, this type of diagnostic can only be used for models for which all native grid information is available. State variables like SSTs, sea ice and salinity are regridded using grid information (i.e., coordinates, bounds, and cell areas) available in the ocean input files of the CMIP5 models. To create difference plots against observations or other models all data are regridded to a common grid (e.g., $1^\circ \times 1^\circ$) using the regridding functionality of the Earth System Modeling Framework (ESMF, <https://www.ncl.ucar.edu/Applications/ESMF.shtml>).

4.2.2. Southern Ocean Diagnostics

4.2.2.1. Southern Ocean mixed layer dynamics and surface turbulent fluxes

Earth system models often show large biases in the Southern Ocean mixed layer. For example, Sterl et al. (2012) showed that in EC-Earth/NEMO the Southern Ocean is too warm and salinity too low, while the mixed-layer is too shallow. These biases are not specific to EC-Earth, but are rather widespread. At the same time, values for Antarctic Circumpolar Current (ACC) transport vary between 90 and 264 Sv in CMIP5 models, with a mean of 155 ± 51 Sv. The differences are associated with differences in the ACC density structure.

1 A namelist has been implemented in the ESMValTool to analyse these biases
 2 [*namelist_SouthernOcean.xml*]. With these diagnostics polar stereographic (difference) maps can be
 3 produced to compare monthly/annual mean model fields with corresponding ERA-Interim data. The
 4 patch recovery technique is applied to regrid data to a common $1^{\circ} \times 1^{\circ}$ grid. There are also scripts to
 5 plot the differences in the area mean vertical profiles of ocean temperature and salinity between
 6 models and data from the World Ocean Atlas (Antonov et al., 2010; Locarnini et al., 2010). The
 7 ocean mixed layer thickness from models can be compared with that obtained from the Argo floats
 8 (Dong et al., 2008). Finally, the ACC strength, as measured by water mass transport through the
 9 Drake Passage, is calculated using the same method as in the CDFTOOLS package (CDFTOOLS,
 10 <http://servforge.legi.grenoble-inp.fr/projects/CDFTOOL>). This diagnostic can be used to calculate
 11 the transport through other sections as well, but is presently only available for NEMO/ORCA1
 12 output, for which all grid information is available. The required variables for the comparison with
 13 ERA-Interim are sea surface temperature (tos), downward heat flux (hfds, calculated from ERA-
 14 Interim by summing the surface latent and sensible heat flux and the net shortwave and longwave
 15 fluxes (hfsls+hfss+rsns+rlns)), water flux (wfpe, calculated by summing precipitation and
 16 evaporation (pr+evspsbl)) and the wind stress components (tauu and tauv). For the comparison with
 17 the World Ocean Atlas 2009 data (WOA09) sea surface salinity (sos), sea water salinity (so) and
 18 temperature (to) are required variables. For the comparison with the Argo floats the ocean mixed
 19 layer thickness (mldst) is required. Finally the two components of sea water velocity (uo and vo)
 20 are required for the volume transport calculation. Some example figures from this set of diagnostic
 21 scripts are shown for EC-Earth in Fig. 14.

22 **4.2.2.2. Atmospheric processes forcing the Southern Ocean**

23 One leading cause of SST biases in the Southern Ocean is systematic biases in surface radiation
 24 fluxes (Trenberth and Fasullo, 2010) coupled with systematic errors in macrophysical (e.g. cloud
 25 amount) and microphysical (e.g. frequency of mixed-phase clouds) cloud properties (Bodas-Salcedo
 26 et al., 2014).

27 A namelist has been implemented into the ESMValTool that compares model estimates of cloud,
 28 radiation and surface turbulent flux variables over the Southern Ocean with suitable observations
 29 [*namelist_SouthernHemisphere.xml*]. Due to the lack of surface/in-situ observations over the
 30 Southern Ocean, remotely sensed data can be subject to considerable uncertainty (Mace, 2010).
 31 While this uncertainty is not explicitly addressed in ESMValTool (v1.0), in future releases we will
 32 include a number of alternative satellite based data sets for cloud variables (e.g., MISR, MODIS,

1 ISCCP) as well as new methods under development to derive surface turbulent flux estimates
2 constrained by observed TOA radiation flux estimates and atmospheric energy divergence derived
3 from reanalysis products (Trenberth and Fasullo, 2008). Inclusion of multiple satellite-based
4 estimates will provide some estimate of observational uncertainty over the region. Variables
5 analysed include (i) total cloud cover (clt), vertically integrated cloud liquid water and cloud ice
6 water (clwvi, clivi) (ii) surface/ (TOA) downward/outgoing total sky and clear sky short wave and
7 longwave radiation fluxes (rsds, rsdcs, rlds, rldscs / rsut, rsutcs, rlut, rlutcs) and (iii) surface
8 turbulent latent and sensible heat fluxes (hfls, hfss). Observational constraints are derived from,
9 respectively; cloud: CloudSat level 3 data (Stephens et al., 2002), radiation: CERES-EBAF level 3
10 Ed2 data and surface turbulent fluxes: WHOI-OAflux (Yu et al., 2008).

11 The following diagnostics are calculated with accompanying plots: (i) Seasonal mean absolute-
12 value and difference maps for model data versus observations covering the Southern Ocean region
13 (30°S-65°S) for all variables. (ii) Mean seasonal cycles using zonal means averaged separately over
14 three latitude bands (i) 30°S-65°S, the entire Southern Ocean, (ii) 30°S-45°S, the sub-tropical
15 Southern Ocean and (iii) 45°S-65°S, the mid-latitude Southern Ocean. (iii) Annual means of each
16 variable (models and observations) plotted as zonal means, over 30°S-65°S, (iv) Scatter plots of
17 seasonal mean downward (surface) and outgoing (TOA) longwave and short wave radiation as a
18 function of total cloud cover, cloud liquid water path or cloud ice water path, calculated for the
19 three regions outlined above. The data are regridded using a cubic interpolation method with the
20 observations grid as target. Figure 15 provides an example diagnostic, with the top panel showing
21 covariability of seasonal mean surface downward short wave radiation as a function of total cloud
22 cover. To construct the figure, grid point values of cloud cover, for each season covering 30°S to
23 65°S, are saved into bins of 5% increasing cloud cover. For each grid point the corresponding
24 seasonal mean radiation value is used to obtain a mean radiation flux for each cloud cover bin. The
25 lower panel plots the fractional occurrence of seasonal mean cloud cover from CloudSat and model
26 data for the same spatial and temporal averaging as used in the upper panel. Observations from
27 CERES-EBAF radiation plotted against CloudSat cloud cover are compared to an example CMIP5
28 model. From the covariability plot we can diagnose whether models exhibit a similar dependency
29 between incoming surface short wave radiation and cloud cover as seen in observations. We can
30 further assess if there is a systematic bias in surface solar radiation and whether this bias occurs at
31 specific values of cloud cover. Similar covariability plots are available for surface incoming
32 longwave radiation and for TOA long and short wave radiation, plotted respectively against cloud
33 cover, cloud liquid water path and cloud ice water path. Combining these diagnostics provides a

comprehensive evaluation of simulated relationships between surface and TOA radiation fluxes and cloud variables.

4.2.3. Simulated tropical ocean climatology

An accurate representation of the tropical climate is fundamental for ESMs. The majority of solar energy received by the Earth is in the tropics and the potential for thermal emission of absorbed energy back to space is also largest in the tropics due to the high column concentrations of water vapor at low latitudes (Pierrehumbert, 1995; Stephens and Greenwald, 1991). Coupled interactions between equatorial SSTs, surface wind stress, precipitation and upper-ocean mixing are central to many tropical biases in ESMs. This is the case both with respect to the mean state and for key modes of variability, influenced by, or interacting with, the mean state (e.g., El Nino Southern Oscillation (ENSO), Choi et al. (2011)). Such biases are often reflected in a “double ITCZ” seen in the majority of CMIP3 and CMIP5 CCMs (Li and Xie, 2014; Oueslati and Bellon, 2015). The double ITCZ bias, present in many ESMs, occurs when models fail to simulate a single, year round, ITCZ rainfall maximum north of the equator. Instead, an unrealistic secondary maximum in models south of the equator is present for some or all of the year. Such biases are particularly prevalent in the tropical Pacific, but can also occur in the Atlantic (Oueslati and Bellon, 2015). This double ITCZ is often accompanied by an overextension of the East Pacific equatorial cold tongue into the Central Pacific, collocated with a positive bias in easterly near-surface wind speeds and a shallow bias in ocean mixed layer depth (Lin, 2007). Such biases can directly impact the ability of an ESM to accurately represent ENSO variability (An et al., 2010; Guilyardi, 2006) and its potential sensitivity to climate change (Chen et al., 2015), with negative consequences for a range of simulated features, such as regional tropical temperature and precipitation variability, monsoon dynamics and ocean and terrestrial carbon uptake (Iguchi, 2011; Jones et al., 2001).

To assess such tropical biases with the ESMValTool, we have implemented a namelist with diagnostics motivated by the work of Li and Xie (2014) [*namelist_TropicalVariability.xml*]. In particular, we reproduce their Fig. 5 for models and observations/reanalyses, calculating equatorial mean (5°N-5°S), longitudinal sections of annual mean precipitation (pr), skin temperature (ts), horizontal winds (ua and va) and 925 hPa divergence (derived from the sum of the partial derivatives of the wind components extracted at the 925 hPa pressure level (that is $du/dx + dv/dy$). Latitude cross sections of the model variables are plotted for the equatorial Pacific, Indian and Atlantic Oceans with observational constraints provided by the TRMM-3B43-v7 for precipitation, the HadISST for SSTs, and ERA-interim reanalysis for temperature and winds. Latitudinal sections

1 of absolute and normalized annual mean SST and precipitation are also calculated, spatially
2 averaged for the three ocean basins. Normalization follows the procedure outlined in Fig. 1 of Li
3 and Xie (2014) whereby values at each latitude are normalized by the tropical mean (20°N-20°S)
4 value of the corresponding parameter (e.g., annual mean precipitation at a given location is divided
5 by the 20°N-20°S annual mean value). Finally, to assess how models capture observed relationships
6 between SST and precipitation we calculate the co-variability of precipitation against SST for
7 specific regions of the tropical Pacific. This analysis includes calculation of the Mean Square Error
8 (MSE) between model SST/precipitation and observational equivalents. A similar regridding
9 procedure as for the Southern Hemisphere diagnostics is applied here, based on a cubic
10 interpolation method and using the observations as target grid. The namelist as included in
11 ESMValTool (v1.0) runs on all CMIP5 models. Figure 16 provides one example of the tropical
12 climate diagnostics, with latitude cross sections of absolute and tropical normalized SST and
13 precipitation from three CMIP5 models (HadGEM2-ES (Collins et al., 2011), MPI-ESM-LR and
14 IPSL-CM5A-MR (Dufresne et al., 2013)) plotted against HadISST and TRMM data.

15 **4.2.4. Sea ice**

16 Sea ice is a key component of the climate system through its effects on radiation and seawater
17 density. A reduction in sea ice area results in increased absorption of shortwave radiation, which
18 warms the sea ice region and contributes to further sea ice loss. This process is often referred to as
19 the sea ice albedo climate feedback which is part of the Arctic amplification phenomena (Curry,
20 2007). CMIP5 models tend to underestimate the decline in summer Arctic sea ice extent observed
21 by satellites during the last decades (Stroeve et al., 2012) which may be related to models'
22 underestimation of the sea ice albedo feedback process (Boé et al., 2009). Conversely in the
23 Antarctic, observations show a small increase in March sea ice extent while the CMIP5 models
24 simulate a small decrease (Flato et al., 2013; Stroeve et al., 2012). It is therefore important that
25 model sea-ice processes are evaluated and improvements regularly assessed. Caveats have been
26 noted with respect to the limitations of using only sea ice extent as a metric of model performance
27 (Notz et al., 2013) as the sea ice concentration, volume, and drift, sea ice thickness and surface
28 albedo, as well as sea ice processes such as melt pond formation or the summer sea ice melt are all
29 important sea ice related quantities. In addition the atmospheric forcings (e.g., wind, clouds, and
30 snow) and ocean forcings (e.g., salinity and ocean transport) impact on the sea ice state and
31 evolution.

32 In ESMValTool (v1.0) the sea ice namelist includes diagnostics that cover sea ice extent and

1 concentration [*namelist_SeaIce.xml*], but work is underway to include other variables and processes
2 in future releases. An example diagnostic produced by the sea ice namelist is given in Figure 17,
3 which shows the timeseries of September Arctic sea ice extent from the CMIP5 historical
4 simulations compared to observations from the National Snow and Ice Data Center (NSIDC)
5 produced by combining concentration estimates created with the NASA Team algorithm and the
6 Bootstrap algorithm (Meier et al., 2013; Peng et al., 2013) and SSTs from the HadISST data set,
7 similar to Figure 9.24 of Flato et al. (2013). Sea ice extent is calculated as the total area (km²) of
8 grid cells over the Arctic or Antarctic with sea-ice concentrations (sic) of at least 15%. The sea ice
9 namelist can also calculate the seasonal cycle of sea ice extent and polar stereographic contour and
10 polar contour difference plots of Arctic and Antarctic sea ice concentration. For the latter
11 diagnostic, data is regridded to a common 1°×1° grid using the patch recovery technique.

12 **4.3. Detection of systematic biases in the physical climate: land**

13 **4.3.1. Continental dry bias**

14 The representation of land surface processes and fluxes in climate models critically affects the
15 simulation of near-surface climate over land. In particular, energy partitioning at the surface
16 strongly influences surface temperature and it has been suggested that temperature biases in ESMs
17 can be in part related to biases in evapotranspiration. The most notable feature in a majority of
18 CMIP3 and CMIP5 models is a tendency to overestimate evapotranspiration globally (Mueller and
19 Seneviratne, 2014).

20 A diagnostic to analyse the representation of evapotranspiration in ESMs has been included in the
21 ESMValTool [*namelist_Evapotranspiration.xml*]. For comparison with the LandFlux-EVAL
22 product (Mueller et al., 2013), the modelled surface latent heat flux (hfls) is converted to
23 evapotranspiration units using the latent heat of vaporization. The diagnostic then produces lat-lon
24 maps of absolute evapotranspiration as well as bias maps (model minus reference product, after
25 regridding data to the coarsest grid using area-conservative interpolation). In Fig. 18, the global
26 pattern of monthly mean evapotranspiration is evaluated against the LandFlux-EVAL product. The
27 evapotranspiration diagnostic is complemented by the Standardized Precipitation Index (SPI)
28 diagnostic [*namelist_SPI.xml*], which gives a measure of drought intensity from an atmospheric
29 perspective and can help relating biases in evapotranspiration to atmospheric causes such as the
30 accumulated precipitation amounts. For each month, precipitation (pr) is summed over the
31 preceding months (options for 3, 6 or 12-monthly SPI). Then a two-parameter Gamma distribution

1 of cumulative probability is fitted to the strictly positive month sums, such that the probability of a
2 non-zero precipitation sum being below a certain value x corresponds to $\text{Gamma}(x)$. The shape and
3 scale parameters of the gamma distribution are estimated with a maximum likelihood approach.
4 Accounting for periods of no precipitation, occurring at a frequency q , the total cumulative
5 probability distribution of a precipitation sum below x , $H(x)$, becomes $H(x) = q + (1 -$
6 $q) * \text{Gamma}(x)$. In the last step, a precipitation sum x is assigned to its corresponding SPI value by
7 computing the quantile $q_N(0,1)$ of the standard normal distribution at probability $H(x)$. The SPI of
8 a precipitation sum x , thus, corresponds to the quantile of the standard normal distribution which is
9 assigned by preserving the probability of the original precipitation sum, $H(x)$. Mean and annual
10 cycle are not meaningful since the SPI accounts for seasonality and transforms the data to a zero
11 average in each month. Therefore the diagnostic focuses on lat-lon maps of annual or seasonal
12 trends in SPI (unitless) when comparing models with observations.

13 **4.3.2. Runoff**

14 Evaluation of precipitation is a challenge due to potentially large errors and uncertainty in observed
15 precipitation data (Biemans et al., 2009; Legates and Willmott, 1990). An alternative or additional
16 option to the direct evaluation of precipitation over land (such as, e.g., included in the global
17 precipitation evaluation in Sect. 4.1.2) is the evaluation of river runoff that can in principle be
18 measured with comparatively small errors for most rivers. Routine measurements are performed for
19 many large rivers, generating a large global database (e.g., available at the Global Runoff Data
20 Centre (GRDC, Dümenil Gates et al. (2000)). The length of available time series, however, varies
21 between the rivers, with large data gaps especially in recent years for many rivers. The evaluation of
22 runoff against river gauge data can provide a useful independent measure of the simulated
23 hydrological cycle. If both river flow and precipitation are given with reasonable accuracy, it will
24 also provide an observational constraint on model surface evaporation, provided that the considered
25 averaging time periods are long enough so that changes in surface water storages are negligible
26 (Hagemann et al., 2013), e.g., by considering climatological means of 20 years or more. For present
27 climate conditions ESMs often exhibit a dry and warm near-surface bias during summer over mid-
28 latitude continents (Hagemann et al., 2004). Continental dry biases in precipitation exist in the
29 majority of CMIP5 models over South America, the Mid-west of US, the Mediterranean region,
30 Central and Eastern Europe, West and South Asia (Fig. 4 and Fig. 9.4 of Flato et al. (2013)). These
31 precipitation biases often transfer into dry biases in runoff, but sometimes dry biases in runoff can
32 be caused by a too large evapotranspiration (Hagemann et al., 2013). In order to relate biases in

runoff to biases in precipitation and evapotranspiration, the catchment oriented evaluation in this section considers biases in all three variables. This means that the respective variables are considered as spatially averages over the drainage basins of large rivers.

Beside bias maps, a set of diagnostics to produce basin-scale comparisons of runoff (mrro), evapotranspiration (evspsbl) and precipitation (pr) have also been implemented in ESMValTool [*namelist_runoff_et.xml*]. This namelist calculates biases in climatological annual means of the three variables for 12 large-scale catchments areas on different continents and for different climates. For total runoff, catchment averaged model values are compared to climatological long-term averages of GRDC observations. Due to the incompleteness of these station data, a year-to-year correspondence of data cannot be achieved so only climatological data are considered, as in Hagemann et al. (2013). Simulated precipitation is compared to catchment-averaged WATCH forcing data based on ERA-Interim (WFDEI) data (Weedon et al., 2014) for the period 1979-2010. Evapotranspiration observations are estimated using the difference of the catchment-averaged WFDEI precipitation minus the climatological GRDC river runoff. As an example, Fig. 19 shows biases in runoff coefficient (runoff/precipitation) against the relative precipitation bias for the historical simulation of one of the CMIP5 models (MPI-ESM-LR).

4.4. Detection of biogeochemical biases: carbon cycle

4.4.1. Terrestrial biogeochemistry

A realistic representation of the global carbon cycle is a fundamental requirement for ESMs. In the past, climate models were directly forced by atmospheric CO₂ concentrations, but since CMIP5, ESMs are routinely forced by anthropogenic CO₂ emissions, the atmospheric concentration being inferred from the difference between these emissions and the ESM simulated land and ocean carbon sinks. These sinks are affected by atmospheric CO₂ and climate change, inducing feedbacks between the climate system and the carbon cycle (Arora et al., 2013; Friedlingstein et al., 2006). Quantification of these feedbacks is critical to estimate the future of these carbon sinks and hence atmospheric CO₂ and climate change (Friedlingstein et al., 2014).

The diagnostics implemented in ESMValTool to evaluate simulated terrestrial biogeochemistry are based on the study of Anav et al. (2013) and span several time-scales: climatological means, intra-annual (seasonal cycle), interannual and long-term trends [*namelist_anav13jclim.xml*]. Further extending these routines, the diagnostics presented in Sect. 4.1.1 are also applied here to calculate quantitative performance metrics. These metrics assess how both the land and ocean

biogeochemical components of ESMs reproduce different aspects of the land and ocean carbon cycle, with an emphasis on variables controlling the exchange of carbon between the atmosphere and these two reservoirs. The analysis indicates some level of compensating errors within the models. Selecting, within the namelist, several specific diagnostics to be applied to more key variables controlling the land or ocean carbon cycle, can help reducing the risk of missing such compensating errors. Figure 20 shows a portrait diagram similar to Fig. 3 of Anav et al. (2013) but for seasonal carbon cycle metrics against suitable reference data sets (see below).

For land, diagnostics of the land carbon sink net biosphere productivity (nbp) are essential. Although direct observations are not available, nbp can be estimated from atmospheric CO₂ inversions (JMA and TRANSCOM) and on the global scale combined with observation-based estimates of the oceanic carbon sink (fgco2 from GCP (Le Quéré et al., 2014)). In addition to net carbon fluxes, diagnostics for gross primary productivity of land (gpp), leaf area index (lai), vegetation (cVeg) and soil carbon pools (cSoil) are also implemented in the ESMValTool to assess possible error compensation in ESMs. Observation-based gpp estimates are derived from Model Tree Ensemble (MTE) upscaling data (Jung et al., 2009) from the network of eddy-covariance flux towers (FLUXNET, Beer et al. (2010)). The leaf area index data set used for evaluation (LAI3g) is derived from the Global Inventory Modeling and Mapping Studies group (GIMMS) AVHRR normalized difference vegetation index (NDVI-017b) data (Zhu et al., 2013). Finally, cSoil and cVeg are assessed as mean annual values over different large sub-domains using the Harmonised World soil Database (HWSD, Nachtergaele et al. (2012)) and the Olson based vegetation carbon data set (Gibbs, 2006; Olson et al., 1985).

4.4.2. Marine biogeochemistry

Marine biogeochemistry models form a core component of ESMs and require evaluation for multiple passive tracers. The increasing availability of quality-controlled global biogeochemical data sets for the historical period (e.g. Surface Ocean CO₂ Atlas Version 2 (SOCAT v2, Bakker et al. (2014)) provides further opportunity to evaluate model performance on multi-decadal timescales. Recent analyses of CMIP5 ESMs indicate that persistent biases exist in simulated biogeochemical variables, for instance as identified in ocean oxygen (Andrews et al., 2013) and carbon cycle (Anav et al., 2013) fields derived from CMIP5 historical experiments. Some systematic biases in biogeochemical tracers can be attributed to physical deficiencies within ocean models (see Section 4.2), motivating further understanding of coupled physical-biogeochemical processes in the current generation of ESMs. For example, erroneous over oxygenation of subsurface waters within the

MPI-ESM-LR CMIP5 model has been attributed to excess ventilation and vertical mixing in mid- to high-latitude regions (Ilyina et al., 2013).

A namelist is provided that includes diagnostics to support the evaluation of ocean biogeochemical cycles at global scales, as simulated by both ocean-only and coupled climate-carbon cycle ESMs [*namelist_GlobalOcean.xml*]. Supported input variables include surface partial pressure of CO₂ (*spco2*), surface chlorophyll concentration (*chl*), surface total alkalinity (*talk*) and dissolved oxygen concentration (*o2*). These variables provide an integrated view of model skill with regard to reproducing bulk marine ecosystem and carbon cycle properties. Observation-based reference data sets include SOCAT v2 and ETH-SOM-FFN (Landschützer et al., 2014a, b) for surface *p*CO₂, Sea-viewing Wide Field-of-view Sensor (SeaWiFS) satellite data for surface chlorophyll (McClain et al., 1998), climatological data for total alkalinity (Takahashi et al., 2014), and World Ocean Atlas 2005 climatological data (WOA05) with in situ corrections following Bianchi et al. (2012) for dissolved oxygen. Diagnostics calculate contour plots for climatological distributions, inter-annual or inter-seasonal (e.g. JJAS) variability together with the difference between each model and a chosen reference data set. Such differences are calculated after regridding the data to the coarsest grid using an area-conservative interpolation. Monthly, seasonal or annual frequency time-series plots can also be produced either globally averaged or for a selected latitude-longitude range. Optional extensions include the ability to mask model data with the same coverage as observations, calculate anomaly fields, and to overlay trend lines, and running or multi-model means. Pre-processing routines are also included to accommodate native curvilinear grids, common in ocean model discretisation (see Section 4.2.1), along with providing the ability to extract depth levels from 3-D input fields. An example plot is presented in Fig. 22, showing inter-annual variability in surface ocean *p*CO₂ as simulated by a subset of CMIP5 ESMs (BNU-ESM, HadGEM2-ES, GFDL-ESM2M), expressed as the standard deviation of de-trended annual averages for the period 1992 – 2005. As an observation-based reference *p*CO₂ field, ETH SOM-FFN (1998-2011) is used, which extrapolates SOCAT v2 data (Bakker et al., 2014) using a 2-step neural network method. As described in Landschützer et al. (2014a), ETH SOM-FFN partitions monthly SOCAT v2 *p*CO₂ observations into discrete biogeochemical provinces by establishing common relationships between independent input parameters using a Self Organising Map (SOM). Non-linear input-target relationships, as derived for each biogeochemical province using a Feed-Forward Network (FFN) method, are then used to extrapolate observed *p*CO₂.

1 A diagnostic for oceanic Net Primary Production (NPP) is also implemented in ESMValTool for
2 climatological annual mean and seasonal cycle, as well as for inter-annual variability over the 1986-
3 2005 period [*namelist_anav13jclim.xml*]. Observations are derived from the SeaWiFS satellite
4 chlorophyll data, using the Vertically Generalized Production Model (VGPM, Behrenfeld and
5 Falkowski (1997)).

6 **4.5. Detection of biogeochemical biases: aerosols and trace gas chemistry**

7 **4.5.1. Tropospheric aerosols**

8 Tropospheric aerosols play a key role in the Earth system and have a strong influence on climate
9 and air pollution. The global aerosol distribution is characterized by a large spatial and temporal
10 variability which makes its representation in ESMs particularly challenging (Ghan and Schwartz,
11 2007). In addition, aerosol interactions with radiation (direct aerosol effect (Schulz et al., 2006))
12 and with clouds (indirect aerosol effects (Lohmann and Feichter, 2005)) need to be accounted for.
13 Model-based estimates of anthropogenic aerosol effects are still affected by large uncertainties,
14 mostly due to an incorrect representation of aerosol processes (Kinne et al., 2006). Myhre et al.
15 (2013) report a substantial spread in simulated aerosol direct effects among 16 global aerosol
16 models and attribute it to diversities in aerosol burden, aerosol optical properties and aerosol optical
17 depth (AOD). Diversities in black carbon (BC) burden up to a factor of three, related to model
18 disagreements in simulating deposition processes were also found by Lee et al. (2013). Model
19 meteorology can be a source of diversity since it impacts on atmospheric transport and aerosol
20 lifetime. This in turn relates to the simulated essential climate variables such as winds, humidity and
21 precipitation (see Section 4.1). Large biases also exist in simulated aerosol indirect effects (IPCC,
22 2013) and are often a result of systematic errors in both model aerosol and cloud fields (see Section
23 4.1.6).

24 To assess current biases in global aerosol models, the aerosol namelist of the ESMValTool
25 comprises several diagnostics to compare simulated aerosol concentrations and optical depth at the
26 surface against station data, motivated by the work of Pringle et al. (2010), Pozzer et al. (2012), and
27 Righi et al. (2013) [*namelist_aerosol_CMIP5.xml*]. Diagnostics include time series of monthly or
28 yearly mean aerosol concentrations, scatter plots with the relevant statistical indicators, and contour
29 maps directly comparing model results against observations. The comparison is performed
30 considering collocated model and observations in space and time. In the current version of
31 ESMValTool, these diagnostics are supplied with observational data from a wide range of station

1 networks, including Interagency Monitoring of Protected Visual Environments (IMPROVE) and
2 CASTNET (North America), European Monitoring and Evaluation Programme (EMEP, Europe)
3 and the recently-established Asian network (EANET). The AERONET data are also available for
4 evaluating aerosol optical depth in continental regions and in a few remote marine locations. For
5 evaluating aerosol optical depth, we also use satellite data, the primary advantage of which is
6 almost-global coverage, particularly over the oceans. Satellite data is however affected by
7 uncertainties related to the algorithm used to process radiances into relevant geophysical state
8 variables. The tool currently implements data from the Multi-angle Imaging SpectroRadiometer
9 (MISR, Stevens and Schwartz (2012)), MODIS and the ESACCI-AEROSOL product (Kinne et al.,
10 2015) which is a combination of ERS2-ATSR2 and ENVISAT-AATSR data. To calculate model
11 biases against satellite data, regridding is performed using a bilinear interpolation to the coarsest
12 grid. Aerosol optical depth time series over the ocean for the period 1850-2010 are shown in Fig. 23
13 for the CMIP5 models in comparison to MODIS and ESACCI-AEROSOL. Finally, more specific
14 aerosol diagnostics have been implemented to compare aerosol vertical profiles of mass and number
15 concentrations and aerosol size distributions, based on the evaluation work by Lauer et al. (2005)
16 and Aquila et al. (2011). These diagnostics, however, use model quantities that were not part of the
17 CMIP5 data request and therefore will not be discussed here.

18 **4.5.2. Tropospheric trace gas chemistry and stratospheric ozone**

19 In the past, climate models were forced with prescribed tropospheric and stratospheric ozone
20 concentration, but since CMIP5 some ESMs include interactive chemistry and are capable of
21 representing prognostic ozone (Eyring et al., 2013; Flato et al., 2013). This allows models to
22 simulate important chemistry-climate interactions and feedback processes. Examples include the
23 increase in oxidation rates in a warmer climate which leads to decreases in methane and its lifetime
24 (Voulgarakis et al., 2013) or the increase in tropical upwelling (associated with the Brewer Dobson
25 circulation) in a warmer climate and corresponding reductions in tropical lower stratospheric ozone
26 as a result of faster transport and less time for ozone production (Butchart et al., 2010; Eyring et al.,
27 2010). It is thus becoming important to evaluate the simulated atmospheric composition in ESMs. A
28 common high bias in the Northern Hemisphere and a low bias in the Southern Hemisphere has been
29 identified in tropospheric column ozone simulated by chemistry-climate models participating in the
30 Atmospheric Chemistry Climate Model Intercomparison Project (ACCMIP), which could partly be
31 related to deficiencies in the ozone precursor emissions (Young et al., 2013). Analysis of CMIP5
32 models with respect to trends in total column ozone show that the multi-model mean of the models

with interactive chemistry is in good agreement with observations, but that significant deviations exist for individual models (Eyring et al., 2013; Flato et al., 2013). Large variations in stratospheric ozone in models with interactive chemistry drive large variations in lower stratospheric temperature trends. The results show that both ozone recovery and the rate of GHG increase determine future Southern Hemisphere summer-time circulation changes and are important to consider in ESMs (Eyring et al., 2013).

The namelists implemented in the ESMValTool to evaluate atmospheric chemistry can reproduce the analysis of tropospheric ozone and precursors of Righi et al. (2015) [*namelist_righi15gmd_tropo3.xml*, *namelist_righi15gmd_Emmons.xml*] and the study by Eyring et al. (2013) [*namelist_eyring13jgr.xml*]. The calculation of the RMSE, mean bias, and Taylor diagrams (see Section 4.1.1) has been extended to tropospheric column ozone (derived from tro3 fields), ozone profiles (tro3) at selected levels, and surface carbon monoxide (vmrco) (see Righi et al. (2015) for details). This enables a consistent calculation of relative performance for the climate parameters and ozone, which is particularly relevant given that biases in climate can impact on biases in chemistry and vice versa. In addition, diagnostics that evaluate tropospheric ozone and its precursors (nitrogen oxides (vmrnox), ethylene (vmrc2h4), ethane (vmrc2h6), propene (vmrc3h6), propane (vmrc3h8) and acetone (vmrch3coch3)) are compared to the observational data of Emmons et al. (2000). A diagnostic to compare tropospheric column ozone from the CMIP5 historical simulations to Aura MLS/OMI observations (Ziemke et al., 2011) is also included and shown as an example in Fig. 24. This diagnostic also remaps the data to the coarsest grid using local area averaging in order to calculate differences. For the stratosphere, total column ozone (toz) diagnostics are implemented. As an example, Figure 25 shows the CMIP5 total column ozone time series compared to the NIWA combined total column ozone database (Bodeker et al., 2005).

4.6. Linking model performance to projections

The relatively new research field of emergent constraints aims to link model performance evaluation with future projection feedbacks. An emergent constraint refers to the use of observations to constrain a simulated future Earth system feedback. It is referred to as emergent, because a relationship between a simulated future projection feedback and an observable element of climate variability emerges from an ensemble of ESM projections, potentially providing a constraint on the future feedback. Emergent constraints can help focus model development and evaluation onto processes underpinning uncertainty in the magnitude and spread of future Earth system change. Systematic model biases in certain forced modes, such as the seasonal cycle of

snow cover or inter-annual variability of tropical land CO₂ uptake appear to project in an understandable way onto the spread of future climate change feedbacks resulting from these phenomena (Cox et al., 2013; Hall and Qu, 2006; Wenzel et al., 2014).

To reproduce the analysis of Wenzel et al. (2014) that provides an emergent constraint on future tropical land carbon uptake, a namelist is included into ESMValTool (v1.0) to perform an emergent constraint analysis of the carbon cycle-climate feedback parameter (γ_{LT}) (Cox et al., 2013; Friedlingstein et al., 2006) [*namelist_wenzel14jgr.xml*]. This namelist only considers the CMIP5 ESMs that have provided the necessary output for the analysis. This criterion precludes most CMIP5 models and only seven ESMs are therefore considered here. The namelist includes diagnostics which analyse the short-term sensitivity of atmospheric CO₂ to temperature variability on interannual time scales (γ_{IAV}) for models and observations, as well as diagnostics for γ_{LT} from the models. The observed sensitivity γ_{IAV} is calculated by summing land (nbp) and ocean (fgco2) carbon fluxes which are correlated to tropical near-surface air temperature (tas). Results from historical model simulations are compared to observational based estimates of carbon fluxes from the Global Carbon project (GCP, Le Quéré et al. (2014)) and reanalysis temperature data from the NOAA National Climate Data Center (NCDC, Smith et al. (2008)). For diagnosing γ_{LT} from the models, nbp from idealized fully coupled and biochemically coupled simulations are used as well as tas from fully coupled idealized simulations (see Fig. 26). Emergent constraints of this type help to understand some of the underlying processes controlling future projection sensitivity and offer a promising approach to reduce uncertainty in multi-model climate projections.

5. Use of the ESMValTool in the model development cycle and evaluation workflow

5.1. Model development

As new model versions are developed, standardized diagnostics suites as presented here allow model developers to compare their results against previous versions of the same model or against other models, e.g. CMIP5 models. Such analyses help to identify different aspects in a model that have either improved or degraded as a result of a particular model development. The benchmarking of ESMs using performance metrics (see Section 4.1.1) provides an overall picture of the quality of the simulation, whereas process-oriented diagnostics help determine whether the simulation quality improvements are for the correct underlying physical reasons and point to paths for further model improvement.

1 The ESMValTool is intended to support modelling centres with quality control of their CMIP
2 DECK experiments and the CMIP6 historical simulation, as well as other experiments from
3 CMIP6-Endorsed Model Intercomparison Projects (Eyring et al., 2015). A significant amount of
4 institutional resources go into running, post-processing, and publishing model results from such
5 experiments. It is important that centres can easily identify and correct potential errors in this
6 process. The standardized analyses contained in the ESMValTool can be used to monitor the
7 progress of CMIP experiments. While the tool is designed to accommodate a wide range of time
8 axes and configurations, and many of the diagnostics may be run on control or future climate
9 experiments, ESMValTool (v1.0) is largely targeted to evaluate AMIP and the CMIP historical
10 simulations.

11 **5.2. Integration into modelling workflows**

12 The ESMValTool can be run as a stand-alone tool, or integrated into existing modelling workflows.
13 The primary challenge is to provide CF/CMOR compliant data. Not all modelling centres produce
14 CF/CMOR compliant data directly as part of their workflow although we note that more are doing
15 so as the potential benefits are being realized. For many groups conversion to CF/CMOR standards
16 involves significant post-processing of native model output. This may require some groups to
17 perform analysis via the ESMValTool on their model output after conversion to CF/CMOR, or to
18 create intermediate “CMOR-like” versions of the data. Users who wish to use native model output
19 can take advantage of the reformatting routine flexibility (see Section 3) to create scripts that
20 convert this data into the CF/CMOR standard. As an example, reformat scripts for the NOAA-
21 GFDL models and the EMAC model are included with the initial release. These scripts are used to
22 convert the native model output for direct use with the ESMValTool. The reformatting routine
23 capability may provide an alternative to more expensive and complete “CMORization” processes
24 that are usually required to formally publish model data on the ESGF.

25 **5.3. Running the ESMValTool alongside the ESGF**

26 Large international model inter-comparison projects such as CMIP stimulated the development of a
27 globally distributed federation of data providers, supporting common data provisioning policies and
28 infrastructures. ESGF is an international open source effort to establish a distributed data and
29 computing platform, enabling world wide access to Peta- (in the future Exa-) byte scale scientific
30 climate data. Data can be searched via a globally distributed search index with access possible via
31 HTTP, OpenDAP and GridFTP. To efficiently run the ESMValTool on CMIP model data and

1 observations alongside the ESGF, the necessary data hosted by the ESGF has to be made locally
2 accessible at the site where ESMValTool is executed. There are various ways this might be
3 achieved. One possibility is to run ESMValTool separately at each site holding datasets required by
4 the analysis, then combine the results. However, this is limited by the extent to which calculations
5 can be performed without requiring data from another site. A more practical possibility is running
6 ESMValTool alongside a large store of replica datasets gathered from across the ESGF, so that all
7 the required data are in one location. Certain large ESGF sites (e.g., DKRZ, BADC, IPSL, PCMDI)
8 provide replica dataset stores, and ESMValTool has been run in such a way at several of these sites.

9 Replica dataset stores do not provide a complete solution however, as it is impossible to replicate all
10 ESGF datasets at one site, so circumstances will arise when one or more required datasets are not
11 available locally. The obvious solution is to download these datasets from elsewhere in the ESGF,
12 and store them locally whilst the analysis is carried out. The indexed search facility provided by the
13 ESGF makes it easy to identify the download URL of such ‘remote’ datasets, and a prototype of
14 ESMValTool (not included in v1.0) has been developed that performs this search automatically
15 using `esgf-pyclient`¹. If the search is successful, the prototype provides the user with the URL of
16 each file in the dataset, and the user (or system administrator) is then responsible for performing the
17 download. The workflow of this prototype is illustrated in Figure 27. It is possible that the fully
18 automated downloading of remote ESGF datasets may be provided by a future version of
19 ESMValTool, but for now it is preferable for a human to manage the process due to large size of the
20 files involved. A more complete coupling to the ESGF was originally planned for version 1.0 but
21 was not possible due to the long down period of the ESGF.

22 23 **6. Summary and Outlook**

24 The Earth System Model Evaluation Tool (ESMValTool) is a diagnostics package for routine
25 evaluation of Earth System Models (ESMs) with observations and reanalyses data or for
26 comparison with results from other models. The ESMValTool has been developed to facilitate the
27 evaluation of complex ESMs at individual modelling centres and to help streamline model
28 evaluation standards within CMIP. Priorities to date that are included in ESMValTool (v1.0)
29 described in this paper concentrate on selected systematic biases that were a focus of the European

¹ <https://pypi.python.org/pypi/esgf-pyclient>

1 Commission's 7th Framework Programme “Earth system Model Bias Reduction and assessing
2 Abrupt Climate change (EMBRACE) project, the DLR Earth System Model Evaluation (ESMVal)
3 project and other collaborative projects, in particular: performance metrics for selected ECVs,
4 coupled tropical climate variability, monsoons, Southern Ocean processes, continental dry biases
5 and soil hydrology-climate interactions, atmospheric CO₂ budgets, ozone, and tropospheric aerosol.
6 We have applied the bulk of the diagnostics of ESMValTool (v1.0) to the entire set of CMIP5
7 historical or AMIP simulations. The namelist on emergent constraints for the carbon cycle has been
8 additionally applied to idealized carbon cycle experiments and the emission driven RCP 8.5
9 simulations.

10 ESMValTool (v1.0) can be used to compare new model simulations against CMIP5 models and
11 observations for the selected scientific themes much faster than this was possible before. Model
12 groups, who wish to do this comparison before submitting their CMIP6 historical simulations or
13 AMIP experiments to the ESGF can do so since the tool is provided as open source software. In
14 order to run the tool locally, observations need to be downloaded and for tiers 2 and 3 reformatted
15 with the help of the reformatting scripts that are included. Model output needs to be either in CF
16 compliant NetCDF or a reformatting routine needs to be written by the modelling group, following
17 given examples for EMAC, GFDL models, and NEMO.

18 Users of the ESMValTool (v1.0) results need to be aware that ESMValTool (v1.0) only includes a
19 subset of the wide behaviour of model performance that the community aims to characterize. The
20 results of running the ESMValTool need to be interpreted accordingly. Over time, the ESMValTool
21 will be extended with additional diagnostics and performance metrics. A particular focus will be to
22 integrate additional diagnostics that can reproduce the analysis of the climate model evaluation
23 chapter of IPCC AR5 (Flato et al., 2013) as well as the projection chapter (Collins et al., 2013). We
24 will also extend the tool with diagnostics to quantify forcings and feedbacks in the CMIP6
25 simulations and to calculate metrics such as the equilibrium climate sensitivity (ECS), transient
26 climate response (TCR), and the transient climate response to cumulative carbon emissions (TCRE)
27 (IPCC, 2013). While inclusion of these diagnostics is straightforward, the evaluation of processes
28 and phenomena to improve understanding about the sources of errors and uncertainties in models
29 that we also plan to enhance remains a scientific challenge. The field of emergent constraints
30 remains in its infancy and more research is required how to better link model performance to
31 projections (Flato et al., 2013). In addition, an improved consideration of the interdependency in the

1 evaluation of a multi-model ensemble (Sanderson et al., 2015a, b) as well as internal variability in
2 ESM evaluation is required.

3 A critical aspect in ESM evaluation is the availability of consistent, error-characterized global and
4 regional Earth observations, as well as accurate globally gridded reanalyses that are constrained by
5 assimilated observations. Additional or longer records of observations and reanalyses will be used
6 as they become available, with a focus on using obs4MIPs - including new contributions from the
7 European Space Agency's Climate Change Initiative (ESA CCI) - and ana4MIPs data. The
8 ESMValTool can consider observational uncertainty in different ways, e.g. through the use of more
9 than one observational data set to directly evaluate the models, by showing the difference between
10 the reference data set and the alternative observations, or by including an observed uncertainty
11 ensemble that spans the observed uncertainty range (e.g., available for the surface temperature data
12 set compiled for HadISST). Often the uncertainties in the observations are not readily available.
13 Reliable and robust error characterization/estimation of observations is a high priority throughout
14 the community, and obs4MIPs and other efforts that create data sets for model evaluation should
15 encourage the inclusion of such uncertainty estimates as part of each data set.

16 The ESMValTool will be contributed to the analysis code catalogue being developed by the
17 WGNE/WGCM climate model metrics panel. The purpose of this catalogue is to make the diversity
18 of existing community-based analysis capabilities more accessible and transparent, and ultimately
19 for developing solutions to ensure they can be readily applied to the CMIP DECK and the CMIP6
20 historical simulation in a coordinated way. We are currently exploring options to interface with
21 complimentary efforts, e.g. the PCMDI Metrics Package (PMP, Gleckler et al. (2016)) and the
22 Auto-Assess package that is under development at the UK Met Office. An international strategy for
23 organising and presenting CMIP results produced by various diagnostic tools is needed, and this
24 will be a priority for the WGNE/WGCM climate metrics panel in collaboration with the CMIP
25 Panel (<http://www.wcrp-climate.org/index.php/wgcm-cmip/about-cmip>).

26 This paper presents ESMValTool (v1.0) which allows users to repeat all the analyses shown.
27 Additional updates and improvements will be included in subsequent versions of the software,
28 which are planned to be released on a regular basis. The ESMValTool works on CMIP5 simulations
29 and, given CMIP DECK and CMIP6 simulations will be in a similar format, it will be
30 straightforward to run the package on these simulations. A limiting factor at present is the need to
31 download all data to a local cache. This limitation has spurred the development allowing
32 ESMValTool to run alongside the ESGF at one of the data nodes. An initial attempt to couple the

1 tool to the ESGF has been made, but this is still at prototype stage (see Section 5.3). An additional
2 limiting factor is that the model output from all CMIP models has to be mirrored to the ESGF data
3 node where the tool is installed. This is facilitated by providing a listing of the variables and time
4 frequencies that are used in ESMValTool (v1.0) which uses a significantly smaller volume than the
5 data request for the CMIP DECK and CMIP6 simulations will include. This reduced set of data
6 could be mirrored with priority.

7 Several technical improvements are required to make the software package more efficient. One
8 current limitation is the lack of a parallelization. Given the huge amount of data involved in a
9 typical CMIP analysis, this can be highly CPU-time-intensive when performed on a single
10 processor. In future releases, the possibility of parallelizing the tool will be explored. Additional
11 development work is ongoing to create a more flexible pre-processing framework, which will
12 include operations like ensemble-averaging and regridding to the current reformatting procedures as
13 well as an improved coupling to the ESGF. Here, future versions of the ESMValTool will build as
14 much as possible on existing efforts for the backend that reads and reformats data. In this regard it
15 would be helpful if an application programming interface (API) could be defined for example by
16 the WGCM Infrastructure Panel (WIP) that allows for flexible integration of diagnostics across
17 different tools and programming languages in CMIP to this backend.

18 We aim to move ESM evaluation beyond the state-of-the-art by investing in operational evaluation
19 of physical and biogeochemical aspects of ESMs, process-oriented evaluation and by identifying
20 processes most important to the magnitude and uncertainty of future projections. Our goal is to
21 support model evaluation in CMIP6 by contributing the ESMValTool as one of the standard
22 documentation functions and by running it alongside the ESGF. In collaboration with similar
23 efforts, we aim for a routine evaluation that provides a comprehensive documentation of broad
24 aspects of model performance and its evolution over time and to make evaluation results available
25 at a timescale that was not possible in CMIP5. This routine evaluation is not meant to replace
26 further in-depth analysis of model performance and can to date not strongly reduce uncertainties in
27 global climate sensitivity which remains an active area of research. However, the ability to routinely
28 perform such evaluation will drive the quality and realism of ESMs forward and will leave more
29 time to develop innovative process-oriented diagnostics - especially those related to feedbacks in
30 the climate system that link to the credibility of model projections.

1 **7. Code availability**

2 ESMValTool (v1.0) is released under the Apache License, VERSION 2.0. The latest version of the
3 ESMValTool is available from the ESMValTool webpage at <http://www.esmvaltool.org/>. Users
4 who apply the Software resulting in presentations or papers are kindly asked to cite this paper
5 alongside with the Software doi (doi:10.17874/ac8548f0315) and version number. In addition,
6 ESMValTool will be further developed in a version controlled repository that is accessible only to
7 the development team. Regular releases are planned for the future. The wider climate community is
8 encouraged to contribute to this effort and to join the ESMValTool development team for
9 contribution of additional more in-depth diagnostics for ESM evaluation. A wiki page for the
10 development that describes ongoing developments is also available. Interested users and developers
11 are welcome to contact the lead author.

13 **Acknowledgements**

14 The development of ESMValTool (v1.0) was funded by the European Commission's 7th
15 Framework Programme, under Grant Agreement number 282672, the “Earth system Model Bias
16 Reduction and assessing Abrupt Climate change (EMBRACE)” project and the DLR “Earth System
17 Model Validation (ESMVal)” and “Klimarelevanz von atmosphärischen Spurengasen, Aerosolen
18 und Wolken: Auf dem Weg zu EarthCARE und MERLIN (KliSAW)” projects. In addition,
19 financial support for the development of ESMValTool (v1.0) was provided by ESA’s Climate
20 Change Initiative Climate Modelling User Group (CMUG). We acknowledge the World Climate
21 Research Program’s (WCRP’s) Working Group on Coupled Modelling (WGCM), which is
22 responsible for CMIP, and we thank the climate modelling groups for producing and making
23 available their model output. For CMIP the U.S. Department of Energy's Program for Climate
24 Model Diagnosis and Intercomparison provides coordinating support and led development of
25 software infrastructure in partnership with the Global Organization for Earth System Science
26 Portals. We thank Björn Brötz (DLR, Germany) for his help with the release of the ESMValTool
27 and Clare Enright (UEA, UK) for support with development of the ocean biogeochemistry
28 diagnostics. We are grateful to Patrick Jöckel (DLR, Germany), Ron Stouffer (GFDL, USA) and to
29 the two anonymous referees for their constructive comments on the manuscript.

30

31

1 **References**

- 2 Adler, R. F., Huffman, G. J., Chang, A., Ferraro, R., Xie, P.-P., Janowiak, J., Rudolf, B., Schneider,
3 U., Curtis, S., Bolvin, D., Gruber, A., Susskind, J., Arkin, P., and Nelkin, E.: The Version-2 Global
4 Precipitation Climatology Project (GPCP) Monthly Precipitation Analysis (1979–Present), *J*
5 *Hydrometeorol*, 4, 1147-1167, 2003.
- 6 Alaka, G. J. and Maloney, E. D.: The Influence of the MJO on Upstream Precursors to African
7 Easterly Waves, *J Climate*, 25, 3219-3236, 2012.
- 8 An, S. I., Ham, Y. G., Kug, J. S., Timmermann, A., Choi, J., and Kang, I. S.: The Inverse Effect of
9 Annual-Mean State and Annual-Cycle Changes on ENSO, *J Climate*, 23, 1095-1110, 2010.
- 10 Anav, A., Friedlingstein, P., Kidston, M., Bopp, L., Ciais, P., Cox, P., Jones, C., Jung, M., Myneni,
11 R., and Zhu, Z.: Evaluating the Land and Ocean Components of the Global Carbon Cycle in the
12 CMIP5 Earth System Models, *J Climate*, 26, 6801-6843, 2013.
- 13 Andrews, O. D., Bindoff, N. L., Halloran, P. R., Ilyina, T., and Le Quere, C.: Detecting an external
14 influence on recent changes in oceanic oxygen using an optimal fingerprinting method,
15 *Biogeosciences*, 10, 1799-1813, 2013.
- 16 Annamalai, H., Hamilton, K., and Sperber, K. R.: The South Asian summer monsoon and its
17 relationship with ENSO in the IPCC AR4 simulations, *J Climate*, 20, 1071-1092, 2007.
- 18 Antonov, J. I., Seidov, D., Boyer, T. P., Locarnini, R. A., Mishonov, A. V., Garcia, H. E.,
19 Baranova, O. K., Zweng, M. M., and Johnson, D. R.: World Ocean Atlas 2009, Volume 2: Salinity.
20 In: NOAA Atlas NESDIS 69, Levitus, S. (Ed.), U.S. Government Printing Office, Washington,
21 D.C., 2010.
- 22 Aquila, V., Hendricks, J., Lauer, A., Riemer, N., Vogel, H., Baumgardner, D., Minikin, A., Petzold,
23 A., Schwarz, J. P., Spackman, J. R., Weinzierl, B., Righi, M., and Dall'Amico, M.: MADE-in: a
24 new aerosol microphysics submodel for global simulation of insoluble particles and their mixing
25 state, *Geosci Model Dev*, 4, 325-355, 2011.
- 26 Arora, V. K., Boer, G. J., Friedlingstein, P., Eby, M., Jones, C. D., Christian, J. R., Bonan, G.,
27 Bopp, L., Brovkin, V., Cadule, P., Hajima, T., Ilyina, T., Lindsay, K., Tjiputra, J. F., and Wu, T.:
28 Carbon-Concentration and Carbon-Climate Feedbacks in CMIP5 Earth System Models, *J Climate*,
29 26, 5289-5314, 2013.
- 30 Ashok, K., Guan, Z. Y., Saji, N. H., and Yamagata, T.: Individual and combined influences of
31 ENSO and the Indian Ocean Dipole on the Indian summer monsoon, *J Climate*, 17, 3141-3155,
32 2004.
- 33 Aumann, H. H., Chahine, M. T., Gautier, C., Goldberg, M. D., Kalnay, E., McMillin, L. M.,
34 Revercomb, H., Rosenkranz, P. W., Smith, W. L., Staelin, D. H., Strow, L. L., and Susskind, J.:
35 AIRS/AMSU/HSB on the Aqua mission: design, science objectives, data products and processing
36 system, *EEE Trans. Geosci. and Remote Sensing*, 41, 253-264, 2003.
- 37 Bakker, D. C. E., Pfeil, B., Smith, K., Hankin, S., Olsen, A., Alin, S. R., Cosca, C., Harasawa, S.,
38 Kozyr, A., Nojiri, Y., O'Brien, K. M., Schuster, U., Telszewski, M., Tilbrook, B., Wada, C., Akl, J.,
39 Barbero, L., Bates, N. R., Boutin, J., Bozec, Y., Cai, W. J., Castle, R. D., Chavez, F. P., Chen, L.,
40 Chierici, M., Currie, K., de Baar, H. J. W., Evans, W., Feely, R. A., Fransson, A., Gao, Z., Hales,
41 B., Hardman-Mountford, N. J., Hoppema, M., Huang, W. J., Hunt, C. W., Huss, B., Ichikawa, T.,
42 Johannessen, T., Jones, E. M., Jones, S. D., Jutterström, S., Kitidis, V., Körtzinger, A.,
43 Landschützer, P., Lauvset, S. K., Lefèvre, N., Manke, A. B., Mathis, J. T., Merlivat, L., Metzl, N.,
44 Murata, A., Newberger, T., Omar, A. M., Ono, T., Park, G. H., Paterson, K., Pierrot, D., Ríos, A. F.,

1 Sabine, C. L., Saito, S., Salisbury, J., Sarma, V. V. S. S., Schlitzer, R., Sieger, R., Skjelvan, I.,
2 Steinhoff, T., Sullivan, K. F., Sun, H., Sutton, A. J., Suzuki, T., Sweeney, C., Takahashi, T.,
3 Tjiputra, J., Tsurushima, N., van Heuven, S. M. A. C., Vandemark, D., Vlahos, P., Wallace, D. W.
4 R., Wanninkhof, R., and Watson, A. J.: An update to the Surface Ocean CO₂ Atlas (SOCAT
5 version 2), *Earth Syst. Sci. Data*, 6, 69-90, 2014.

6 Barkstrom, B. R.: The Earth Radiation Budget Experiment (ERBE), *B Am Meteorol Soc*, 65, 1170-
7 1185, 1984.

8 Becker, A., Finger, P., Meyer-Christoffer, A., Rudolf, B., Schamm, K., Schneider, U., and Ziese,
9 M.: A description of the global land-surface precipitation data products of the Global Precipitation
10 Climatology Centre with sample applications including centennial (trend) analysis from 1901–
11 present, *Earth Syst. Sci. Data*, 5, 71-99, 2013.

12 Beer, C., Reichstein, M., Tomelleri, E., Ciais, P., Jung, M., Carvalhais, N., Rodenbeck, C., Arain,
13 M. A., Baldocchi, D., Bonan, G. B., Bondeau, A., Cescatti, A., Lasslop, G., Lindroth, A., Lomas,
14 M., Luyssaert, S., Margolis, H., Oleson, K. W., Rouspard, O., Veenendaal, E., Viovy, N., Williams,
15 C., Woodward, F. I., and Papale, D.: Terrestrial Gross Carbon Dioxide Uptake: Global Distribution
16 and Covariation with Climate, *Science*, 329, 834-838, 2010.

17 Behrenfeld, M. J. and Falkowski, P. G.: Photosynthetic rates derived from satellite-based
18 chlorophyll concentration, *Limnol Oceanogr*, 42, 1-20, 1997.

19 Bianchi, D., Dunne, J. P., Sarmiento, J. L., and Galbraith, E. D.: Data-based estimates of suboxia,
20 denitrification, and N₂O production in the ocean and their sensitivities to dissolved O₂, *Global*
21 *Biogeochem Cy*, 26, 2012.

22 Biasutti, M.: Forced Sahel rainfall trends in the CMIP5 archive, *Journal of Geophysical Research:*
23 *Atmospheres*, 118, 1613-1623, 2013.

24 Biemans, H., Hutjes, R. W. A., Kabat, P., Strengers, B. J., Gerten, D., and Rost, S.: Effects of
25 Precipitation Uncertainty on Discharge Calculations for Main River Basins, *J Hydrometeorol*, 10,
26 1011-1025, 2009.

27 Bodas-Salcedo, A., Williams, K. D., Ringer, M. A., Beau, I., Cole, J. N. S., Dufresne, J. L.,
28 Koshiro, T., Stevens, B., Wang, Z., and Yokohata, T.: Origins of the Solar Radiation Biases over
29 the Southern Ocean in CFMIP2 Models, *J Climate*, 27, 41-56, 2014.

30 Bodeker, G. E., Shiona, H., and Eskes, H.: Indicators of Antarctic ozone depletion, *Atmos Chem*
31 *Phys*, 5, 2603-2615, 2005.

32 Boé, J., Hall, A., and Qu, X.: Current GCMs' Unrealistic Negative Feedback in the Arctic, *J*
33 *Climate*, 22, 4682-4695, 2009.

34 Bollasina, M. and Nigam, S.: Indian Ocean SST, evaporation, and precipitation during the South
35 Asian summer monsoon in IPCC-AR4 coupled simulations, *Clim Dynam*, 33, 1017-1032, 2009.

36 Bollasina, M. A. and Ming, Y.: The general circulation model precipitation bias over the
37 southwestern equatorial Indian Ocean and its implications for simulating the South Asian monsoon,
38 *Clim Dynam*, 40, 823-838, 2013.

39 Boucher, O., D. Randall, P. Artaxo, C. Bretherton, G. Feingold, P. Forster, V.-M. Kerminen, Y.
40 Kondo, H. Liao, U. Lohmann, P. Rasch, S.K. Satheesh, S. Sherwood, Stevens, B., and Zhang, X.
41 Y.: Clouds and Aerosols. In: *Climate Change 2013: The Physical Science Basis. Contribution of*
42 *Working Group I to the Fifth Assessment Report of the Intergovernmental Panel on Climate*
43 *Change*, Stocker, T. F., D. Qin, G.-K. Plattner, M. Tignor, S.K. Allen, J. Boschung, A. Nauels, Y.

1 Xia, V. Bex and P.M. Midgley (Ed.), Cambridge University Press, Cambridge, United Kingdom
2 and New York, NY, USA, 2013.

3 Butchart, N., Cionni, I., Eyring, V., Shepherd, T. G., Waugh, D. W., Akiyoshi, H., Austin, J., Brühl,
4 C., Chipperfield, M. P., Cordero, E., Dameris, M., Deckert, R., Dhomse, S., Frith, S. M., Garcia, R.
5 R., Gettelman, A., Giorgetta, M. A., Kinnison, D. E., Li, F., Mancini, E., McLandress, C., Pawson,
6 S., Pitari, G., Plummer, D. A., Rozanov, E., Sassi, F., Scinocca, J. F., Shibata, K., Steil, B., and
7 Tian, W.: Chemistry–Climate Model Simulations of Twenty-First Century Stratospheric Climate
8 and Circulation Changes, *J Climate*, 23, 5349-5374, 2010.

9 Chen, L., Li, T., and Yu, Y. Q.: Causes of Strengthening and Weakening of ENSO Amplitude under
10 Global Warming in Four CMIP5 Models, *J Climate*, 28, 3250-3274, 2015.

11 Chen, W. T., Woods, C. P., Li, J. L. F., Waliser, D. E., Chern, J. D., Tao, W. K., Jiang, J. H., and
12 Tompkins, A. M.: Partitioning CloudSat ice water content for comparison with upper tropospheric
13 ice in global atmospheric models, *J Geophys Res-Atmos*, 116, 2011.

14 Cherchi, A. and Navarra, A.: Influence of ENSO and of the Indian Ocean Dipole on the Indian
15 summer monsoon variability, *Clim Dynam*, 41, 81-103, 2013.

16 Cheruy, F., Dufresne, J. L., Hourdin, F., and Ducharne, A.: Role of clouds and land-atmosphere
17 coupling in midlatitude continental summer warm biases and climate change amplification in
18 CMIP5 simulations, *Geophys Res Lett*, 41, 6493-6500, 2014.

19 Choi, J., An, S. I., Kug, J. S., and Yeh, S. W.: The role of mean state on changes in El Nio's flavor,
20 *Clim Dynam*, 37, 1205-1215, 2011.

21 Collins, M., R. Knutti, J. Arblaster, J.-L. Dufresne, T. Fichet, P. Friedlingstein, X. Gao, W.J.
22 Gutowski, T. Johns, G. Krinner, M. Shongwe, C. Tebaldi, A.J. Weaver, and Wehner, M.: Long-
23 term Climate Change: Projections, Commitments and Irreversibility. In: *Climate Change 2013: The*
24 *Physical Science Basis. Contribution of Working Group I to the Fifth Assessment Report of the*
25 *Intergovernmental Panel on Climate Change*, Stocker, T. F., D. Qin, G.-K. Plattner, M. Tignor, S.K.
26 Allen, J. Boschung, A. Nauels, Y. Xia, V. Bex and P.M. Midgley (Ed.), Cambridge University
27 Press, Cambridge, United Kingdom and New York, NY, USA, 2013.

28 Collins, W. J., Bellouin, N., Doutriaux-Boucher, M., Gedney, N., Halloran, P., Hinton, T., Hughes,
29 J., Jones, C. D., Joshi, M., Liddicoat, S., Martin, G., O'Connor, F., Rae, J., Senior, C., Sitch, S.,
30 Totterdell, I., Wiltshire, A., and Woodward, S.: Development and evaluation of an Earth-System
31 model-HadGEM2, *Geosci Model Dev*, 4, 1051-1075, 2011.

32 Compo, G. P., Whitaker, J. S., Sardeshmukh, P. D., Matsui, N., Allan, R. J., Yin, X., Gleason, B.
33 E., Vose, R. S., Rutledge, G., Bessemoulin, P., Bronnimann, S., Brunet, M., Crouthamel, R. I.,
34 Grant, A. N., Groisman, P. Y., Jones, P. D., Kruk, M. C., Kruger, A. C., Marshall, G. J., Maugeri,
35 M., Mok, H. Y., Nordli, O., Ross, T. F., Trigo, R. M., Wang, X. L., Woodruff, S. D., and Worley, S.
36 J.: The Twentieth Century Reanalysis Project, *Q J Roy Meteor Soc*, 137, 1-28, 2011.

37 Connolley, W. M. and Bracegirdle, T. J.: An Antarctic assessment of IPCC AR4 coupled models,
38 *Geophys. Res. Lett.*, 34, L22505, 2007.

39 Cook, K. H. and Vizy, E. K.: Coupled model simulations of the west African monsoon system:
40 Twentieth- and Twenty-First-century simulations, *J Climate*, 19, 3681-3703, 2006.

41 Cox, P. M., Pearson, D., Booth, B. B., Friedlingstein, P., Huntingford, C., Jones, C. D., and Luke,
42 C. M.: Sensitivity of tropical carbon to climate change constrained by carbon dioxide variability,
43 *Nature*, 494, 341-344, 2013.

1 Curry, C. L.: Modeling the soil consumption of atmospheric methane at the global scale, *Global*
2 *Biogeochem. Cycles*, 21, GB4012, 2007.

3 Danabasoglu, G., Bates, S. C., Briegleb, B. P., Jayne, S. R., Jochum, M., Large, W. G., Peacock, S.,
4 and Yeager, S. G.: The CCSM4 Ocean Component, *J Climate*, 25, 1361-1389, 2012.

5 Dee, D. P., Uppala, S. M., Simmons, A. J., Berrisford, P., Poli, P., Kobayashi, S., Andrae, U.,
6 Balmaseda, M. A., Balsamo, G., Bauer, P., Bechtold, P., Beljaars, A. C. M., van de Berg, L., Bidlot,
7 J., Bormann, N., Delsol, C., Dragani, R., Fuentes, M., Geer, A. J., Haimberger, L., Healy, S. B.,
8 Hersbach, H., Holm, E. V., Isaksen, I., Kallberg, P., Kohler, M., Matricardi, M., McNally, A. P.,
9 Monge-Sanz, B. M., Morcrette, J. J., Park, B. K., Peubey, C., de Rosnay, P., Tavolato, C., Thepaut,
10 J. N., and Vitart, F.: The ERA-Interim reanalysis: configuration and performance of the data
11 assimilation system, *Q J Roy Meteor Soc*, 137, 553-597, 2011.

12 Deser, C., Alexander, M. A., Xie, S. P., and Phillips, A. S.: Sea Surface Temperature Variability:
13 Patterns and Mechanisms, *Annu Rev Mar Sci*, 2, 115-143, 2010.

14 Deser, C., Knutti, R., Solomon, S., and Phillips, A. S.: Communication of the role of natural
15 variability in future North American climate, *Nat Clim Change*, 2, 775-779, 2012.

16 Deser, C., Phillips, A. S., Alexander, M. A., and Smoliak, B. V.: Projecting North American
17 Climate over the Next 50 Years: Uncertainty due to Internal Variability*, *J Climate*, 27, 2271-2296,
18 2014.

19 Dong, S., Sprintall, J., Gille, S. T., and Talley, L.: Southern Ocean mixed-layer depth from Argo
20 float profiles, *J. Geophys. Res.*, 113, C06013, 2008.

21 Donner, L. J., Wyman, B. L., Hemler, R. S., Horowitz, L. W., Ming, Y., Zhao, M., Golaz, J. C.,
22 Ginoux, P., Lin, S. J., Schwarzkopf, M. D., Austin, J., Alaka, G., Cooke, W. F., Delworth, T. L.,
23 Freidenreich, S. M., Gordon, C. T., Griffies, S. M., Held, I. M., Hurlin, W. J., Klein, S. A., Knutson,
24 T. R., Langenhorst, A. R., Lee, H. C., Lin, Y. L., Magi, B. I., Malyshev, S. L., Milly, P. C. D., Naik,
25 V., Nath, M. J., Pincus, R., Ploshay, J. J., Ramaswamy, V., Seman, C. J., Shevliakova, E., Sirutis, J.
26 J., Stern, W. F., Stouffer, R. J., Wilson, R. J., Winton, M., Wittenberg, A. T., and Zeng, F. R.: The
27 Dynamical Core, Physical Parameterizations, and Basic Simulation Characteristics of the
28 Atmospheric Component AM3 of the GFDL Global Coupled Model CM3, *J Climate*, 24, 3484-
29 3519, 2011.

30 Dufresne, J. L., Foujols, M. A., Denvil, S., Caubel, A., Marti, O., Aumont, O., Balkanski, Y.,
31 Bekki, S., Bellenger, H., Benshila, R., Bony, S., Bopp, L., Braconnot, P., Brockmann, P., Cadule,
32 P., Cheruy, F., Codron, F., Cozic, A., Cugnet, D., Noblet, N., Duvel, J. P., Ethé, C., Fairhead, L.,
33 Fichefet, T., Flavoni, S., Friedlingstein, P., Grandpeix, J. Y., Guez, L., Guilyardi, E., Hauglustaine,
34 D., Hourdin, F., Idelkadi, A., Ghattas, J., Joussaume, S., Kageyama, M., Krinner, G., Labetoulle, S.,
35 Lahellec, A., Lefebvre, M. P., Lefevre, F., Levy, C., Li, Z. X., Lloyd, J., Lott, F., Madec, G.,
36 Mancip, M., Marchand, M., Masson, S., Meurdesoif, Y., Mignot, J., Musat, I., Parouty, S., Polcher,
37 J., Rio, C., Schulz, M., Swingedouw, D., Szopa, S., Talandier, C., Terray, P., Viovy, N., and
38 Vuichard, N.: Climate change projections using the IPSL-CM5 Earth System Model: from CMIP3
39 to CMIP5, *Clim Dynam*, doi: 10.1007/s00382-012-1636-1, 2013. 1-43, 2013.

40 Dümenil Gates, L., Hagemann, S., and Golz, C.: Observed historical discharge data from major
41 rivers for climate model validation, Max Planck Institute for Meteorology, Report 307, 2000.

42 Dunne, J. P., John, J. G., Adcroft, A. J., Griffies, S. M., Hallberg, R. W., Shevliakova, E., Stouffer,
43 R. J., Cooke, W., Dunne, K. A., Harrison, M. J., Krasting, J. P., Malyshev, S. L., Milly, P. C. D.,
44 Phillipps, P. J., Sentman, L. T., Samuels, B. L., Spelman, M. J., Winton, M., Wittenberg, A. T., and

1 Zadeh, N.: GFDL's ESM2 Global Coupled Climate-Carbon Earth System Models. Part I: Physical
2 Formulation and Baseline Simulation Characteristics, *J Climate*, 25, 6646-6665, 2012.

3 Dunne, J. P., John, J. G., Shevliakova, E., Stouffer, R. J., Krasting, J. P., Malyshev, S. L., Milly, P.
4 C. D., Sentman, L. T., Adcroft, A. J., Cooke, W., Dunne, K. A., Griffies, S. M., Hallberg, R. W.,
5 Harrison, M. J., Levy, H., Wittenberg, A. T., Phillips, P. J., and Zadeh, N.: GFDL's ESM2 Global
6 Coupled Climate-Carbon Earth System Models. Part II: Carbon System Formulation and Baseline
7 Simulation Characteristics, *J Climate*, 26, 2247-2267, 2013.

8 Edgerton, E., Lavery, T., Hodges, M., and Bowser, J.: National dry deposition network: Second
9 annual progress report, Tech. rep, 1990.

10 Emmons, L. K., Hauglustaine, D. A., Müller, J.-F., Carroll, M. A., Brasseur, G. P., Brunner, D.,
11 Staehelin, J., Thouret, V., and Marenco, A.: Data composites of airborne observations of
12 tropospheric ozone and its precursors, *J. Geophys. Res.*, 105, 20497-20538, 2000.

13 Eyring, V., Arblaster, J. M., Cionni, I., Sedlacek, J., Perliwitz, J., Young, P. J., Bekki, S.,
14 Bergmann, D., Cameron-Smith, P., Collins, W. J., Faluvegi, G., Gottschaldt, K. D., Horowitz, L.
15 W., Kinnison, D. E., Lamarque, J. F., Marsh, D. R., Saint-Martin, D., Shindell, D. T., Sudo, K.,
16 Szopa, S., and Watanabe, S.: Long-term ozone changes and associated climate impacts in CMIP5
17 simulations, *J Geophys Res-Atmos*, 118, 5029-5060, 2013.

18 Eyring, V., Bony, S., Meehl, G. A., Senior, C., Stevens, B., Stouffer, R. J., and Taylor, K. E.:
19 Overview of the Coupled Model Intercomparison Project Phase 6 (CMIP6) experimental design and
20 organisation, *Geosci. Model Dev. Discuss.*, 8, 10539-10583, 2015.

21 Eyring, V., Cionni, I., Bodeker, G. E., Charlton-Perez, A. J., Kinnison, D. E., Scinocca, J. F.,
22 Waugh, D. W., Akiyoshi, H., Bekki, S., Chipperfield, M. P., Dameris, M., Dhomse, S., Frith, S. M.,
23 Garny, H., Gettelman, A., Kubin, A., Langematz, U., Mancini, E., Marchand, M., Nakamura, T.,
24 Oman, L. D., Pawson, S., Pitari, G., Plummer, D. A., Rozanov, E., Shepherd, T. G., Shibata, K.,
25 Tian, W., Braesicke, P., Hardiman, S. C., Lamarque, J. F., Morgenstern, O., Pyle, J. A., Smale, D.,
26 and Yamashita, Y.: Multi-model assessment of stratospheric ozone return dates and ozone recovery
27 in CCMVal-2 models, *Atmos. Chem. Phys.*, 10, 9451-9472, 2010.

28 Feng, J., Liu, P., Chen, W., and Wang, X. C.: Contrasting Madden-Julian Oscillation activity during
29 various stages of EP and CP El Ninos, *Atmos Sci Lett*, 16, 32-37, 2015.

30 Ferraro, R., Waliser, D. E., Gleckler, P., Taylor, K. E., and Eyring, V.: Evolving obs4MIPs to
31 Support the Sixth Coupled Model Intercomparison Project (CMIP6), *B Am Meteorol Soc*, doi:
32 10.1175/BAMS-D-14-00216.1, 2015. 2015.

33 Flato, G., Marotzke, J., Abiodun, B., Braconnot, P., Chou, S. C., Collins, W., Cox, P., Driouech, F.,
34 Emori, S., Eyring, V., Forest, C., Gleckler, P., Guilyardi, E., Jakob, C., Kattsov, V., Reason, C., and
35 Rummukainen, M.: Evaluation of Climate Models. In: *Climate Change 2013: The Physical Science
36 Basis. Contribution of Working Group I to the Fifth Assessment Report of the Intergovernmental
37 Panel on Climate Change*, Stocker, T. F., D. Qin, G.-K. Plattner, M. Tignor, S.K. Allen, J.
38 Boschung, A. Nauels, Y. Xia, V. Bex and P.M. Midgley (Ed.), Cambridge University Press,
39 Cambridge, United Kingdom and New York, NY, USA, 2013.

40 Friedlingstein, P., Cox, P., Betts, R., Bopp, L., von Bloh, W., Brovkin, V., Cadule, P., Doney, S.,
41 Eby, M., Fung, I., Bala, G., John, J., Jones, C., Joos, F., Kato, T., Kawamiya, M., Knorr, W.,
42 Lindsay, K., Matthews, H. D., Raddatz, T., Rayner, P., Reick, C., Roeckner, E., Schnitzler, K.-G.,
43 Schnur, R., Strassmann, K., Weaver, A. J., Yoshikawa, C., and Zeng, N.: Climate–Carbon Cycle
44 Feedback Analysis: Results from the C4MIP Model Intercomparison, *J Climate*, 19, 3337-3353,
45 2006.

1 Friedlingstein, P., Meinshausen, M., Arora, V. K., Jones, C. D., Anav, A., Liddicoat, S. K., and
2 Knutti, R.: Uncertainties in CMIP5 Climate Projections due to Carbon Cycle Feedbacks, *J Climate*,
3 27, 511-526, 2014.

4 Frolicher, T. L., Sarmiento, J. L., Paynter, D. J., Dunne, J. P., Krasting, J. P., and Winton, M.:
5 Dominance of the Southern Ocean in Anthropogenic Carbon and Heat Uptake in CMIP5 Models, *J*
6 *Climate*, 28, 862-886, 2015.

7 GCOS: Implementation Plan for the Global Observing System for Climate in Support of the
8 UNFCCC, August 2010, 2010. 2010.

9 Gettelman, A., Eyring, V., Fischer, C., Shiona, H., Cionni, I., Neish, M., Morgenstern, O., Wood, S.
10 W., and Li, Z.: A community diagnostic tool for chemistry climate model validation, *Geosci. Model*
11 *Dev.*, 5, 1061-1073, 2012.

12 GEWEX-news, Vol. 21, No. 1, February 2011.

13 Ghan, S. J. and Schwartz, S. E.: Aerosol properties and processes - A path from field and laboratory
14 measurements to global climate models, *B Am Meteorol Soc*, 88, 1059-+, 2007.

15 Gibbs, H. K.: Olson's Major World Ecosystem Complexes Ranked by Carbon in Live Vegetation:
16 An Updated Database Using the GLC2000 Land Cover Product (NDP-017b), doi: DOI:
17 10.3334/CDIAC/lue.ndp017.2006, 2006. 2006.

18 Giorgetta, M. A., Jungclaus, J., Reick, C. H., Legutke, S., Bader, J., Bottinger, M., Brovkin, V.,
19 Crueger, T., Esch, M., Fieg, K., Glushak, K., Gayler, V., Haak, H., Hollweg, H. D., Ilyina, T.,
20 Kinne, S., Kornblueh, L., Matei, D., Mauritsen, T., Mikolajewicz, U., Mueller, W., Notz, D.,
21 Pithan, F., Raddatz, T., Rast, S., Redler, R., Roeckner, E., Schmidt, H., Schnur, R., Segschneider, J.,
22 Six, K. D., Stockhause, M., Timmreck, C., Wegner, J., Widmann, H., Wieners, K. H., Claussen, M.,
23 Marotzke, J., and Stevens, B.: Climate and carbon cycle changes from 1850 to 2100 in MPI-ESM
24 simulations for the Coupled Model Intercomparison Project phase 5, *Journal of Advances in*
25 *Modeling Earth Systems*, 5, 572-597, 2013.

26 Gleckler, P. J., Doutriaux, C., Durack P. J., Taylor K. E. , Zhang, Y., Williams, D. N., Mason, E.,
27 and Servonnat, J.: A More Powerful Reality Test for Climate Models, *Eos Trans. AGU*, in press,
28 2016.

29 Gleckler, P. J., Taylor, K. E., and Doutriaux, C.: Performance metrics for climate models, *J.*
30 *Geophys. Res.*, 113, D06104, 2008.

31 GLOBALVIEW-CO2: Cooperative Atmospheric Data Integration Project - Carbon Dioxide, CD-
32 ROM, NOAA ESRL, Boulder, Colorado, 2008.

33 Goswami, B. N., Krishnamurthy, V., and Annamalai, H.: A broad-scale circulation index for the
34 interannual variability of the Indian summer monsoon, *Q J Roy Meteor Soc*, 125, 611-633, 1999.

35 Guilyardi, E.: El Nino-mean state-seasonal cycle interactions in a multi-model ensemble, *Clim*
36 *Dynam*, 26, 329-348, 2006.

37 Gurney, K. R., Law, R. M., Denning, A. S., Rayner, P. J., Pak, B. C., Baker, D., Bousquet, P.,
38 Bruhwiler, L., Chen, Y. H., Ciais, P., Fung, I. Y., Heimann, M., John, J., Maki, T., Maksyutov, S.,
39 Peylin, P., Prather, M., and Taguchi, S.: Transcom 3 inversion intercomparison: Model mean results
40 for the estimation of seasonal carbon sources and sinks, *Global Biogeochem Cy*, 18, 2004.

41 Hagemann, S., Loew, A., and Andersson, A.: Combined evaluation of MPI-ESM land surface water
42 and energy fluxes, *Journal of Advances in Modeling Earth Systems*, 5, 259-286, 2013.

1 Hagemann, S., Machenhauer, B., Jones, R., Christensen, O. B., Deque, M., Jacob, D., and Vidale,
2 P. L.: Evaluation of water and energy budgets in regional climate models applied over Europe, *Clim*
3 *Dynam*, 23, 547-567, 2004.

4 Hall, A. and Qu, X.: Using the current seasonal cycle to constrain snow albedo feedback in future
5 climate change, *Geophys Res Lett*, 33, 2006.

6 Hazeleger, W., Wang, X., Severijns, C., Stefanescu, S., Bintanja, R., Sterl, A., Wyser, K., Semmler,
7 T., Yang, S., van den Hurk, B., van Noije, T., van der Linden, E., and van der Wiel, K.: EC-Earth
8 V2.2: description and validation of a new seamless earth system prediction model, *Clim Dynam*, 39,
9 2611-2629, 2012.

10 Held, I. M., Delworth, T. L., Lu, J., Findell, K. L., and Knutson, T. R.: Simulation of Sahel drought
11 in the 20th and 21st centuries, *P Natl Acad Sci USA*, 102, 17891-17896, 2005.

12 Hoell, A., Barlow, M., Wheeler, M. C., and Funk, C.: Disruptions of El Niño–Southern Oscillation
13 Teleconnections by the Madden–Julian Oscillation, *Geophys Res Lett*, 41, 998-1004, 2014.

14 Holben, B. N., Eck, T. F., Slutsker, I., Tanre, D., Buis, J. P., Setzer, A., Vermote, E., Reagan, J. A.,
15 Kaufman, Y. J., Nakajima, T., Lavenue, F., Jankowiak, I., and Smirnov, A.: AERONET - A
16 federated instrument network and data archive for aerosol characterization, *Remote Sens Environ*,
17 66, 1-16, 1998.

18 Huffman, G. J., Adler, R. F., Bolvin, D. T., Gu, G. J., Nelkin, E. J., Bowman, K. P., Hong, Y.,
19 Stocker, E. F., and Wolff, D. B.: The TRMM multisatellite precipitation analysis (TMPA): Quasi-
20 global, multiyear, combined-sensor precipitation estimates at fine scales, *J Hydrometeorol*, 8, 38-
21 55, 2007.

22 Huffman, G. J., Adler, R. F., Morrissey, M. M., Bolvin, D. T., Curtis, S., Joyce, R., McGavock, B.,
23 and Susskind, J.: Global Precipitation at One-Degree Daily Resolution from Multisatellite
24 Observations, *J Hydrometeorol*, 2, 36-50, 2001.

25 Hung, M. P., Lin, J. L., Wang, W. Q., Kim, D., Shinoda, T., and Weaver, S. J.: MJO and
26 Convectively Coupled Equatorial Waves Simulated by CMIP5 Climate Models, *J Climate*, 26,
27 6185-6214, 2013.

28 Hurrell, J. W. and Deser, C.: North Atlantic climate variability: The role of the North Atlantic
29 Oscillation, *J Marine Syst*, 78, 28-41, 2009.

30 Iguchi, T.: Correlations between interannual variations of simulated global and regional CO₂ fluxes
31 from terrestrial ecosystems and El Nino Southern Oscillation, *Tellus B*, 63, 196-204, 2011.

32 Ihaka, R. and Gentleman, R.: R: A Language for Data Analysis and Graphics, *Journal of*
33 *Computational and Graphical Statistics*, 5, 299-314, 1996.

34 Ilyina, T., Six, K. D., Segschneider, J., Maier-Reimer, E., Li, H. M., and Nunez-Riboni, I.: Global
35 ocean biogeochemistry model HAMOCC: Model architecture and performance as component of the
36 MPI-Earth system model in different CMIP5 experimental realizations, *Journal of Advances in*
37 *Modeling Earth Systems*, 5, 287-315, 2013.

38 IPCC: Climate Change 2013: The Physical Science Basis. Contribution of Working Group I to the
39 Fifth Assessment Report of the Intergovernmental Panel on Climate Change, Cambridge University
40 Press, Cambridge, United Kingdom and New York, NY, USA, 2013.

41 Jiang, J. H., Su, H., Zhai, C. X., Shen, T. J., Wu, T. W., Zhang, J., Cole, J. N. S., von Salzen, K.,
42 Donner, L. J., Seman, C., Del Genio, A., Nazarenko, L. S., Dufresne, J. L., Watanabe, M.,
43 Morcrette, C., Koshiro, T., Kawai, H., Gettelman, A., Millan, L., Read, W. G., Livesey, N. J.,

1 Kasai, Y., and Shiotani, M.: Evaluating the Diurnal Cycle of Upper-Tropospheric Ice Clouds in
2 Climate Models Using SMILES Observations, *J Atmos Sci*, 72, 1022-1044, 2015.

3 Jöckel et al., P.: Earth System Chemistry Integrated Modelling (ESCiMo) with the Modular Earth
4 Submodel System (MESSy, version 2.51), *Geosci. Model Dev.*, submitted, 2015.

5 Jöckel, P., Kerkweg, A., Pozzer, A., Sander, R., Tost, H., Riede, H., Baumgaertner, A., Gromov, S.,
6 and Kern, B.: Development cycle 2 of the Modular Earth Submodel System (MESSy2), *Geosci*
7 *Model Dev*, 3, 717-752, 2010.

8 Jones, C. D., Collins, M., Cox, P. M., and Spall, S. A.: The carbon cycle response to ENSO: A
9 coupled climate-carbon cycle model study, *J Climate*, 14, 4113-4129, 2001.

10 Jung, M., Reichstein, M., and Bondeau, A.: Towards global empirical upscaling of FLUXNET eddy
11 covariance observations: validation of a model tree ensemble approach using a biosphere model,
12 *Biogeosciences*, 6, 2001-2013, 2009.

13 Kanamitsu, M., Ebisuzaki, W., Woollen, J., Yang, S. K., Hnilo, J. J., Fiorino, M., and Potter, G. L.:
14 Ncep-Doe Amip-Ii Reanalysis (R-2), *B Am Meteorol Soc*, 83, 1631-1643, 2002.

15 Kim, D., Sperber, K., Stern, W., Waliser, D., Kang, I. S., Maloney, E., Wang, W., Weickmann, K.,
16 Benedict, J., Khairoutdinov, M., Lee, M. I., Neale, R., Suarez, M., Thayer-Calder, K., and Zhang,
17 G.: Application of MJO Simulation Diagnostics to Climate Models, *J Climate*, 22, 6413-6436,
18 2009.

19 King, M. D., Menzel, W. P., Kaufman, Y. J., Tanre, D., Gao, B. C., Platnick, S., Ackerman, S. A.,
20 Remer, L. A., Pincus, R., and Hubanks, P. A.: Cloud and aerosol properties, precipitable water, and
21 profiles of temperature and water vapor from MODIS, *Ieee T Geosci Remote*, 41, 442-458, 2003.

22 Kinne, S., Schulz, M., Litvinov, P., Stebel, K., Holzer-Popp, T., and de Leeuw, G.: ATSR Climate
23 Data Record Evaluation Report, version 1.2, ESA, Aerosol_cci,, http://www.esa-aerosol-cci.org/?q=webfm_send/836, 2015.

24 Kinne, S., Schulz, M., Textor, C., Guibert, S., Balkanski, Y., Bauer, S. E., Bernsten, T., Berglen, T.
25 F., Boucher, O., Chin, M., Collins, W., Dentener, F., Diehl, T., Easter, R., Feichter, J., Fillmore, D.,
26 Ghan, S., Ginoux, P., Gong, S., Grini, A., Hendricks, J. E., Herzog, M., Horowitz, L., Isaksen, I.,
27 Iversen, T., Kirkavag, A., Kloster, S., Koch, D., Kristjansson, J. E., Krol, M., Lauer, A., Lamarque,
28 J. F., Lesins, G., Liu, X., Lohmann, U., Montanaro, V., Myhre, G., Penner, J. E., Pitari, G., Reddy,
29 S., Seland, O., Stier, P., Takemura, T., and Tie, X.: An AeroCom initial assessment - optical
30 properties in aerosol component modules of global models, *Atmos Chem Phys*, 6, 1815-1834, 2006.

31 Kistler, R., Collins, W., Saha, S., White, G., Woollen, J., Kalnay, E., Chelliah, M., Ebisuzaki, W.,
32 Kanamitsu, M., Kousky, V., van den Dool, H., Jenne, R., and Fiorino, M.: The NCEP-NCAR 50-
33 Year Reanalysis: Monthly Means CD-ROM and Documentation, *B Am Meteorol Soc*, 82, 247-267,
34 2001.

35 Klein, S. A., Zhang, Y. Y., Zelinka, M. D., Pincus, R., Boyle, J., and Gleckler, P. J.: Are climate
36 model simulations of clouds improving? An evaluation using the ISCCP simulator, *J Geophys Res-*
37 *Atmos*, 118, 1329-1342, 2013.

38 Klotzbach, P. J.: The Madden-Julian Oscillation's Impacts on Worldwide Tropical Cyclone
39 Activity, *J Climate*, 27, 2317-2330, 2014.

40 Krishnamurthy, V. and Misra, V.: Daily atmospheric variability in the South American monsoon
41 system, *Clim Dynam*, 37, 803-819, 2011.

1 Landschützer, P., Gruber, N., Bakker, D. C. E., and Schuster, U.: An observation-based global
2 monthly gridded sea surface pCO₂ product from 1998 through 2011 and its monthly climatology.
3 Carbon Dioxide Information Analysis Center, O. R. N. L., US Department of Energy (Ed.), Oak
4 Ridge, Tennessee, 2014a.

5 Landschützer, P., Gruber, N., Bakker, D. C. E., and Schuster, U.: Recent variability of the global
6 ocean carbon sink, *Global Biogeochem Cy*, 28, 927-949, 2014b.

7 Lauer, A. and Hamilton, K.: Simulating Clouds with Global Climate Models: A Comparison of
8 CMIP5 Results with CMIP3 and Satellite Data, *J Climate*, 26, 3823-3845, 2013.

9 Lauer, A., Hendricks, J., Ackermann, I., Schell, B., Hass, H., and Metzger, S.: Simulating aerosol
10 microphysics with the ECHAM/MADE GCM - Part I: Model description and comparison with
11 observations, *Atmos Chem Phys*, 5, 3251-3276, 2005.

12 Le Quéré, C., Moriarty, R., Andrew, R. M., Peters, G. P., Ciais, P., Friedlingstein, P., Jones, S. D.,
13 Sitch, S., Tans, P., Arneeth, A., Boden, T. A., Bopp, L., Bozec, Y., Canadell, J. G., Chevallier, F.,
14 Cosca, C. E., Harris, I., Hoppema, M., Houghton, R. A., House, J. I., Jain, A., Johannessen, T.,
15 Kato, E., Keeling, R. F., Kitidis, V., Klein Goldewijk, K., Koven, C., Landa, C. S., Landschützer,
16 P., Lenton, A., Lima, I. D., Marland, G., Mathis, J. T., Metzl, N., Nojiri, Y., Olsen, A., Ono, T.,
17 Peters, W., Pfeil, B., Poulter, B., Raupach, M. R., Regnier, P., Rödenbeck, C., Saito, S., Salisbury,
18 J. E., Schuster, U., Schwinger, J., Séférian, R., Segschneider, J., Steinhoff, T., Stocker, B. D.,
19 Sutton, A. J., Takahashi, T., Tilbrook, B., van der Werf, G. R., Viovy, N., Wang, Y. P.,
20 Wanninkhof, R., Wiltshire, A., and Zeng, N.: Global carbon budget 2014, *Earth Syst. Sci. Data*
21 *Discuss.*, 7, 521-610, 2014.

22 Lee, Y. H., Lamarque, J. F., Flanner, M. G., Jiao, C., Shindell, D. T., Bernsten, T., Bisiaux, M. M.,
23 Cao, J., Collins, W. J., Curran, M., Edwards, R., Faluvegi, G., Ghan, S., Horowitz, L. W.,
24 McConnell, J. R., Ming, J., Myhre, G., Nagashima, T., Naik, V., Rumbold, S. T., Skeie, R. B.,
25 Sudo, K., Takemura, T., Thevenon, F., Xu, B., and Yoon, J. H.: Evaluation of preindustrial to
26 present-day black carbon and its albedo forcing from Atmospheric Chemistry and Climate Model
27 Intercomparison Project (ACCMIP), *Atmos Chem Phys*, 13, 2607-2634, 2013.

28 Legates, D. R. and Willmott, C. J.: Mean seasonal and spatial variability in gauge-corrected, global
29 precipitation, *International Journal of Climatology*, 10, 111-127, 1990.

30 Levine, R. C., Turner, A. G., Marathayil, D., and Martin, G. M.: The role of northern Arabian Sea
31 surface temperature biases in CMIP5 model simulations and future projections of Indian summer
32 monsoon rainfall, *Clim Dynam*, 41, 155-172, 2013.

33 Li, G. and Xie, S. P.: Tropical Biases in CMIP5 Multimodel Ensemble: The Excessive Equatorial
34 Pacific Cold Tongue and Double ITCZ Problems, *J Climate*, 27, 1765-1780, 2014.

35 Liebmann, B. and Smith, C. A.: Description of a complete (interpolated) outgoing longwave
36 radiation dataset, *B Am Meteorol Soc*, 77, 1275-1277, 1996.

37 Lin, J. L.: The double-ITCZ problem in IPCC AR4 coupled GCMs: Ocean-atmosphere feedback
38 analysis, *J Climate*, 20, 4497-4525, 2007.

39 Lin, J. L., Kiladis, G. N., Mapes, B. E., Weickmann, K. M., Sperber, K. R., Lin, W., Wheeler, M.
40 C., Schubert, S. D., Del Genio, A., Donner, L. J., Emori, S., Gueremy, J. F., Hourdin, F., Rasch, P.
41 J., Roeckner, E., and Scinocca, J. F.: Tropical intraseasonal variability in 14 IPCC AR4 climate
42 models. Part I: Convective signals, *J Climate*, 19, 2665-2690, 2006.

1 Lin, J. L., Weickman, K. M., Kiladis, G. N., Mapes, B. E., Schubert, S. D., Suarez, M. J.,
2 Bacmeister, J. T., and Lee, M. I.: Subseasonal variability associated with Asian summer monsoon
3 simulated by 14 IPCC AR4 coupled GCMs, *J Climate*, 21, 4541-4567, 2008.

4 Locarnini, R. A., Mishonov, A. V., Antonov, J. I., Boyer, T. P., Garcia, H. E., Baranova, O. K.,
5 Zweng, M. M., and Johnson, D. R.: World Ocean Atlas 2009, Volume 1: Temperature. In: NOAA
6 Atlas NESDIS 68,, Levitus, S. (Ed.), U.S. Government Printing Office, Washington, D.C., 2010.

7 Loeb, N. G., Lyman, J. M., Johnson, G. C., Allan, R. P., Doelling, D. R., Wong, T., Soden, B. J.,
8 and Stephens, G. L.: Observed changes in top-of-the-atmosphere radiation and upper-ocean heating
9 consistent within uncertainty, *Nat Geosci*, 5, 110-113, 2012.

10 Loeb, N. G., Wielicki, B. A., Doelling, D. R., Smith, G. L., Keyes, D. F., Kato, S., Manalo-Smith,
11 N., and Wong, T.: Toward Optimal Closure of the Earth's Top-of-Atmosphere Radiation Budget, *J*
12 *Climate*, 22, 748-766, 2009.

13 Lohmann, U. and Feichter, J.: Global indirect aerosol effects: a review, *Atmos. Chem. Phys.*, 5,
14 715-737, 2005.

15 Mace, G. G.: Cloud properties and radiative forcing over the maritime storm tracks of the Southern
16 Ocean and North Atlantic derived from A-Train, *J Geophys Res-Atmos*, 115, 2010.

17 Madden, R. A. and Julian, P. R.: Detection of a 40–50 Day Oscillation in the Zonal Wind in the
18 Tropical Pacific, *J. Atmos. Sci.*, 28, 702–708, 1971.

19 Madec, G.: NEMO ocean engine. Note du Pole de modélisation, Institut Pierre-Simon Laplace
20 (IPSL), France, No 27, ISSN No 1288-1619, 2008.

21 Mantua, N. J., Hare, S. R., Zhang, Y., Wallace, J. M., and Francis, R. C.: A Pacific interdecadal
22 climate oscillation with impacts on salmon production, *B Am Meteorol Soc*, 78, 1069-1079, 1997.

23 Mathon, V., Laurent, H., and Lebel, T.: Mesoscale convective system rainfall in the Sahel, *J Appl*
24 *Meteorol*, 41, 1081-1092, 2002.

25 McClain, C. R., Cleave, M. L., Feldman, G. C., Gregg, W. W., Hooker, S. B., and Kuring, N.:
26 Science quality SeaWiFS data for global biosphere research, *Sea Technol*, 39, 10-16, 1998.

27 Meier, W., Fetterer, F., Savoie, M., Mallory, S., Duerr, R., and Stroeve, J.: NOAA/NSIDC Climate
28 Data Record of Passive Microwave Sea Ice Concentration. Version 2. [sea ice concentration].
29 Center, N. S. a. I. D. (Ed.), Boulder, Colorado USA, 2013.

30 Mitchell, T. D. and Jones, P. D.: An improved method of constructing a database of monthly
31 climate observations and associated high-resolution grids, *International Journal of Climatology*, 25,
32 693-712, 2005.

33 Mueller, B., Hirschi, M., Jimenez, C., Ciais, P., Dirmeyer, P. A., Dolman, A. J., Fisher, J. B., Jung,
34 M., Ludwig, F., Maignan, F., Miralles, D. G., McCabe, M. F., Reichstein, M., Sheffield, J., Wang,
35 K., Wood, E. F., Zhang, Y., and Seneviratne, S. I.: Benchmark products for land evapotranspiration:
36 LandFlux-EVAL multi-data set synthesis, *Hydrol. Earth Syst. Sci.*, 17, 3707-3720, 2013.

37 Mueller, B. and Seneviratne, S. I.: Systematic land climate and evapotranspiration biases in CMIP5
38 simulations, *Geophys Res Lett*, 41, 128-134, 2014.

39 Myhre, G., D. Shindell, F.-M. Bréon, W. Collins, J. Fuglestad, J. Huang, D. Koch, J.-F.
40 Lamarque, D. Lee, B. Mendoza, T. Nakajima, A. Robock, G. Stephens, Takemura, T., and Zhang,
41 H.: Anthropogenic and Natural Radiative Forcing. In: *Climate Change 2013: The Physical Science*
42 *Basis. Contribution of Working Group I to the Fifth Assessment Report of the Intergovernmental*
43 *Panel on Climate Change*, Stocker, T. F., D. Qin, G.-K. Plattner, M. Tignor, S.K. Allen, J.

1 Boschung, A. Nauels, Y. Xia, V. Bex and P.M. Midgley (Ed.), Cambridge University Press,
2 Cambridge, United Kingdom and New York, NY, USA, 2013.

3 Nachtergaele, F., H. v. V., L. Verekst, and Widberg, D.: Harmonized World Soil Database v 1.2.,
4 2012.

5 Nam, C., Bony, S., Dufresne, J. L., and Chepfer, H.: The 'too few, too bright' tropical low-cloud
6 problem in CMIP5 models, *Geophys Res Lett*, 39, 2012.

7 NCL: The NCAR Command Language (Version 6.3.0) [Software], Boulder, Colorado,
8 UCAR/NCAR/CISL/TDD, <http://dx.doi.org/10.5065/D6WD3XH5>. 2016.

9 Nicholson, S. E., Some, B., and Kone, B.: An analysis of recent rainfall conditions in West Africa,
10 including the rainy seasons of the 1997 El Nino and the 1998 La Nina years, *J Climate*, 13, 2628-
11 2640, 2000.

12 Notz, D., Haumann, F. A., Haak, H., Jungclaus, J. H., and Marotzke, J.: Arctic sea-ice evolution as
13 modeled by Max Planck Institute for Meteorology's Earth system model, *Journal of Advances in*
14 *Modeling Earth Systems*, doi: 10.1002/jame.20016, 2013. n/a-n/a, 2013.

15 O'Dell, C. W., Wentz, F. J., and Bennartz, R.: Cloud liquid water path from satellite-based passive
16 microwave observations: A new climatology over the global oceans, *J Climate*, 21, 1721-1739,
17 2008.

18 Olson, J. S., Watts, J. A., and Allison, L. J.: Major world ecosystem complexes ranked by carbon in
19 live vegetation: A database (NDP-017). Carbon Dioxide Information Analysis Center. 1985.

20 Orłowsky, B. and Seneviratne, S. I.: Elusive drought: uncertainty in observed trends and short- and
21 long-term CMIP5 projections, *Hydrol Earth Syst Sc*, 17, 1765-1781, 2013.

22 Oueslati, B. and Bellon, G.: The double ITCZ bias in CMIP5 models: interaction between SST,
23 large-scale circulation and precipitation, *Clim Dynam*, 44, 585-607, 2015.

24 Pai, D. S., Bhate, J., Sreejith, O. P., and Hatwar, H. R.: Impact of MJO on the intraseasonal
25 variation of summer monsoon rainfall over India, *Clim Dynam*, 36, 41-55, 2011.

26 Peng, G., Meier, W. N., Scott, D. J., and Savoie, M. H.: A long-term and reproducible passive
27 microwave sea ice concentration data record for climate studies and monitoring, *Earth Syst. Sci.*
28 *Data*, 5, 311-318, 2013.

29 Phillips, A. S., Deser, C., and Fasullo, J.: Evaluating Modes of Variability in Climate Models, *Eos*
30 *Trans. AGU*, 95(49), 453-455, 2014.

31 Pierrehumbert, R. T.: Thermostats, Radiator Fins, and the Local Runaway Greenhouse, *J Atmos*
32 *Sci*, 52, 1784-1806, 1995.

33 Pincus, R., Batstone, C. P., Hofmann, R. J. P., Taylor, K. E., and Glecker, P. J.: Evaluating the
34 present-day simulation of clouds, precipitation, and radiation in climate models, *J. Geophys. Res.*,
35 113, D14209, 2008.

36 Pincus, R., Platnick, S., Ackerman, S. A., Hemler, R. S., and Hofmann, R. J. P.: Reconciling
37 Simulated and Observed Views of Clouds: MODIS, ISCCP, and the Limits of Instrument
38 Simulators, *J Climate*, 25, 4699-4720, 2012.

39 Pozzer, A., de Meij, A., Pringle, K. J., Tost, H., Doering, U. M., van Aardenne, J., and Lelieveld, J.:
40 Distributions and regional budgets of aerosols and their precursors simulated with the EMAC
41 chemistry-climate model, *Atmos Chem Phys*, 12, 961-987, 2012.

1 Pringle, K. J., Tost, H., Message, S., Steil, B., Giannadaki, D., Nenes, A., Fountoukis, C., Stier, P.,
2 Vignati, E., and Leueved, J.: Description and evaluation of GMXe: a new aerosol submodel for
3 global simulations (v1), *Geosci Model Dev*, 3, 391-412, 2010.

4 Rayner, N. A., Parker, D. E., Horton, E. B., Folland, C. K., Alexander, L. V., Rowell, D. P., Kent,
5 E. C., and Kaplan, A.: Global analyses of sea surface temperature, sea ice, and night marine air
6 temperature since the late nineteenth century, *J. Geophys. Res.*, 108, 4407, 2003.

7 Redelsperger, J. L., Thorncroft, C. D., Diedhiou, A., Lebel, T., Parker, D. J., and Polcher, J.:
8 African monsoon multidisciplinary analysis - An international research project and field campaign,
9 *B Am Meteorol Soc*, 87, 1739-+, 2006.

10 Reichler, T. and Kim, J.: How Well Do Coupled Models Simulate Today's Climate?, *B Am*
11 *Meteorol Soc*, 89, 303-311, 2008.

12 Richter, I., Behera, S. K., Doi, T., Taguchi, B., Masumoto, Y., and Xie, S. P.: What controls
13 equatorial Atlantic winds in boreal spring?, *Clim Dynam*, 43, 3091-3104, 2014.

14 Rienecker, M. M., Suarez, M. J., Gelaro, R., Todling, R., Bacmeister, J., Liu, E., Bosilovich, M. G.,
15 Schubert, S. D., Takacs, L., Kim, G. K., Bloom, S., Chen, J. Y., Collins, D., Conaty, A., Da Silva,
16 A., Gu, W., Joiner, J., Koster, R. D., Lucchesi, R., Molod, A., Owens, T., Pawson, S., Pegion, P.,
17 Redder, C. R., Reichle, R., Robertson, F. R., Ruddick, A. G., Sienkiewicz, M., and Woollen, J.:
18 MERRA: NASA's Modern-Era Retrospective Analysis for Research and Applications, *J Climate*,
19 24, 3624-3648, 2011.

20 Righi, M., Eyring, V., Gottschaldt, K. D., Klinger, C., Frank, F., Jöckel, P., and Cionni, I.:
21 Quantitative evaluation of ozone and selected climate parameters in a set of EMAC simulations,
22 *Geosci. Model Dev.*, 8, 733-768, 2015.

23 Righi, M., Hendricks, J., and Sausen, R.: The global impact of the transport sectors on atmospheric
24 aerosol: simulations for year 2000 emissions, *Atmos Chem Phys*, 13, 9939-9970, 2013.

25 Rio, C., Hourdin, F., Grandpeix, J. Y., and Lafore, J. P.: Shifting the diurnal cycle of parameterized
26 deep convection over land, *Geophys Res Lett*, 36, 2009.

27 Rodenbeck, C., Bakker, D. C. E., Metzl, N., Olsen, A., Sabine, C., Cassar, N., Reum, F., Keeling,
28 R. F., and Heimann, M.: Interannual sea-air CO₂ flux variability from an observation-driven ocean
29 mixed-layer scheme, *Biogeosciences*, 11, 4599-4613, 2014.

30 Roehrig, R., Bouniol, D., Guichard, F., Hourdin, F., and Redelsperger, J. L.: The Present and Future
31 of the West African Monsoon: A Process-Oriented Assessment of CMIP5 Simulations along the
32 AMMA Transect, *J Climate*, 26, 6471-6505, 2013.

33 Rossow, W. B. and Schiffer, R. A.: Advances in Understanding Clouds from ISCCP, *B Am*
34 *Meteorol Soc*, 80, 2261-2287, 1999.

35 Rossow, W. B. and Schiffer, R. A.: ISCCP Cloud Data Products, *B Am Meteorol Soc*, 72, 2-20,
36 1991.

37 Sabeerali, C. T., Dandi, A., Dhakate, A., Salunke, K., Mahapatra, S., and Rao, S. A.: Simulation of
38 boreal summer intraseasonal oscillations in the latest CMIP5 coupled GCMs, *J Geophys Res-*
39 *Atmos*, 118, 4401-4420, 2013.

40 Sanderson, B. M., Knutti, R., and Caldwell, P.: Addressing Interdependency in a Multimodel
41 Ensemble by Interpolation of Model Properties, *J Climate*, 28, 5150-5170, 2015a.

42 Sanderson, B. M., Knutti, R., and Caldwell, P.: A Representative Democracy to Reduce
43 Interdependency in a Multimodel Ensemble, *J Climate*, 28, 5171-5194, 2015b.

1 Schulz, M., Textor, C., Kinne, S., Balkanski, Y., Bauer, S., Bernsten, T., Berglen, T., Boucher, O.,
2 Dentener, F., Guibert, S., Isaksen, I. S. A., Iversen, T., Koch, D., Kirkevåg, A., Liu, X., Montanaro,
3 V., Myhre, G., Penner, J. E., Pitari, G., Reddy, S., Seland, O., Stier, P., and Takemura, T.: Radiative
4 forcing by aerosols as derived from the AeroCom present-day and pre-industrial simulations, *Atmos*
5 *Chem Phys*, 6, 5225-5246, 2006.

6 Shi, Y., Zhang, J., Reid, J. S., Holben, B., Hyer, E. J., and Curtis, C.: An analysis of the collection 5
7 MODIS over-ocean aerosol optical depth product for its implication in aerosol assimilation, *Atmos*
8 *Chem Phys*, 11, 557-565, 2011.

9 Smith, G. L., Mlynchak, P. E., Rutan, D. A., and Wong, T.: Comparison of the Diurnal Cycle of
10 Outgoing Longwave Radiation from a Climate Model with Results from ERBE, *J Appl Meteorol*
11 *Clim*, 47, 3188-3201, 2008.

12 SPARC-CCMVal: SPARC Report on the Evaluation of Chemistry-Climate Models, V. Eyring, T.
13 G. Shepherd, D. W. Waugh (Eds.). SPARC Report No. 5, WCRP-132, WMO/TD-No. 1526., 2010.

14 Sperber, K., Annamalai, H., Kang, I. S., Kitoh, A., Moise, A., Turner, A., Wang, B., and Zhou, T.:
15 The Asian summer monsoon: an intercomparison of CMIP5 vs. CMIP3 simulations of the late 20th
16 century, *Clim Dynam*, 41, 2711-2744, 2013.

17 Stephens, G. L. and Greenwald, T. J.: The Earth's Radiation Budget and Its Relation to Atmospheric
18 Hydrology .1. Observations of the Clear Sky Greenhouse-Effect, *J Geophys Res-Atmos*, 96, 15311-
19 15324, 1991.

20 Stephens, G. L., Vane, D. G., Boain, R. J., Mace, G. G., Sassen, K., Wang, Z. E., Illingworth, A. J.,
21 O'Connor, E. J., Rossow, W. B., Durden, S. L., Miller, S. D., Austin, R. T., Benedetti, A., Mitrescu,
22 C., and Team, C. S.: The cloudsat mission and the a-train - A new dimension of space-based
23 observations of clouds and precipitation, *B Am Meteorol Soc*, 83, 1771-1790, 2002.

24 Sterl, A., Bintanja, R., Brodeau, L., Gleeson, E., Koenigk, T., Schmith, T., Semmler, T., Severijns,
25 C., Wyser, K., and Yang, S. T.: A look at the ocean in the EC-Earth climate model, *Clim Dynam*,
26 39, 2631-2657, 2012.

27 Stevens, B. and Schwartz, S. E.: Observing and Modeling Earth's Energy Flows, *Surv Geophys*, 33,
28 779-816, 2012.

29 Stowasser, M., Annamalai, H., and Hafner, J.: Response of the South Asian Summer Monsoon to
30 Global Warming: Mean and Synoptic Systems, *J Climate*, 22, 1014-1036, 2009.

31 Stroeve, J., Holland, M. M., Meier, W., Scambos, T., and Serreze, M.: Arctic sea ice decline: Faster
32 than forecast, *Geophys. Res. Lett.*, 34, L09501, 2007.

33 Stroeve, J. C., Kattsov, V., Barrett, A., Serreze, M., Pavlova, T., Holland, M., and Meier, W. N.:
34 Trends in Arctic sea ice extent from CMIP5, CMIP3 and observations, *Geophys Res Lett*, 39, 2012.

35 Takahashi, T., Sutherland, S. C., Chipman, D. W., Goddard, J. G., Ho, C., Newberger, T., Sweeney,
36 C., and Munro, D. R.: Climatological distributions of pH, pCO₂, total CO₂, alkalinity, and CaCO₃
37 saturation in the global surface ocean, and temporal changes at selected locations, *Mar. Chem.*, 164,
38 95-125, 2014.

39 Taylor, K. E.: Summarizing multiple aspects of model performance in a single diagram., *J Geophys*
40 *Res-Atmos*, 106, 7183-7192, 2001.

41 Taylor, K. E., Stouffer, R. J., and Meehl, G. A.: An Overview of Cmp5 and the Experiment
42 Design, *B Am Meteorol Soc*, 93, 485-498, 2012.

1 Teixeira, J., Waliser, D., Ferraro, R., Gleckler, P., Lee, T., and Potter, G.: Satellite Observations for
2 CMIP5: The Genesis of Obs4MIPs, *B Am Meteorol Soc*, 95, 1329-1334, 2014.

3 Thompson, D. W. J. and Wallace, J. M.: Annular modes in the extratropical circulation. Part I:
4 Month-to-month variability, *J Climate*, 13, 1000-1016, 2000.

5 Tilmes, S., Lamarque, J. F., Emmons, L. K., Conley, A., Schultz, M. G., Saunois, M., Thouret, V.,
6 Thompson, A. M., Oltmans, S. J., Johnson, B., and Tarasick, D.: Ozone-sonde climatology between
7 1995 and 2009: description, evaluation and applications, *Atmos. Chem. Phys. Discuss.*, 11, 28747-
8 28796, 2011.

9 Totsuka, T., Sase, H., and Shimizu, H.: Major activities of acid deposition monitoring network in
10 East Asia (EANET) and related studies. In: *Plant Responses to Air Pollution and Global Change*,
11 Omasa, K., Nouchi, I., and De Kok, L. (Eds.), Springer Japan, 2005.

12 Trenberth, K. E. and Fasullo, J. T.: An observational estimate of inferred ocean energy divergence,
13 *J Phys Oceanogr*, 38, 984-999, 2008.

14 Trenberth, K. E. and Fasullo, J. T.: Simulation of Present-Day and Twenty-First-Century Energy
15 Budgets of the Southern Oceans, *J Climate*, 23, 440-454, 2010.

16 Trenberth, K. E. and Shea, D. J.: Atlantic hurricanes and natural variability in 2005, *Geophys Res*
17 *Lett*, 33, 2006.

18 Turner, A. G., Inness, P. M., and Slingo, J. M.: The role of the basic state in the ENSO-monsoon
19 relationship and implications for predictability, *Q J Roy Meteor Soc*, 131, 781-804, 2005.

20 Voulgarakis, A., Naik, V., Lamarque, J. F., Shindell, D. T., Young, P. J., Prather, M. J., Wild, O.,
21 Field, R. D., Bergmann, D., Cameron-Smith, P., Cionni, I., Collins, W. J., Dalsoren, S. B., Doherty,
22 R. M., Eyring, V., Faluvegi, G., Folberth, G. A., Horowitz, L. W., Josse, B., MacKenzie, I. A.,
23 Nagashima, T., Plummer, D. A., Righi, M., Rumbold, S. T., Stevenson, D. S., Strode, S. A., Sudo,
24 K., Szopa, S., and Zeng, G.: Analysis of present day and future OH and methane lifetime in the
25 ACCMIP simulations, *Atmos Chem Phys*, 13, 2563-2587, 2013.

26 Waliser, D., Sperber, K., Hendon, H., Kim, D., Wheeler, M., Weickmann, K., Zhang, C., Donner,
27 L., Gottschalck, J., Higgins, W., Kang, I. S., Legler, D., Moncrieff, M., Vitart, F., Wang, B., Wang,
28 W., Woolnough, S., Maloney, E., Schubert, S., Stern, W., and Oscillation, C. M.-J.: MJO
29 Simulation Diagnostics, *J Climate*, 22, 3006-3030, 2009.

30 Wang, B. and Fan, Z.: Choice of south Asian summer monsoon indices, *B Am Meteorol Soc*, 80,
31 629-638, 1999.

32 Wang, B., Liu, J., Kim, H. J., Webster, P. J., and Yim, S. Y.: Recent change of the global monsoon
33 precipitation (1979-2008), *Clim Dynam*, 39, 1123-1135, 2012.

34 Waugh, D. W. and Eyring, V.: Quantitative performance metrics for stratospheric-resolving
35 chemistry-climate models, *Atmos. Chem. Phys.*, 8, 5699-5713, 2008.

36 Webster, P. J. and Yang, S.: Monsoon and Enso - Selectively Interactive Systems, *Q J Roy Meteor*
37 *Soc*, 118, 877-926, 1992.

38 Weedon, G. P., Balsamo, G., Bellouin, N., Gomes, S., Best, M. J., and Viterbo, P.: The WFDEI
39 meteorological forcing data set: WATCH Forcing Data methodology applied to ERA-Interim
40 reanalysis data, *Water Resour Res*, 50, 7505-7514, 2014.

41 Wenzel, S., Cox, P. M., Eyring, V., and Friedlingstein, P.: Emergent constraints on climate-carbon
42 cycle feedbacks in the CMIP5 Earth system models, *Journal of Geophysical Research:*
43 *Biogeosciences*, 119, 2013JG002591, 2014.

1 Wielicki, B. A., Barkstrom, B. R., Harrison, E. F., Lee, R. B., Louis Smith, G., and Cooper, J. E.:
2 Clouds and the Earth's Radiant Energy System (CERES): An Earth Observing System Experiment,
3 B Am Meteorol Soc, 77, 853-868, 1996.

4 Williams, K. and Webb, M.: A quantitative performance assessment of cloud regimes in climate
5 models, Clim Dynam, 33, 141-157, 2009.

6 Xie, P. and Arkin, P. A.: Global Precipitation: A 17-Year Monthly Analysis Based on Gauge
7 Observations, Satellite Estimates, and Numerical Model Outputs, B Am Meteorol Soc, 78, 2539-
8 2558, 1997.

9 Young, P. J., Archibald, A. T., Bowman, K. W., Lamarque, J. F., Naik, V., Stevenson, D. S.,
10 Tilmes, S., Voulgarakis, A., Wild, O., Bergmann, D., Cameron-Smith, P., Cionni, I., Collins, W. J.,
11 Dalsøren, S. B., Doherty, R. M., Eyring, V., Faluvegi, G., Horowitz, L. W., Josse, B., Lee, Y. H.,
12 MacKenzie, I. A., Nagashima, T., Plummer, D. A., Righi, M., Rumbold, S. T., Skeie, R. B.,
13 Shindell, D. T., Strode, S. A., Sudo, K., Szopa, S., and Zeng, G.: Pre-industrial to end 21st century
14 projections of tropospheric ozone from the Atmospheric Chemistry and Climate Model
15 Intercomparison Project (ACCMIP), Atmos Chem Phys, 13, 2063-2090, 2013.

16 Yu, L., Xiangze Jin, and Weller, R. A.: Multidecade Global Flux Datasets from the Objectively
17 Analyzed Air-sea Fluxes (OAFlux) Project: Latent and Sensible Heat Fluxes, Ocean Evaporation,
18 and Related Surface Meteorological Variables. (OA-2008-01), W. H. O. I. O. P. T. R. (Ed.), 2008.

19 Zhang, Y. C., Rossow, W. B., Lacis, A. A., Oinas, V., and Mishchenko, M. I.: Calculation of
20 radiative fluxes from the surface to top of atmosphere based on ISCCP and other global data sets:
21 Refinements of the radiative transfer model and the input data, J Geophys Res-Atmos, 109, 2004.

22 Zhu, Z. C., Bi, J., Pan, Y. Z., Ganguly, S., Anav, A., Xu, L., Samanta, A., Piao, S. L., Nemani, R.
23 R., and Myneni, R. B.: Global Data Sets of Vegetation Leaf Area Index (LAI)3g and Fraction of
24 Photosynthetically Active Radiation (FPAR)3g Derived from Global Inventory Modeling and
25 Mapping Studies (GIMMS) Normalized Difference Vegetation Index (NDVI3g) for the Period
26 1981 to 2011, Remote Sens-Basel, 5, 927-948, 2013.

27 Ziemke, J. R., Chandra, S., Labow, G. J., Bhartia, P. K., Froidevaux, L., and Witte, J. C.: A global
28 climatology of tropospheric and stratospheric ozone derived from Aura OMI and MLS
29 measurements, Atmos Chem Phys, 11, 9237-9251, 2011.

30

31

1 Table 1. Overview of standard namelists implemented in ESMValTool (v1.0) along with the
2 quantity and ESMValTool variable name for which the namelist is tested, the corresponding
3 observations or reanalyses, the section and example figure in this paper, and references for the
4 namelist. When the namelist is named with a specific paper (naming convention:
5 *namelist_SurnameYearJournalabbreviation.xml*), it can be used to reproduce in general all or in
6 some cases only a subset of the figures published in that paper. Otherwise the namelists group a set
7 of diagnostics and performance metrics for a specific scientific topic (e.g.,
8 *namelist_aerosol_CMIP5.xml*). Observations and reanalyses are listed together with their Tier, type
9 (e.g., reanalysis, satellite or in situ observations), the time period used, and a reference. Tier 1
10 includes observations from obs4MIPs or reanalyses from ana4MIPs. Tier 2 and tier 3 indicate
11 freely-available and restricted data sets, respectively. For these observations, reformatting routines
12 are provided to bring the original data in the CF/CMOR standard format so that they can directly be
13 used in the ESMValTool.

<i>xml namelist</i>	Tested Quantity (CMOR units)	ESMValTool Variable Name	Tested Observations /Reanalyses (Tier, type, time period, reference)	Section / Example Figure(s)	References for namelist
Section 4.1: Detection of systematic biases in the physical climate: atmosphere					
<i>namelist_performance_CMI_P5</i>	Temperature (K)	ta	ERA-Interim (Tier 3, reanalysis, 1979-2014 (Dee et al., 2011))	Section 4.1.1. / Fig. 2 and Fig. 3	Gleckler et al. (2008); Taylor (2001); Fig. 9.7 of Flato et al. (2013) Righi et al. (2015)
<i>namelist_righi_15gmd_ECVs</i>	Eastward wind (m s^{-1})	ua			
	Northward wind (m s^{-1})	va			
	Near-surface air temperature (K)	tas	NCEP (Tier 2, reanalysis, 1948-2012 (Kistler et al., 2001))		
	Geopotential height (m)	zg			
	Specific Humidity (1)	hus	AIRS (Tier 1, satellite, 2003-2010 (Aumann et al., 2003))		
	Precipitation ($\text{kg m}^{-2} \text{ s}^{-1}$)	pr	GPCP-SG (Tier 1, satellite & rain gauge, 1979-near-present (Adler et al., 2003))		
	TOA outgoing shortwave radiation (W m^{-2})	rsut	CERES-EBAF (Tier 1, satellite, 2001-2011 (Wielicki et al., 1996))		
	TOA outgoing longwave radiation (W m^{-2})	rlut			
	TOA outgoing clear sky longwave radiation (W m^{-2})	rlutes			

	Shortwave cloud radiative effect (W m^{-2})	SW_CRE			
	Longwave cloud radiative effect (W m^{-2})	LW_CRE			
	Aerosol optical depth at 550 nm (1)	od550aer	MODIS (Tier 1, satellite, 2001-2012 (King et al., 2003)) ESACCI-AEROSOL (Tier 2, satellite, 1996-2012 (Kinne et al., 2015))		
	Total cloud amount (%)	clt	MODIS (Tier 1, satellite, 2001-2012 (King et al., 2003))		
<i>namelist_flato13ipcc</i>	Near-surface air temperature (K)	tas	ERA-Interim (Tier 3, reanalysis, 1979-2014 (Dee et al., 2011))	Section 4.1.2 / Fig. 4	Fig. 9.2 and Fig. 9.4 of Flato et al. (2013)
	Precipitation ($\text{kg m}^{-2} \text{s}^{-1}$)	pr	GPCP-1DD (Tier 1, satellite, 1997-2010 (Huffman et al., 2001))		
<i>namelist_SAMonsoon</i> <i>namelist_SAMonsoon_AMIP</i> <i>namelist_SAMonsoon_daily</i>	Eastward wind (m s^{-1})	ua	ERA-Interim (Tier 3, reanalysis, 1979-2014 (Dee et al., 2011)) MERRA (Tier 1, reanalysis, 1979-2011 (Rienecker et al., 2011))	Section 4.1.3.1 / Fig. 5 and Fig. 6	Goswami et al. (1999) Sperber et al. (2013) Wang and Fan (1999) Wang et al. (2012) Webster and Yang (1992) Lin et al. (2008); Fig. 9.32 of Flato et al. (2013)
	Northward wind (m s^{-1})	va			
	Precipitation ($\text{kg m}^{-2} \text{s}^{-1}$)	pr	TRMM-3B42-v7 (Tier 1, satellite, 1998-near-present (Huffman et al., 2007)) GPCP-1DD 1DD (Tier 1, satellite, 1997-2010 (Huffman et al., 2001)) CMAP (Tier 2, satellite & rain gauge, 1979-near-present (Xie and Arkin, 1997)) MERRA (Tier 1, reanalysis, 1979-2011 (Rienecker		

			et al., 2011))		
			ERA-Interim (Tier 3, reanalysis, 1979-2014 (Dee et al., 2011))		
	Skin temperature (K)	ts	HadISST (Tier 2, reanalysis, 1870-2014 (Rayner et al., 2003))		
<i>namelist_WA Monsoon</i>	Eastward wind (m s^{-1})	ua	ERA-Interim (Tier 3, reanalysis, 1979-2014 (Dee et al., 2011))	Section 4.1.3.2 / Fig. 7	Roehrig et al. (2013); Cook and Vizy (2006)
<i>namelist_WA Monsoon_daily</i>	Northward wind (m s^{-1})	va			
	Temperature (K)	ta			
	Near-surface air temperature (K)	tas			
	Precipitation ($\text{kg m}^{-2} \text{s}^{-1}$)	pr	GPCP-1DD (Tier 1, satellite, 1997-2010 (Huffman et al., 2001)) TRMM (Tier 1, satellite, 1998-present (Huffman et al., 2007))		
	TOA outgoing shortwave radiation (W m^{-2})	rsut	CERES-EBAF (Tier 1, satellite, 2001-2011 (Wielicki et al., 1996))		
	TOA outgoing longwave radiation (W m^{-2})	rlut			
	TOA outgoing clear sky shortwave radiation (W m^{-2})	rsutcs			
	TOA outgoing clear sky longwave radiation (W m^{-2})	rlutcs			
	Shortwave cloud radiative effect (W m^{-2})	SW_CRE			
	Longwave cloud radiative effect (W m^{-2})	LW_CRE			
	Shortwave downwelling radiation at surface (W m^{-2})	rsds			
	Longwave downwelling radiation at surface (W m^{-2})	rlds			
	TOA outgoing longwave radiation (W m^{-2})	rlut	NOAA polar-orbiting satellites (Tier 2, satellite, 1974-2013 (Liebmann and Smith, 1996))		
<i>namelist_CV</i>	Precipitation ($\text{kg m}^{-2} \text{s}^{-1}$)	pr	GPCP-SG (Tier 1,	Section	Phillips et

<i>DP</i>			satellite & rain gauge, 1979-near-present (Adler et al., 2003)) TRMM (Tier 1, satellite, 1998-near-present (Huffman et al., 2007))	4.1.4.1 / Fig. 8 and Fig. 9	al. (2014)
	Air pressure at sea level (Pa)	psl	NOAA-CIRES Twentieth Century Reanalysis Project (Tier 1, reanalysis, 1900-2012 (Compo et al., 2011))		
	Near-surface air temperature (K)	tas	NCEP (Tier 2, reanalysis, 1948-2012 (Kistler et al., 2001))		
	Skin temperature (K)	ts	HadISST (Tier 2, satellite-based, 1870-2014 (Rayner et al., 2003))		
	Snow depth (m)	snd	without obs		
	Ocean meridional overturning mass streamfunction (kg s^{-1})	msftmyz	without obs		
<i>namelist_mjo_daily</i> <i>namelist_mjo_mean_state</i>	Eastward wind (m s^{-1})	ua	ERA-Interim (Tier 3, reanalysis, 1979-2014 (Dee et al., 2011))	Section 4.1.4.2 / Fig. 10	Waliser et al. (2009); Kim et al. (2009)
	Northward wind (m s^{-1})	va	NCEP (Tier 2, reanalysis, 1979-2013 (Kistler et al., 2001))		
	Precipitation ($\text{kg m}^{-2} \text{s}^{-1}$)	pr	GPCP-1DD (Tier 1, satellite, 1997-2010 (Huffman et al., 2001))		
	TOA longwave radiation (W m^{-2})	rlut	NOAA polar-orbiting satellites (Tier 2, satellite, 1974-2013 (Liebmann and Smith, 1996))		
<i>namelist_DiurnalCycle</i>	Precipitation ($\text{kg m}^{-2} \text{s}^{-1}$)	pr	TRMM (Tier 1, satellite, 1998-near-present (Huffman et al., 2007))	Section 4.1.5 / Fig. 11	Rio et al. (2009)
	Convective Precipitation ($\text{kg m}^{-2} \text{s}^{-1}$)	prc			

	TOA outgoing longwave radiation (W m^{-2})	rlut	CERES-SYN1deg (Tier 1, satellite, 2001-2011		
	TOA outgoing shortwave radiation (W m^{-2})	rsut	(Wielicki et al., 1996))		
	TOA outgoing clear sky longwave radiation (W m^{-2})	rlutcs			
	TOA outgoing clear sky shortwave radiation (W m^{-2})	rsutcs			
	Surface downwelling shortwave radiation (W m^{-2})	rsds			
	Surface downwelling clear sky sky shortwave radiation (W m^{-2})	rsdscs			
	Surface upwelling shortwave radiation (W m^{-2})	rsus			
	Surface upwelling clear sky shortwave radiation (W m^{-2})	rsuscs			
	Surface upwelling longwave radiation (W m^{-2})	rlus			
	Surface upwelling clear sky longwave radiation (W m^{-2})	rluscs			
	Surface downwelling shortwave radiation (W m^{-2})	rlds			
	Surface downwelling clear sky longwave radiation (W m^{-2})	rldscs			
<i>namelist_laue</i> <i>r13jclim</i>	Atmosphere cloud condensed water content (kg m^{-2})	clwvi	UWisc: SSM/I, TMI, AMSR-E (Tier 3, satellite, 1988-2007 (O'Dell et al., 2008))	Section 4.1.6.1 / Fig. 12	Lauer and Hamilton (2013); Fig. 9.5 of Flato et al. (2013)
	Atmosphere cloud ice content (kg m^{-2})	clivi	MODIS-CFMIP (Tier 2, satellite, 2003-2014 (King et al., 2003; Pincus et al., 2012))		
	Total cloud amount (%)	clt	MODIS (Tier 1, satellite, 2001-2012 (King et al., 2003))		
	TOA outgoing longwave radiation (W m^{-2})	rlut	CERES-EBAF (Tier 1, satellite, 2001-2011		
	TOA outgoing longwave radiation (clear sky) (W m^{-2})	rlutcs	(Wielicki et al., 1996))		
	TOA outgoing shortwave radiation (W m^{-2})	rsut	SRB (Tier 2, satellite, 1984-		

	TOA outgoing shortwave radiation (clear sky) (W m ⁻²)	rsutcs	2007 (GEWEX-news, February 2011))		
	Precipitation (kg m ⁻² s ⁻¹)	pr	GPCP-SG (Tier 1, satellite & rain gauge, 1979-near-present (Adler et al., 2003))		
namelist_williams09climdyn_CREM	ISCCP mean cloud albedo (1)	albiscpp	ISCCP (Tier 1, satellite, 1985-1990 (Rossow and Schiffer, 1991)) ISCCP-FD (Tier 2, satellite, 1985-1990 (Zhang et al., 2004))	Section 4.1.6.2 / Fig. 13	Williams and Webb (2009)
	ISCCP mean cloud top pressure (Pa)	pctisccp			
	ISCCP total cloud fraction (%)	cltisccp			
	TOA outgoing shortwave radiation (W m ⁻²)	rsut			
	TOA outgoing longwave radiation (W m ⁻²)	rlut			
	TOA outoing clear sky shortwave radiation (W m ⁻²)	rsutcs			
	TOA outoing clear sky longwave radiation (W m ⁻²)	rlutcs			
	Surface snow area fraction (%)	snc			
	Surface snow amount (kg m ⁻²)	snw			
	Sea ice area fraction (%)	sic			
Section 4.2: Detection of systematic biases in the physical climate: ocean					
namelist_SouthernOcean	Ocean Mixed Layer Thickness Defined by Sigma T (m)	mloitst	ARGO (Tier 2, Buoy, Monthly mean climatology 2001-2006 (Dong et al., 2008))	Section 4.2.2.1 / Fig. 14	CDFTOOLS
	Sea surface temperature (K)	tos	ERA-Interim (Tier 3, reanalysis, 1979-2014 (Dee et al., 2011))		
	Downward heat flux at sea water surface (W m ⁻²)	hfds (hfsl + hfss + rsns + rlms)			
	Surface Downward Eastward Wind Stress (Pa)	tauu			
	Surface Downward Nordward Wind Stress (Pa)	tauv			
	Water Flux from precipitation and evaporation (kg m-2 s ⁻¹)	wfpe (pr + evspsbl)			
	Sea water salinity (psu)	so	WOA09 (Tier 2, in-situ, climatology, (Antonov et al., 2010; Locarnini et al., 2010))		
	Sea surface salinity (psu)	sos			
	Sea Water Temperature (K)	to			
	Sea Water X Velocity (m s ⁻¹)	uo	without obs		

	Sea Water Y Velocity (m s^{-1})	vo			
<i>namelist_SouthernHemisphere</i>	Total Cloud Fraction (%)	clt	CloudSat (Tier 1, satellite, 2000-2005 (Stephens et al., 2002))	Section 4.2.2.2 / Fig. 15	Frolicher et al. (2015)
	Atmosphere cloud ice content (kg m^{-2})	clivi			
	Atmosphere cloud condensed water content (kg m^{-2})	clwvi			
	Surface upward latent heat flux (W m^{-2})	hfls	WHOI-OAflux (Tier 2, satellite-based, 2000-2005 (Yu et al., 2008))		
	Surface upward sensible heat flux (W m^{-2})	hfss			
	TOA outgoing longwave radiation (W m^{-2})	rlut	CERES-EBAF (Tier 1, satellite, 2001-2011 (Wielicki et al., 1996))		
	TOA outgoing clear sky longwave radiation (W m^{-2})	rlutcs			
	TOA outgoing shortwave radiation (W m^{-2})	rsut	SRB (Tier 2, satellite, 1984-2007 (GEWEX-news, February 2011))		
	TOA outgoing clear sky shortwave radiation (W m^{-2})	rsutcs			
	Surface downwelling shortwave radiation (W m^{-2})	rlds			
<i>namelist_TropicalVariability</i>	Surface downwelling clear sky longwave radiation (W m^{-2})	rldscs		Section 4.2.3 / Fig. 16	Choi et al. (2011); Li and Xie (2014)
	Surface downwelling shortwave radiation (W m^{-2})	rsds			
	Surface downwelling clear sky shortwave radiation (W m^{-2})	rsdscs			
	Precipitation ($\text{kg m}^{-2} \text{ s}^{-1}$)	pr	TRMM (Tier 1, satellite, 1998-present (Huffman et al., 2007))		
<i>namelist_Sealce</i>	Sea surface temperature (K)	ts	HadISST (Tier 2, satellite-based, 1870-2014 (Rayner et al., 2003))	Section 4.2.4 / Fig. 17	Stroeve et al. (2007) Stroeve et al. (2012); Fig. 9.24 of Flato et al. (2013)
	Eastward wind (m s^{-1})	ua	ERA-Interim (Tier 3, reanalysis, 1979-2014 (Dee et al., 2011))		
	Northward wind (m s^{-1})	va			
	Sea ice area fraction (%)	sic	HadISST (Tier 2, satellite-based, 1870-2014 (Rayner et al., 2003)) NSIDC (Tier 2,		

			satellite, 1978-2010 (Meier et al., 2013; Peng et al., 2013))		
Section 4.3: Detection of systematic biases in the physical climate: land					
<i>namelist_Evapotranspiration</i>	Surface upward latent heat flux (W m^{-2})	hfls	LandFlux-EVAL (Tier 3, ground, 1989-2004 (Mueller et al., 2013)) GPCC (Tier 2, Rain gauge analysis, 1901-2010 (Becker et al., 2013))	Section 4.3.1 / Fig. 18	Mueller and Seneviratne (2014); Orłowsky and Seneviratne (2013)
<i>namelist_SPI</i>	Precipitation ($\text{kg m}^{-2} \text{s}^{-1}$)	pr	CRU (Tier 2, Rain gauge analysis, 1901-2010 (Mitchell and Jones, 2005))		
<i>namelist_runoff_et</i>	Total runoff ($\text{kg m}^{-2} \text{s}^{-1}$) Evaporation ($\text{kg m}^{-2} \text{s}^{-1}$) Precipitation ($\text{kg m}^{-2} \text{s}^{-1}$)	mrro evspsbl pr	GRDC (Tier 2, river runoff gauges, varying periods (Dümenil Gates et al., 2000)) WFDEI (Tier 2, Reanalysis, 1979-2010 (Weedon et al., 2014))	Section 4.3.2 / Fig. 19	Dümenil Gates et al. (2000); Hagemann et al. (2013); Weedon et al. (2014)
Section 4.4: Detection of biogeochemical biases: carbon cycle					
<i>namelist_anav13jclim</i>	Net biosphere production of carbon ($\text{kg m}^{-2} \text{s}^{-1}$)	nbp	TRANSCOM (Tier 2, Reanalysis, 1985 - 2008 (Gurney et al., 2004))	Section 4.4.1 / Fig. 20 and Fig. 21	Anav et al. (2013)
	Gross primary production of carbon ($\text{mol m}^{-2} \text{s}^{-1}$)	gpp	MTE (Tier 2, Reanalysis, 1982 - 2008 (Jung et al., 2009))		
	Leaf area index ($\text{mol m}^{-2} \text{s}^{-1}$)	lai	LAI3g (Tier 2, Reanalysis, 1981 - 2008 (Zhu et al., 2013))		
	Carbon mass in vegetation (kg m^{-2})	cVeg	NDP-017b (Tier 2, remote sensing 2000 (Gibbs, 2006))		
	Carbon mass in soil pool (kg m^{-2})	cSoil	HWSD (Tier 2, reanalysis, climatology (Nachtergaele et al., 2012))		
	Primary organic Carbon Production by all types of phytoplankton ($\text{mol m}^{-2} \text{s}^{-1}$)	intPP	SeaWiFS (Tier 2, satellite, 1998-2010 (Behrenfeld and Falkowski,		

			1997; McClain et al., 1998))		
	Near-surface air temperature (K)	tas	CRU (Tier 3, near-surface temperature analysis, 1901-2006)		
	Precipitation ($\text{kg m}^{-2} \text{ s}^{-1}$)	pr	CRU (Tier 2, rain gauge analysis, 1901-2010 (Mitchell and Jones, 2005))		
namelist_GlobalOcean	Surface partial pressure of CO ₂ (Pa)	spco2	SOCAT v2 (Tier 2, in-situ, 1968 - 2011 (Bakker et al., 2014)) ETH SOM-FFN (Tier 2, extrapolated in situ, 1998 - 2011, (Landschützer et al., 2014a, b))	Section 4.4.2 / Fig. 22	
	Total chlorophyll mass concentration at surface (kg m^{-3})	chl	SeaWiFS (Tier 2, satellite, 1997 - 2010 (Behrenfeld and Falkowski, 1997; McClain et al., 1998))		
	Dissolved oxygen concentration (mol m^{-3})	o2	WOA05 (Tier 2, in situ, climatology 1950-2004 (Bianchi et al., 2012))		
	Total alkalinity at surface (mol m^{-3})	talk	T14 (Tier 2, in situ, 2005 (Takahashi et al., 2014))		
Section 4.5: Detection of biogeochemical biases: chemistry and aerosols					
namelist_aerosol_CMIP5	Surface concentration of SO ₄ (kg m^{-3})	sconcs04	CASTNET (Tier 2, Ground, 1987-2012 (Edgerton et al., 1990))	Section 4.5.1 / Fig. 23	Lauer et al. (2005) Aquila et al. (2011) Righi et al. (2013); Fig. 9.29 of Flato et al. (2013)
	Surface concentration of NO ₃ (kg m^{-3})	sconcn03	EANET (Tier 2, Ground, 2001-2005 (Totsuka et al., 2005))		
	Surface concentration of NH ₄ (kg m^{-3})	sconcnh4			
	Surface concentration of black carbon aerosol (kg m^{-3})	sconcbc	EMEP (Tier 2, Ground, 1970-2014)		
	Surface concentration of dry aerosol organic matter (kg m^{-3})	sconcoa			
	Surface concentration of PM10 aerosol (kg m^{-3})	sconcpm10	IMPROVE (Tier 2, Ground, 1988-2014)		
	Surface concentration of PM2.5				

	aerosol (kg m ⁻³)	sconcpm2 p5			
	Aerosol Number Concentration (m ⁻³)	conccn	Aircraft campaigns (Tier 3, aircraft, various)		
	BC Mass Mixing Ratio (kg kg ⁻¹)	mrbc			
	Aerosol mass mixing ration (kg kg ⁻¹)	mmraer			
	BC-Free Mass Mixing Ratio (kg kg ⁻¹)	mmrbcfre e			
	Aerosol Optical Depth at 550 nm (1)	od550aer	AERONET (Tier 2, Ground, 1992-2015 (Holben et al., 1998)) MODIS (Tier 1, satellite, 2001-2012 (King et al., 2003)) MISR (Tier 1, Satellite, 2001-2012 (Stevens and Schwartz, 2012)) ESACCI-AEROSOL (Tier 2, satellite, 1998-2011 (Kinne et al., 2015))		
<i>namelist_righi</i> <i>15gmd_tropo3</i> <i>namelist_righi</i> <i>15gmd_Emmons</i>	Ozone (nmol mol ⁻¹)	tro3	Aura MLS-OMI (Tier 2, satellite, 2005-2013 (Ziemke et al., 2011)) Ozone sondes (Tier 2, sondes, 1995-2009 (Tilmes et al., 2011))	Section 4.5.2 / Fig. 24	Emmons et al. (2000) Righi et al. (2015)
	Carbon Monoxide (mol mol ⁻¹)	vmrco	GLOBALVIEW (Tier 2, ground, 1991-2008, (GLOBALVIEW-CO2, 2008))		
	Nitrogen Dioxide (NO _x = NO + NO ₂) (mol mol ⁻¹)	vmrnox	Emmons (Tier 2, aircraft, various campaign (Emmons et al., 2000))		
	C ₂ H ₄ Propane (mol mol ⁻¹)	vmrc2h4			
	C ₂ H ₆ Propane (mol mol ⁻¹)	vmrc2h6			
	C ₃ H ₆ Propane (mol mol ⁻¹)	vmrc3h6			

	C3H8 Propane (mol mol ⁻¹)	vmrc3h8			
	CH3COCH3 Acetone (mol mol ⁻¹)	vmrch3co ch3			
<i>namelist_eyri ng13jgr</i>	Temperature (K)	ta	ERA-Interim (Tier 3, reanalysis, 1979- 2014 (Dee et al., 2011))	Section 4.5.2 / Fig. 25	Eyring et al. (2013): Fig. 9.10 of Flato et al. (2013)
	Eastward wind (m s ⁻¹)	ua	NCEP (Tier 2, reanalysis, 1948- 2012 (Kistler et al., 2001))		
	Total Column Ozone (DU)	toz	NIWA (Tier 3, sondes, climatology, Bodeker et al., 2005)		
	Tropospheric column ozone (DU)	tropoz	AURA-MLS- OMI (Tier 2, satellite, 2005- 2013 (Ziemke et al., 2011))		
	Ozone (nmol mol ⁻¹)	tro3			
Section 4.6: Linking model performance to projections					
<i>namelist_wen zell4jgr</i>	Near-surface air temperature (K)	tas	NCDC (Tier 2, reanalysis, 1880- 2001 (Smith et al., 2008))	Section 4.6 / Fig. 26	Wenzel et al. (2014); Fig. 9.45 of Flato et al. (2013)
	Net biosphere production of carbon (kg m ⁻² s ⁻¹)	nbp	GCP (Tier 2, reanalysis, 1959- present, (Le Quéré et al., 2014))		
	Carbon Dioxide (mol mol ⁻¹)	co2			
	Surface Downward CO ₂ Flux into ocean (kg m ⁻² s ⁻¹)	fgco2			

1

2

1 Table 2. Overview of the diagnostics included for each namelist along with specific calculations,
2 the plot type, settings in the configuration file (cfg-file), and comments. See also Annex C in the
3 Supplement for additional information.

<i>xml</i> namelist	Diagnostics included	Specific Calculations (e.g., statistical measures, regridding)	Plot Types	Settings in cfg-file	Comments
Section 4.1: Detection of systematic biases in the physical climate: atmosphere					
<i>namelist_perfmetrics_CMI P5</i> <i>namelist_right15gmd_ECVs</i>	perfmetrics_main.ncl	Time averages, Regional weighted averages, t-test for difference plots	Annual cycle line plot, zonal mean plot, lat-lon map plot	Specific plot type, time averaging (e.g. annual, seasonal and monthly climatologies, annual and multi-year monthly means), region, target grid, pressure level, reference model, difference plot (True/False), statistical significance level of t-test for difference plot, multi model mean/median	The results of the analysis are saved to a netCDF file for each model to be read by perfmetrics_grading.ncl or perfmetrics_taylor.ncl.
	perfmetrics_grading.ncl	Grading metric, normalization	No plot	Time averaging, region, pressure level, reference model, type of metric for grading models (RMSE, Bias) type of normalization (mean, median, centered median)	For tractability the filename for every diagnostic is written into a temporary file, which then is read by the perfmetrics_XXX_collect.ncl scripts. Additional metric and normalization methods can be added.
	perfmetrics_taylor.ncl	Taylor metrics	No plot	Time averaging, region, pressure level, reference model	
	perfmetrics_grading_collect.ncl	Collection of model grades from pre-calculated netCDF files	Portrait diagram		If individual models did not provide output for all variables or are compared to a different number of observations, the code will recognize this and return a blank array entry, producing a white box in the portrait diagram;

					produces Figure 9.7 included in <i>namelist_flato13ipcc</i>
	perfmetrics_taylor_collect.ncl	Collection of model grades from precalculated netCDF files	Taylor diagram		
<i>namelist_flato13ipcc</i>	clouds_ipcc.ncl	Multi-model means, linear regridding to the grid of the reference data set	Zonal mean plots, global map	Map projection (CylindricalEquidistant, Mercator, Mollweide), selection of target grid, time mean (annualclim, seasonal-clim), reference data set	Produces Figure 9.5 of Flato et al. (2013) with <i>namelist_flato13ipcc</i>
	clouds_bias.ncl	Multi-model means, linear regridding to the grid of the reference data set	Global map	map projection (CylindricalEquidistant, Mercator, Mollweide), selection of target grid, time mean (annualclim, seasonal-clim), reference data set	Produces Figures 9.2 and 9.4 of Flato et al. (2013) with <i>namelist_flato13ipcc</i>
<i>namelist_SAMonsoon</i>	SAMonsoon_wind_basic.ncl	Mean and interannual standard deviation	Map contour plot, regional mean, RMSE and spatial correlation are given in plot titles	Region (latitude, longitude), season (consecutive month), contour levels	Zonal and meridional wind fields are used; mean and standard deviation (across all years) for each model. This diagnostic also plots the difference of the mean/standard deviation with respect to a reference data set. Mean contour plots include wind vectors.
	SAMonsoon_wind_seasonal.ncl	Climatology, seasonal anomalies and interannual variability	Annual cycle	Region (latitude, longitude), season (consecutive month), line colours, multi model mean (y/n)	Dynamical indices calculated from zonal and meridional wind fields are used. Wind levels are selected by input quantity (e.g. ua-200-850 and va-200-850)
	SAMonsoon_precip_basic.ncl	Mean and interannual standard deviation	Map contour plot, regional mean, RMSE and spatial correlation are given in plot titles	Region (latitude, longitude), season (consecutive month), contour levels	Similar to SAMonsoon_wind_basic.ncl
	SAMonsoon_precip_seasonal.ncl	Climatology, seasonal anomalies and interannual variability	Annual cycle	Region (latitude, longitude), season (consecutive month), line colours, multi model mean (y/n)	Similar to SAMonsoon_wind_seasonal.ncl

	SAMonsoon_precip_domain.ncl	Mean and standard deviation	Map contour plot	Region (latitude, longitude), season (consecutive month), contour levels	Domain and intensity defined using summer and winter precipitation defined appropriately for each hemisphere. Differences from reference data set also plotted. Produces Figure 9.32 included in <i>namelist_flato13ipcc</i>
	SAMonsoon_teleconnections.ncl	Correlation between interannual seasonal mean Nino3.4 SST timeseries (5S-5N, 190-240E) and precipitation over monsoon region.	Map contour plot, regional mean, RMSE and spatial correlation are given in plot titles	Region (latitude, longitude), season (consecutive month), contour levels	pr and ts are used to calculate teleconnections between precip and interannual Nino3.4 SSTs. Differences from reference data set also plotted.
<i>namelist_SAMonsoon_AMIP</i>	SAMonsoon_wind_IAV.ncl	Mean and standard deviation	Time-series line plot	Region (latitude, longitude), season (consecutive month), multi model mean (y/n)	Seasonal means of dynamical indices calculated for each year from zonal and meridional wind fields are used.
	SAMonsoon_precip_IAV.ncl	Mean and standard deviation	Time-series line plot	Region (latitude, longitude), season (consecutive month), multi model mean (y/n)	Seasonal means of precipitation for each year are used. Note that the scripts in <i>namelist_SAMonsoon</i> and <i>namelist_SAMonsoon_daily</i> can be used for coupled and atmosphere-only models alike, but this namelist allows year-to-year variations to be examined only for atmosphere-only simulations forced by observed SSTs.
<i>namelist_SAMonsoon_daily</i>	SAMonsoon_precip_daily.ncl	Standard deviation of filtered daily precipitation rates for each season	Map contour plot. Regional mean, spatial correlation and averages for Bay of Bengal (10-20N, 80-100E) and E. Eq. Indian Ocean (10S-10N, 80-100E) are	Region (latitude, longitude), season (consecutive month), contour levels	Both, actual standard deviations and standard deviations normalized by a climatology (with masking for precipitation rates < 1mm/day) are plotted.

			given in plot titles.		
	SAMonsoon_precip_propagation.ncl	Regional averages, lagged correlations, band-pass filtering of daily precipitation rates	Hovmöller diagrams: (lag, lat) and (lag, lon)	Regions (latitude, longitude), season (consecutive months), filter settings	Similar to <i>namelist_mjo_daily_propagation</i> but using 30-80 day band-pass filtering and regions appropriate for SASM.
<i>namelist_WAMonsoon</i> <i>namelist_WAMonsoon_daily</i>	WAMonsoon_contour_basic.ncl	Mean and standard deviation	Map contour plot	Region (latitude, longitude), season (consecutive months), specific contour levels	Similar to SAMonsoon_wind_basic.ncl
	WAMonsoon_wind_basic.ncl	Mean and standard deviation	Map contour and vector plot	Region (latitude, longitude), season (consecutive months), contour levels, reference vector length	Mean wind contour and vector plots at selected pressure level. Similar to SAMonsoon_wind_basic.ncl
	WAMonsoon_10W10E_1D_basic.ncl	Zonal average over 10°W-10°E	Latitude line plot	Region (latitude), season (consecutive month)	Only 2 dimensional fields
	WAMonsoon_10W10E_3D_basic.ncl	Zonal average over 10°W-10°E	Vertical profile (latitude vs. level) contour plot	Region (latitude, pressure level), season (consecutive month), contour levels	Only 3 dimensional fields
	WAMonsoon_precip_IAV.ncl	Seasonal anomalies and interannual variability	Time-series line plot	Region (latitude, longitude)	Similar to SAMonsoon_wind_IAV.ncl
	WAMonsoon_precip_seasonal.ncl	Mean annual cycle	Time-series line plot	Region (latitude, longitude)	Similar to SAMonsoon_wind_seasonal.ncl
	WAMonsoon_autocorr.ncl	1-day autocorrelation of 1-90d (intraseasonal) anomalies	Map contour plot	Region (latitude, longitude), season (consecutive months), filtering properties, contour levels	
	WAMonsoon_isv_filtered.ncl	Intra-seasonal variance (time filtering)	Map contour plot	Region (latitude, longitude), season (consecutive months), filtering properties, contour levels	
<i>namelist_CVDP</i>	cvdp_atmos.ncl	Renaming climo files to CVDP naming convention, Generates CVDP namelist with all models	No plot		Needed for the CVDP coupling to the ESMValTool.
	cvdp_ocean.ncl	Renaming climo files to CVDP naming convention	No plot		
	cvdp_obs.ncl	Generates	No plot	Reference model(s)	Needed for the CVDP

		CVDP name-list with all observations		for each variable	coupling to the ESMValTool.
	cvdp_driver.ncl	Calls the CVDP	No plot		Needed for the CVDP coupling to the ESMValTool. Flexible implementation for easy update-processes, Results of the analysis are saved in netCDF files for each model/observation
	amo.ncl	Area-weighted average, linear regression, spectral analysis, regridding for area-weighted pattern correlation and RMS difference	Lat-lon contour plots, time-series, spectral plots		Original CVDP diagnostic
	amoc.ncl	Mean, standard deviation, EOF, linear regression, lag correlations, spectral analysis	Pattern plots, spectral plots, time-series		Original CVDP diagnostic
	pdo.ncl	EOF, linear regression, spectral analysis	Lat-lon contour plots, time-series, spectral plots		Original CVDP diagnostic
	pr.mean_stddev.ncl	Global means, standard deviation	Lat-lon contour plots		Original CVDP diagnostic
	pr.trends_timeseries.ncl	Global trends	Lat-lon contour plots, time-series		Original CVDP diagnostic
	psl.mean_stddev.ncl	Global means, standard deviation	Lat-lon contour plots		Original CVDP diagnostic
	psl.modes_indices.ncl	EOF, linear regression,	Lat-lon contour plots, time series		Original CVDP diagnostic
	psl.trends.ncl	Global trends	Lat-lon contour plots		Original CVDP diagnostic
	snd.trends.ncl	Global trends	Lat-lon contour plots		Original CVDP diagnostic
	sst.indices.ncl	Area-weighted average, standard deviation, spectral analysis	Spatial composites, Hovmöller diagram, time-series, spectral plots		Original CVDP diagnostic
	sst.mean_stddev.ncl	Global means, standard deviation	Lat-lon contour plots		Original CVDP diagnostic

	sst.trends_timeseries.ncl	Global trends	Lat-lon contour plots, time-series		Original CVDP diagnostic
	tas.mean_stddev.ncl	Global means, standard deviation	Lat-lon contour plots		Original CVDP diagnostic
	tas.trends_timeseries.ncl	Global trends	Lat-lon contour plots, timeseries		Original CVDP diagnostic
	metrics.ncl	Collect all area-weighted pattern correlations and RMS differences created by the various scripts, calculates total score	txt-file		Original CVDP diagnostic
	webpage.ncl	Creates webpages to display CVDP results	.html files		Original CVDP diagnostic
namelist_mjo_daily	mjo_wave_freq.ncl	Meridional averaged over 10°S-10°N, wavenumber-frequency	Wavenumber -frequency contour plot	Season (summer, winter), daily max/min, region (latitude)	
	mjo_univariate_eof.ncl	Conventional (covariance) univariate EOF analysis	Lat-lon contour plot	Region (latitude, longitude), number and name of EOF modes, contour levels	EOF for 20-100 day band-pass filtered daily anomaly data
	mjo_precip_u850-200_propagation.ncl	Correlation, zonal average over 80°E-100°E, meridional average over 10°S-10°N, reference region over 75°E-100°E, 10°S-5°N	Lag-longitude and lag-latitude diagram	Season(summer, winter, annual), region(latitude, longitude)	Lead/lag correlation of two variables with daily time resolution
	mjo_precip_uwnd_variance.ncl	Variance	Lat-lon contour plot	Season (summer, winter), region (latitude, longitude), contour levels	20-100 day bandpass filtered variance for two variables with daily time resolution
	mjo_olr_u850-200_cross_spectra.ncl	Coherence squared and phase lag	Wavenumber -frequency contour plot	Region (latitude), segments length and overlapped segments length, spectra type	Missing values are not allowed in the input data
	mjo_olr_u850_200_ceof.ncl	CEOF	Line plot	Region(latitude), number and names of CEOF modes, y-axis limit	the first two CEOF modes (PC1 and PC2) are retained for the MJO composite life cycle analysis
	mjo_olr_uv850_c	Calculate mean	Lat-lon	Season (summer,	The appropriate MJO

	eof_life_cycle.ncl	value for each phase category	contour plot	winter), region (latitude, longitude)	phase categories are derived from PC1 and PC2 of CEOF analysis
<i>namelist_mjo_mean_state</i>	mjo_precip_u850_basic.ncl	Season mean	Lat-lon contour plot	Season (summer, winter), region (latitude, longitude)	Based on monthly data
<i>namelist_DiurnalCycle</i>		Mean diurnal cycle computation, regridding of observations and models over a specific grid and first harmonic analysis to derive amplitude and phase of maximum rainfall	Composites of diurnal cycles over specific regions and seasons, global maps of maximum precipitation phase and amplitude		A prerequisite to use this namelist is to check the time axis of high frequency data from models and observations to be sure of what is provided. One should check in particular if it is instantaneous or averaged values, and if the time provided corresponds to the middle or the end of the 3h interval. Note that timeaxis is modified in the namelist to make data coherent.
<i>namelist_laurel13jclim</i>	clouds.ncl	Multi-model mean	Lat-lon contour plot	map projection (CylindricalEquidistant, Mercator, Mollweide), destination grid	Produces Figure 9.5 included in <i>namelist_flato13ipcc</i>
	clouds_taylor.ncl	Multi-model mean	Taylor diagram		Taylor diagrams
	clouds_interannual.ncl	Interannual variability, multi-model mean	Lat-lon contour plot	Map projection (CylindricalEquidistant, Mercator, Mollweide), destination grid, reference data sets	
<i>namelist_williams09climdyn_CREM</i>	ww09_ESMValTool.py	Model data assigned to observed cloud regimes and regime frequency and mean radiative properties calculated.	Bar graph		
Section 4.2: Detection of systematic biases in the physical climate: ocean					
<i>namelist_SouthernOcean</i>	SeaIce_polcon.ncl		Polar stereographic maps	contour values	
	SeaIce_polcon_diff.ncl	Regridding (ESMF)	Polar stereographic maps	contour values, reference model	
	SouthernOcean_vector_polcon_diff.ncl	Vector overlay (magnitude and direction)	Polar stereographic maps	contour plot scales, reference model	based on SeaIce_polcon_diff.ncl, variables with u and v components
	SouthernOcean_a	Regridding	Zonal mean	coordinates of	based on CDFTOOLS

	reamean_vertconplot.ncl	(ESMF)	vertical profiles (Hovmöller diagrams)	subdomain	package
	SouthernOcean_transport.ncl	Sea water volume transport calculation	Line plot	coordinates of subdomain	
namelist_SouthernHemisphere	SouthernHemisphere.py	Regridding (interpolation to common grid), Temporal and zonal averages, RMSEs	Seasonal cycle line plot with calculated RMSEs and zonal mean contour plot	Masking of unwanted values (limits), region (coordinates) and season (months) specification, plotting limits, contour colourmap	
	SouthernHemisphere_scatter.py	Covariability of radiation fluxes as function of cloud metrics	Scatter plot of values with line plot of value distribution		
namelist_TropicalVariability	TropicalVariability.py	Temporal and zonal averages, RMSEs, normalization, co-variability	Annual cycles, seasonal scatter plots with calculated RMSEs	Masking of unwanted values (limits), Region (coordinates) and season (months), plotting limits	Fig. 5 of Lie and Xie, 2014
	TropicalVariability_EQ.py	Temporal and zonal averages, RMSEs, normalization, co-variability	Latitude cross sections of equatorial variables		
	TropicalVariability_wind.py	Regridding (interpolation)	Wind divergence plots		
namelist_Sealce	SeaIce_tsline.ncl	Sea-ice area and extent, regridding (ESMF)	Time series	Selection of Arctic/Antarctic,	Produces Figure 9.24 included in <i>namelist_flato13ipcc</i>
	SeaIce_ancyc.ncl	Sea-ice area and extent, regridding (ESMF)	Annual cycle line plot	Selection of Arctic/Antarctic	
	SeaIce_polcon.ncl	Sea-ice area and extent, regridding (ESMF)	Polar stereographic maps	Selection of Arctic/Antarctic, optional red line depicting edges of sea-ice extent	
	SeaIce_polcon_diff.ncl	Sea-ice area and extent, regridding (ESMF)	Polar stereographic maps	Selection of Arctic/Antarctic, optional red line depicting edges of sea-ice extent	
Section 4.3: Detection of systematic biases in the physical climate: land					
namelist_Evapotranspiration	Evapotranspiration.ncl	Conversion to evapotranspiration units, global average, RMSE	Lat-lon contour plot	Time period	
namelist_SPI	SPI.r	SPI calculation	Lat-lon contour plot	Time period, time scale (3, 6 or 12	May require manual installation of certain

				monthly)	R-packages to run
<i>namelist_runoff_et</i>	catchment_analysis_val.py	Temporal and spatial mean for 12 large river catchments, regridding to 0.5x0.5 lat-lon grid	Bar plots of evapotranspiration and runoff bias against observation, scatter plots of runoff bias against the biases of evapotranspiration precipitation	(no cfg. file)	Three variables are read by this diagnostic.
Section 4.4: Detection of biogeochemical biases: carbon cycle					
<i>namelist_anav13jclim</i>	Anav_MVI_IAV_Trend_Plot.ncl	Regridding to common grid, monthly and annual special averages, variability (MVI = (model/reference - reference/model) 2)	Scatter plot	Region (latitude), resolution size for regridding (e.g., 0.5°, 1°, 2°)	All carbon flux variables were corrected for the exact amount of carbon in the coastal regions by applying the models land-ocean fraction to the variables.
	Anav_Mean_IAV_ErrorBars_Seasonal_cycle_plots.ncl	Regridding to common grid Monthly and annual special averages	Seasonal cycle line plot, scatter plot, error-bar plot	Region (latitude), resolution size for regridding (e.g., 0.5°, 1°, 2°)	
	Anav_cSoil-cVeg_Scatter.ncl	Regridding to common grid annual special averages	Scatter plot	Region (latitude), resolution size for regridding (e.g., 0.5°, 1°, 2°)	Two variables are read by this diagnostic
	perfmetrics_grading.ncl	RMSE, PDF-skill score	No plot		See details in <i>namelist_perfmetrics_CMIP5</i>
	perfmetrics_grading_collect.ncl		Portrait diagram		See details in <i>namelist_perfmetrics_CMIP5</i>
<i>namelist_GlobalOcean</i>	GO_tsline.ncl	Multi-model mean	Time-series line plot	Region (lat/lon), pressure levels, optional smoothing, anomaly calculations, overlaid trend lines, and masking of model data according to observations	
	GO_comp_map.ncl	Mean, standard deviation, and difference to reference model	Lat-lon contour plot (for specified z-level)	Region (Lat/lon), ocean depth, contour levels	Actual metrics ported from UK MetOffice IDL-monsoon evaluation scripts
Section 4.5: Detection of biogeochemical biases: chemistry and aerosols					

<i>namelist_aerosol_CMIP5</i>	aerosol_stations.ncl	Collocation of model and observational data	Time series, scatter plot, map plot	Time averaging, station data network	All available observational data in the selected time period, on a monthly-mean basis is considered. The model data is extracted in the grid boxes where the respective observational stations are located (collocated model and observational data).
	aerosol_satellite.ncl	Regridding to coarsest grid	Map plots and difference plots	Target grid	
	aerosol_profiles.ncl	Mean, standard deviation, median, 5-10-25-75-90-95 percentiles	Vertical profiles		The model data are extracted based on the campaign/station location (lat-lon box) and time period (on a climatological basis, i.e. selecting the same days/months, but regardless of the year). Rather specific variables are required (i.e., aerosol number concentration for particles with diameter larger than 14 nm) to match the properties of the instruments used during the campaign.
	tsline.ncl		Line plot	Time averaging (annual, seasonal and monthly climatologies, annual and multi-year monthly means), region (latitude, longitude)	
<i>namelist_right15gmd_tropo3</i>	ancyc_lat.ncl	Regridding to reference global (area-weighted) average, zonal mean	Seasonal Hovmöller (month vs. latitude)		global (area-weighted) average is calculated only for grid cells with available observational data
	lat_long.ncl	Regridding to coarsest grid global (area-weighted) average			global (area-weighted) average is calculated only for grid cells with available observational data
	perfmetrics_main.ncl		Annual cycle line plot, zonal mean		See details in <i>namelist_perfmetrics_CMIP5</i>

			plot, lat-lon map plot		
	perfmetrics_gradimg.ncl		No plot		See details in <i>namelist_perfmetrics_CMIP5</i>
	perfmetrics_taylor.ncl		No plot		See details in <i>namelist_perfmetrics_CMIP5</i>
	perfmetrics_gradimg_collect.ncl		Portrait diagram		See details in <i>namelist_perfmetrics_CMIP5</i>
	perfmetrics_taylor_collect.ncl		Taylor diagram		See details in <i>namelist_perfmetrics_CMIP5</i>
<i>namelist_righi15gmd_Emmmons</i>	Emmons.ncl	Percentiles (5,25,75,95)%	Vertical profiles	Name(s) of the observational campaign(s)	
<i>namelist_eyring13jgr</i>	anyc_lat.ncl		Seasonal Hovmöller (month vs. latitude)		See details in <i>namelist_righi15gmd_tropo3</i>
	eyring13jgr_fig01.ncl		Seasonal Hovmöller (month vs. latitude)	Multi model mean (True/False), regions (latitude, longitude), time averaging (annual, individual month, seasons)	
	eyring13jgr_fig02.ncl		Time series	Multi model mean (True/False), regions (latitude, longitude), time averaging (annual, individual month, seasons)	Produces Figure 9.10 of Flato et al. (2013) included in <i>namelist_flato13ipcc</i>
	eyring13jgr_fig04.ncl	Tropospheric column ozone	Global maps		
	eyring13jgr_fig06.ncl	Anomalies with respect to a specifiable base line, mean and standard deviation (95% confidence) for simulation experiment	Time series	Multi model mean (True/False), regions (latitude, longitude), time averaging (annual, individual month, seasons)	
	eyring13jgr_fig07.ncl	Mean simulation experiments, differences between future scenario simulations and historical simulations	Vertical profile	Multi model mean (True/False), regions (latitude, longitude), time averaging (annual, individual month, seasons), list of models w/o interactive chemistry	
	eyring13jgr_fig10	Time averages,	Error bar plot	Multi model mean	

	.ncl	linear trends		(True/False), regions (latitude, longitude), height (in km), time averaging (annual, individual month, seasons)	
	eyring13jgr_fig11.ncl	Correlations and correlation coefficient	Scatterplot	Multi model mean (True/False), regions (latitude, longitude), time averaging (annual, individual month, seasons)	Two quantities are compared to each other for individual models and simulations at once. Simulations are indicated by different marker types.

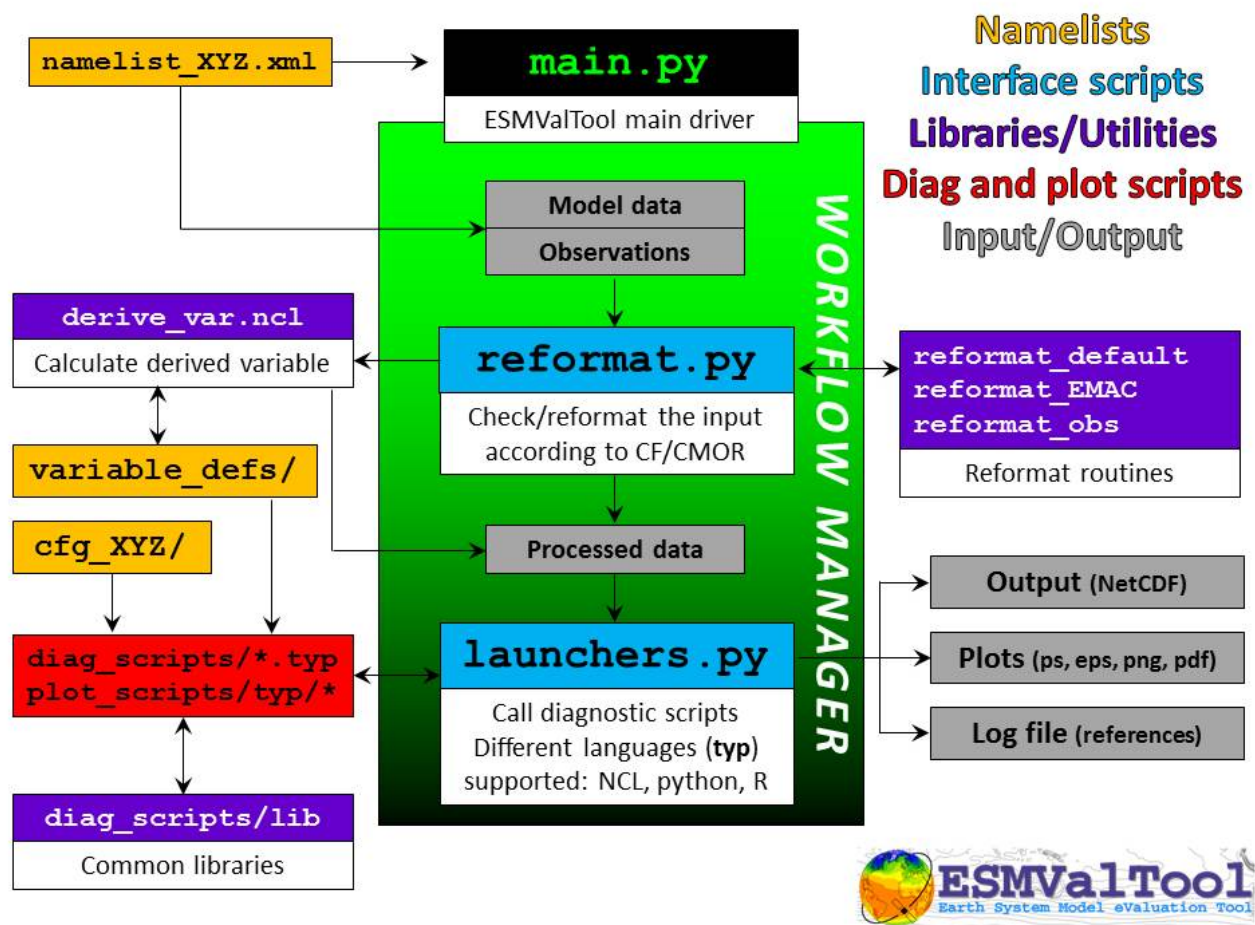
Section 4.6: Linking model performance to projections

namelist_wen zel14jgr	tsline.ncl	Cosine weighting for latitude averaging, anomaly with respect to first 10 years	Line plot	Multi model mean (True/False), anomaly (True/False), regions (latitude, longitude), time averaging (annual, individual month, seasons)	
	carbon_corr_2vars.ncl	Linear regression	Scatter plot and correlation coefficient	Exclude two years after volcanic eruptions (True/False: Mount Agung, 1963; El Chichon, 1982; and Mount Pinatubo, 1991)	Two variables are read. The gradient of the linear regression and the prediction error of the fit, giving γ_{IAV} , are saved in an external netCDF file to be read by the <i>carbon_constraint.ncl</i> script.
	carbon_constraint.ncl	$\gamma_{LT} = \frac{\Delta \text{mbp}^c - \Delta \text{mbp}^u}{\Delta \text{time}^c}$ <p>‘c’ coupled simulation ‘u’ biocemically coupled simulation Gaussian-Normal PDF Conditional PDF</p>	Scatter plot and correlation coefficient	Time period, region (latitude)	Three variables are read. (1) γ_{LT} is diagnosed from the models (2) the previously saved netCDF files containing γ_{IAV} values are read and correlated to γ_{LT} (3) normal and conditional PDFs for the pure model ensemble and the constraint γ_{LT} values are calculated Produces Figure 9.45 included in <i>namelist_flato13ipcc</i>

1

2

1 **FIGURES**



2

3 Figure 1. Schematic overview of the ESMValTool (v1.0) structure. The primary input to the
4 workflow manager is a user-configurable text namelist file (orange). Standardized libraries/utilities
5 (purple) available to all diagnostics scripts are handled through common interface scripts (blue).
6 The workflow manager runs diagnostic scripts (red) that can be written in several freely-available
7 scripting languages. The output of the ESMValTool (gray) includes figures, binary files (netCDF),
8 and a log-file with a list of relevant references and processed input files for each diagnostic.

9

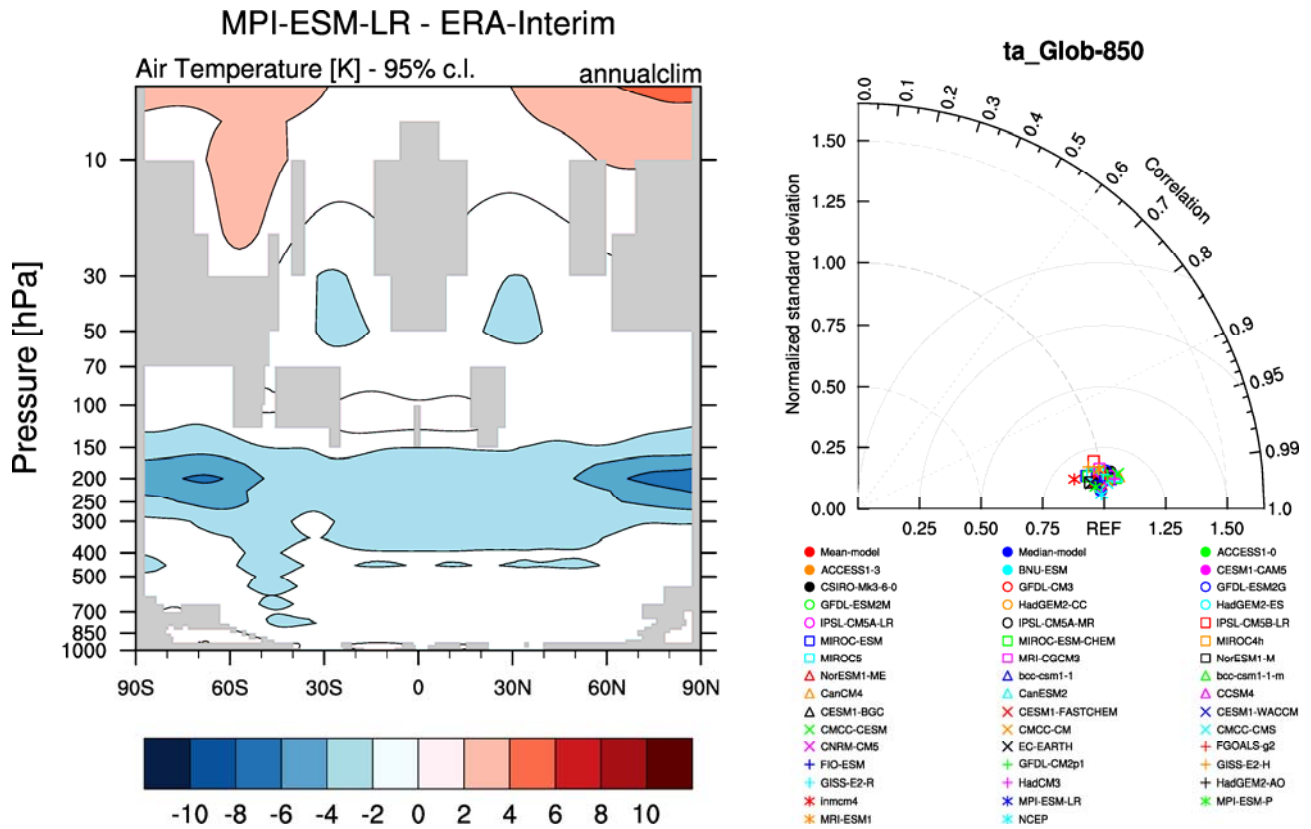
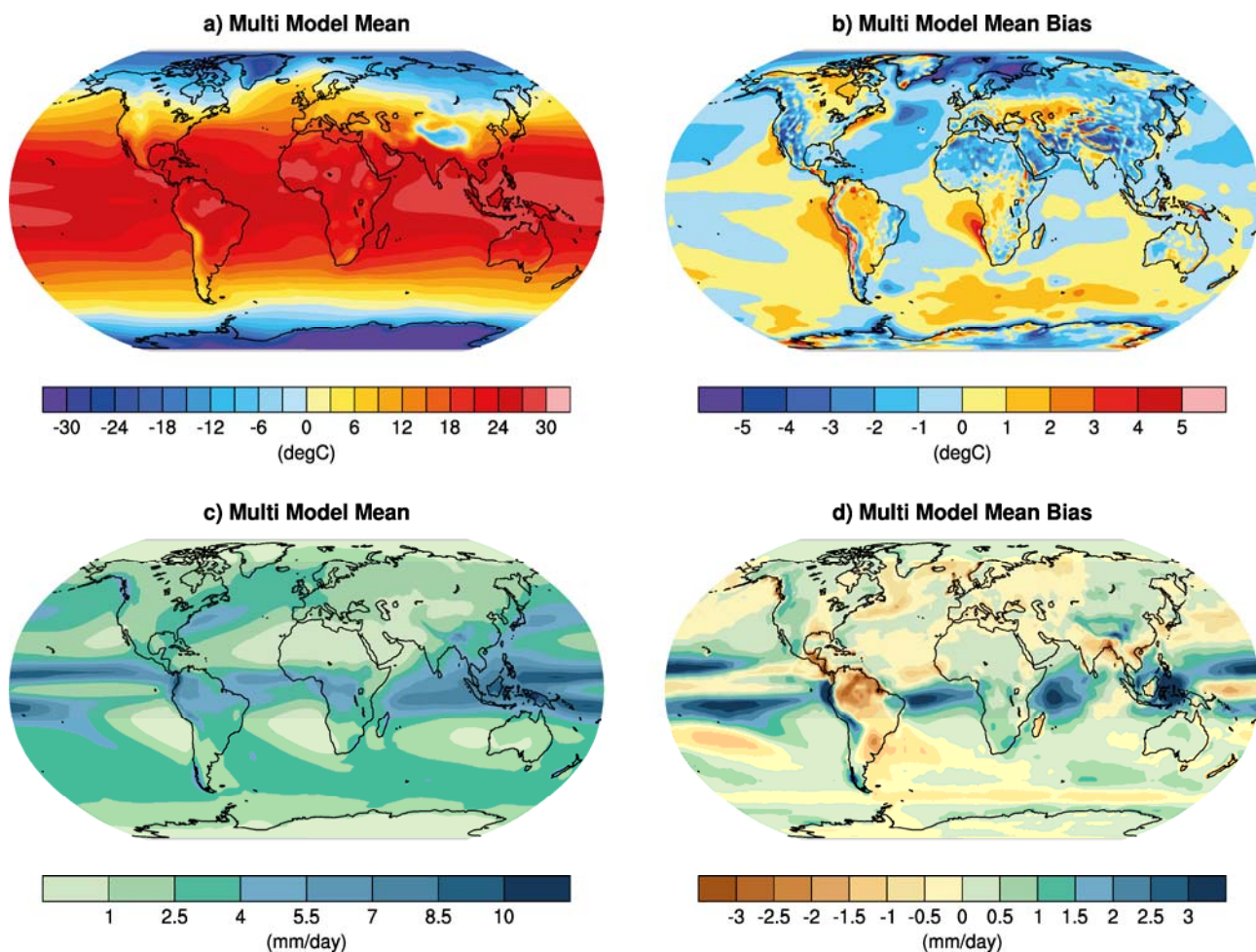


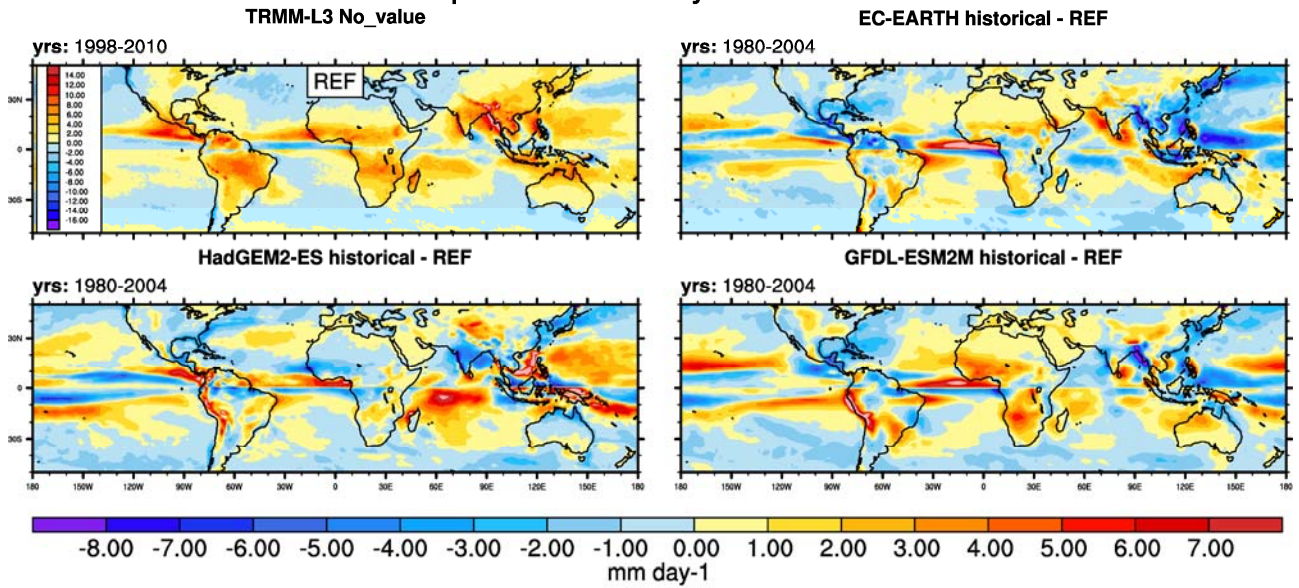
Figure 3. *Left*: Zonally averaged temperature profile difference between MPI-ESM-LR and the ERA-Interim reanalysis data with masked non-significant values. MPI-ESM-LR has generally small biases in the troposphere ($< 1\text{--}2$ K), but a cold bias in the tropopause region that is particularly strong in the extratropical lower stratosphere. This is a systematic bias present in many of the CMIP3 and CCMVal models (IPCC, 2007; SPARC-CCMVal, 2010), related to an overestimation of the water vapour concentrations in that region. *Right*: Taylor diagram for temperature at 850 hPa from CMIP5 models compared with ERA-Interim (reference observation-based data set) and NCEP (alternate observation-based data set) showing a very high correlation or $R > 0.98$ with the reanalyses demonstrating very good performance in this quantity. Both figures produced with *namelist_perfmetrics_CMIP5.xml*.



1
 2 Figure 4. Annual-mean surface air temperature (upper row) and precipitation rate (mm day^{-1}) for
 3 the period 1980–2005. The left panels show the multi-model mean and the right panels the bias as
 4 the difference between the CMIP5 multi-model mean and the climatology from ERA-Interim (Dee
 5 et al., 2011) and the Global Precipitation Climatology Project (Adler et al., 2003) for surface air
 6 temperature and precipitation rate, respectively. The multi-model mean near-surface temperature
 7 agrees with ERA-Interim mostly within $\pm 2^\circ\text{C}$. Larger biases can be seen in regions with sharp
 8 gradients in temperature, for example in areas with high topography such as the Himalaya, the sea
 9 ice edge in the North Atlantic, and over the coastal upwelling regions in the subtropical oceans.
 10 Biases in the simulated multi-model mean precipitation include too low precipitation along the
 11 equator in the western Pacific and too high precipitation amounts in the tropics south of the equator.
 12 Similar to Figures 9.2 and 9.4 of Flato et al. (2013) and produced with *namelist_flato13ipcc.xml*.

13

Monsoon Precipitation Intensity: Model minus Reference



Monsoon Global Domain: Model minus Reference

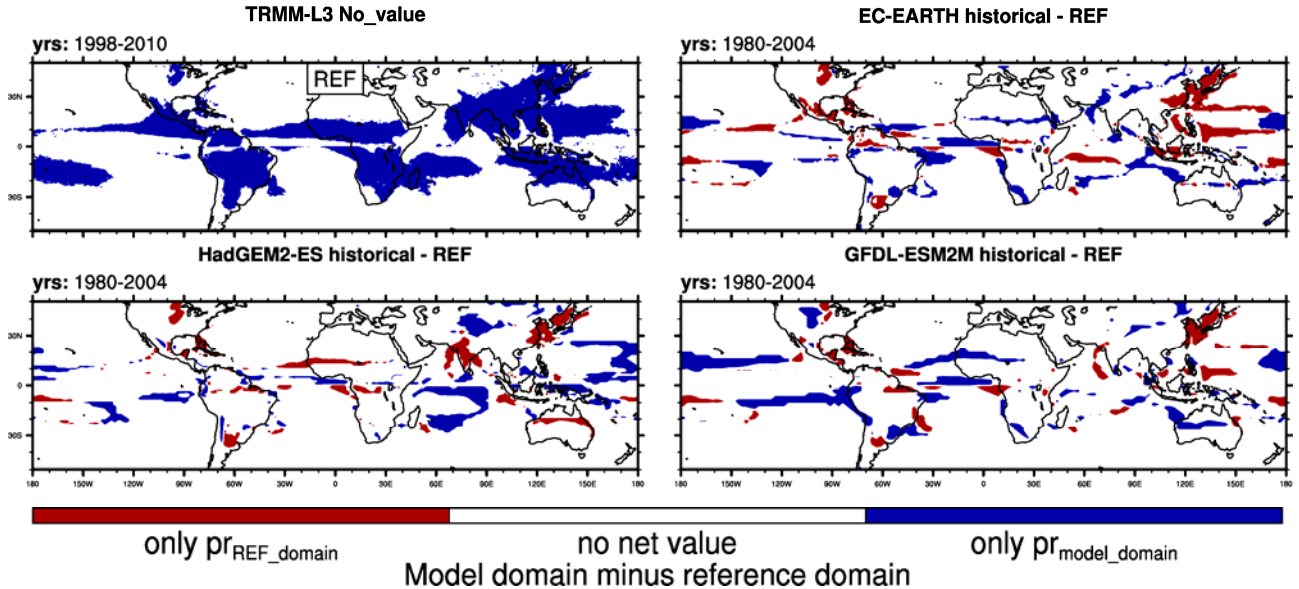
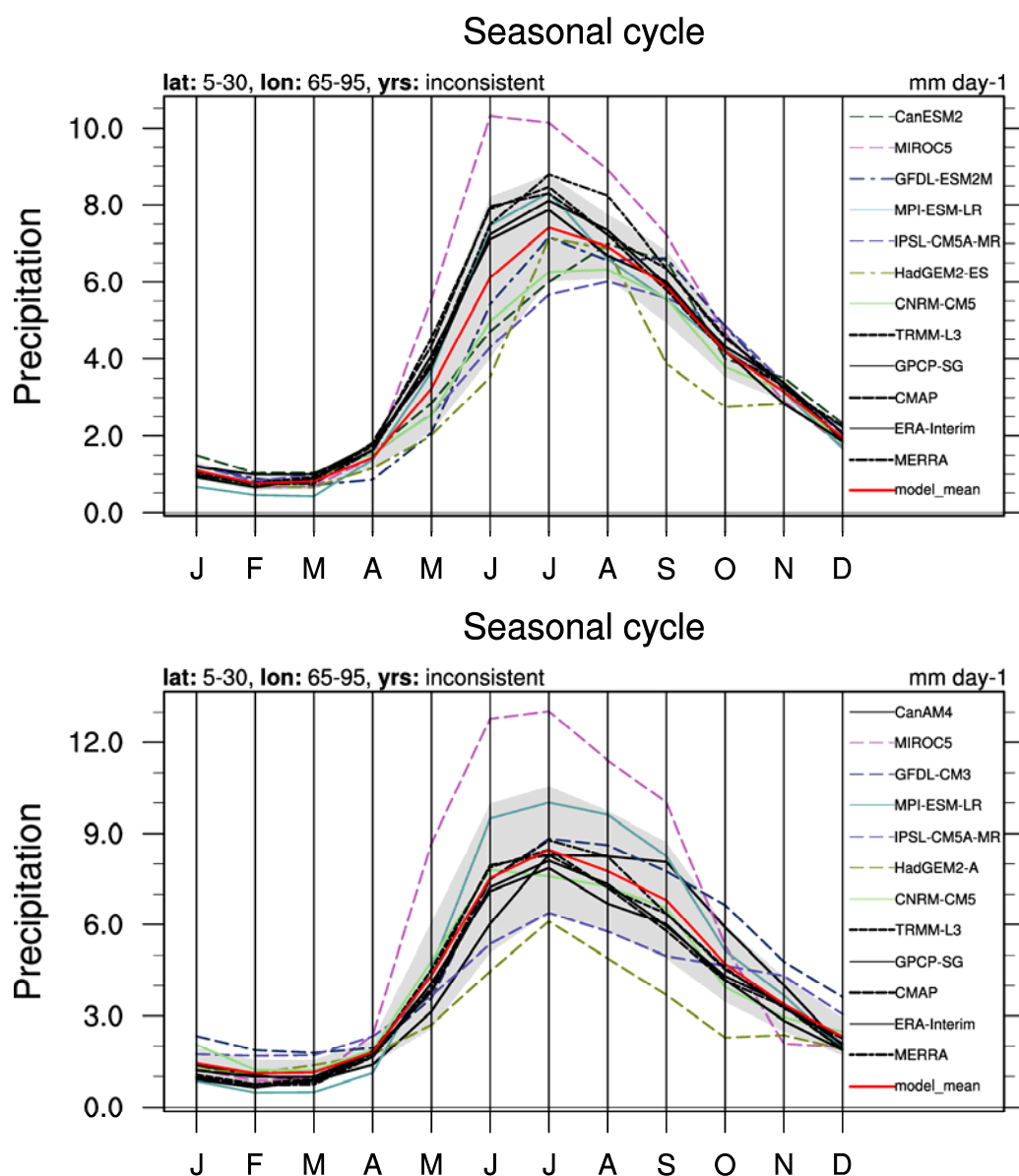
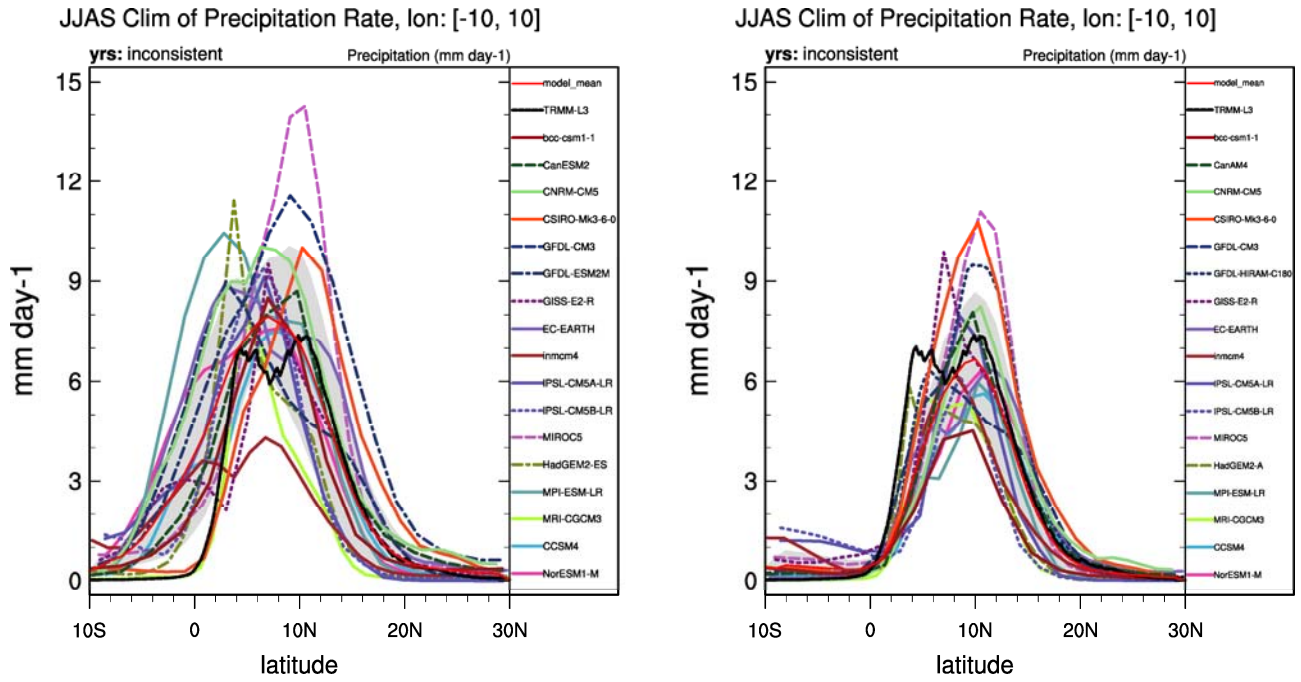


Figure 5. Monsoon precipitation intensity (upper panels) and monsoon precipitation domain (lower panels) for TRMM and an example of deviations from observations from three CMIP5 models (EC-Earth, HadGEM2-ES, and GFDL-ESM2M). The models have difficulties representing the eastward extent of the monsoon domain over the South China Sea and western Pacific, and several models (e.g., HadGEM2-ES) underestimate the latitudinal extent of most of the monsoon regions. The monsoon precipitation intensity tends to be underestimated in the South Asian, East Asian and Australian monsoon regions while in the African and American monsoon regions the sign of the intensity bias varies between models. Similar to Figure 9.32 of Flato et al. (2013) and produced with *namelist_SAMonsoon.xml*.



1
2 Figure 6. Seasonal cycle of monthly rainfall averaged over the Indian region (5-30°N, 65-95°E) for
3 a range of CMIP5 coupled models (upper panel) and their AMIP counterparts (lower panel),
4 averaged over available years (models: 1980-2004, observations: 1998-2010). The grey area in each
5 panel indicates standard deviation from the model mean, to indicate the spread between models
6 (observations/reanalyses are not included in this spread). These illustrate the range of rainfall
7 simulated particularly in AMIP experiments where there is no feedback between precipitation and
8 SST biases that might moderate the rainfall biases (Bollasina and Ming, 2013; Levine et al., 2013).
9 Some of the CMIP5 coupled models (e.g., HadGEM2-ES, IPSL-CM5A-MR) show a delayed
10 monsoon onset that is not apparent in their AMIP configurations. This is related to cold SST biases
11 in the Arabian Sea which develop during boreal winter and spring (Levine et al., 2013). Produced
12 with *namelist_SAMonsoon.xml*.



1
2 Figure 7. Precipitation (mm day⁻¹) averaged over 10°W-10°E for the JJAS season for the years
3 1979-2005 for CMIP5 historical simulations (left) and 1979-2008 for CMIP5 AMIP simulations
4 (right) compared to 1998-2008 for TRMM 3B43 Version 7 data set. The results illustrate the inter-
5 model spread in the mean position and intensity of the WAM among the CMIP5 models. The
6 spread is slightly reduced in AMIP simulations, as the warm SST bias in the equatorial Atlantic is
7 removed. The WAM mean structure, however, is not captured by many models. Produced with
8 *namelist_WAMonsoon.xml*.

PDO (Monthly)

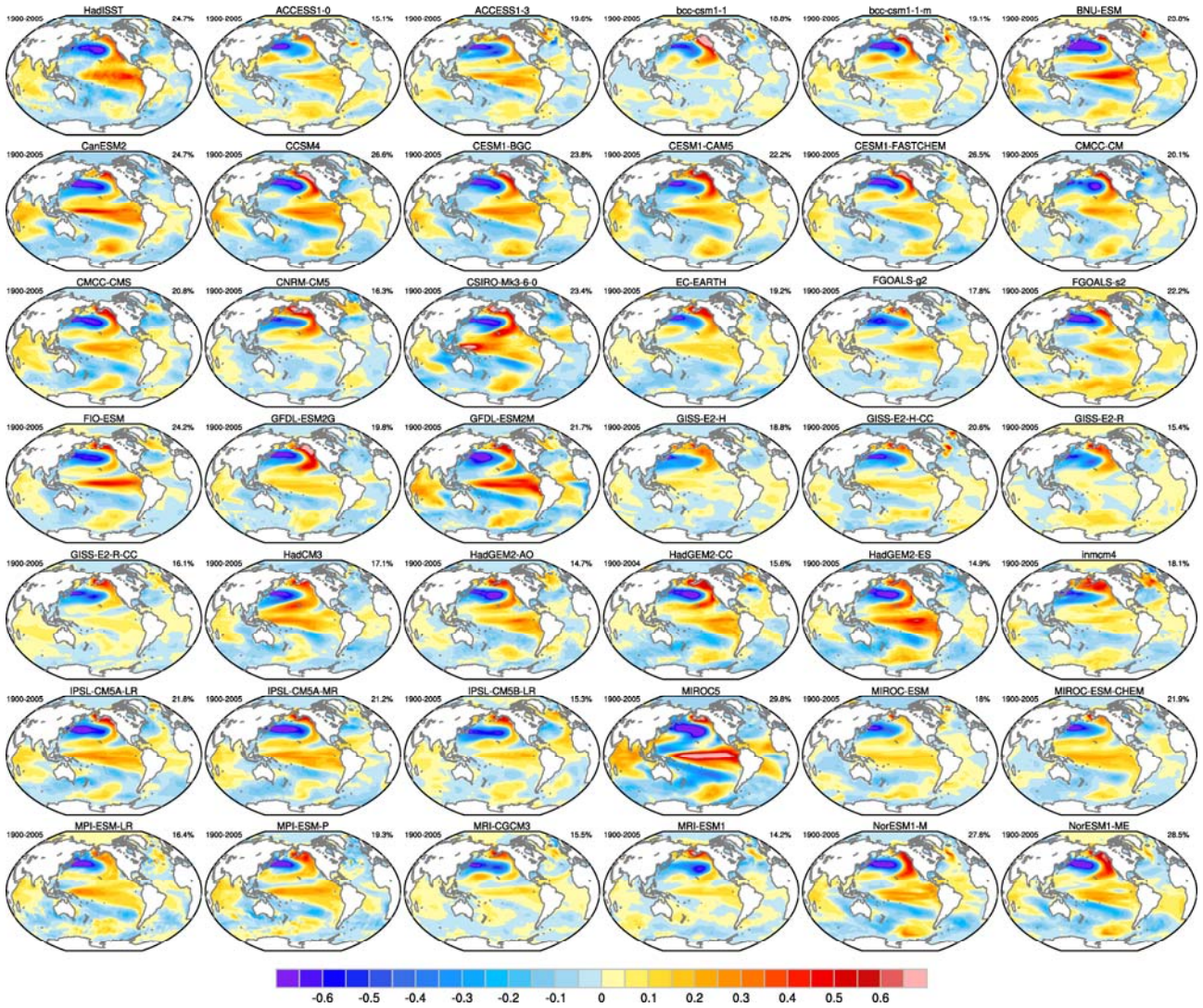
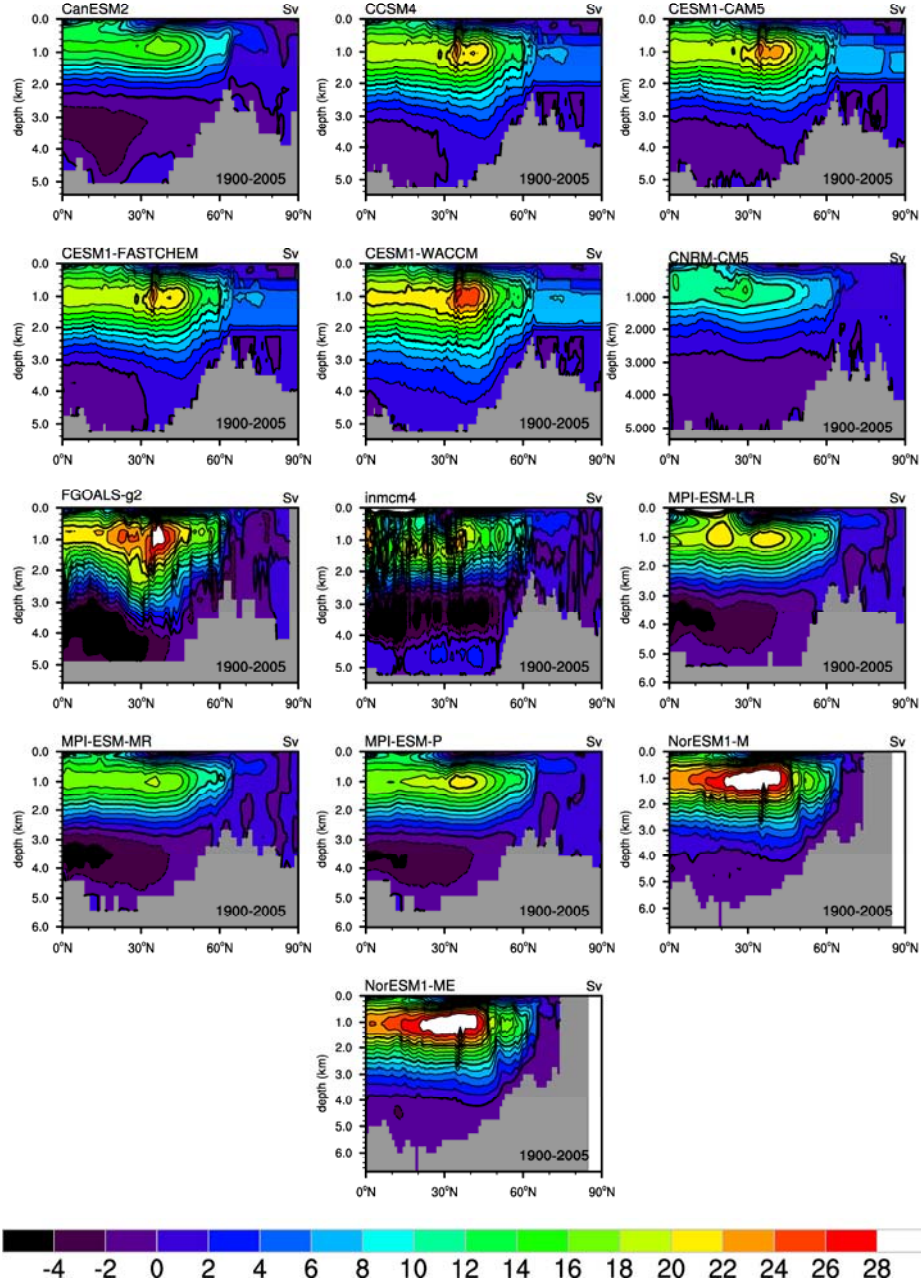


Figure 8. The PDO as simulated by 41 CMIP5 models (individual panels labelled by model name) and observations (upper left panel) for the historical period 1900-2005. These patterns show the global SST anomalies (°C) associated with a one standard deviation change in the normalized principal component (PC) time series. The percent variance accounted by the PDO is given in the upper right of each panel. The PDO is defined as the leading empirical orthogonal function of monthly SST anomalies (minus the global mean SST) over the North Pacific (20-70°N, 110°E-100°W). The global patterns (°C) are formed by regressing monthly SST anomalies at each grid point onto the PC time series. Most CMIP5 models show realistic patterns in the North Pacific. However, linkages with the tropics and the tropical Pacific in particular, vary across models. The lack of a strong tropical expression of the PDO is a major shortcoming in many CMIP5 models (Flato et al., 2013). Figure produced with *namelist_CVDP.xml*.

AMOC Means (Annual)



1

2 Figure 9. Long-term annual mean Atlantic Meridional Overturning Streamfunction (AMOC; Sv) as
 3 simulated by 13 CMIP5 models (individual panels labelled by model name) for the historical period
 4 1900-2005. AMOC annual averages are formed, weighted by the cosine of the latitude and by the
 5 depth of the vertical layer, and then the data is masked by setting all those areas to missing where
 6 the variance is less than $1.e^{-6}$. The figure shows that there is a wide spread among the CMIP5
 7 models, with maximal AMOC strength ranging from ~ 13 Sv (CanESM2) to over ~ 28 Sv
 8 (NorESM1), while the models agree generally well on the position of maximal AMOC strength.
 9 Figure produced with *namelist_CVDP.xml*.

10

summer

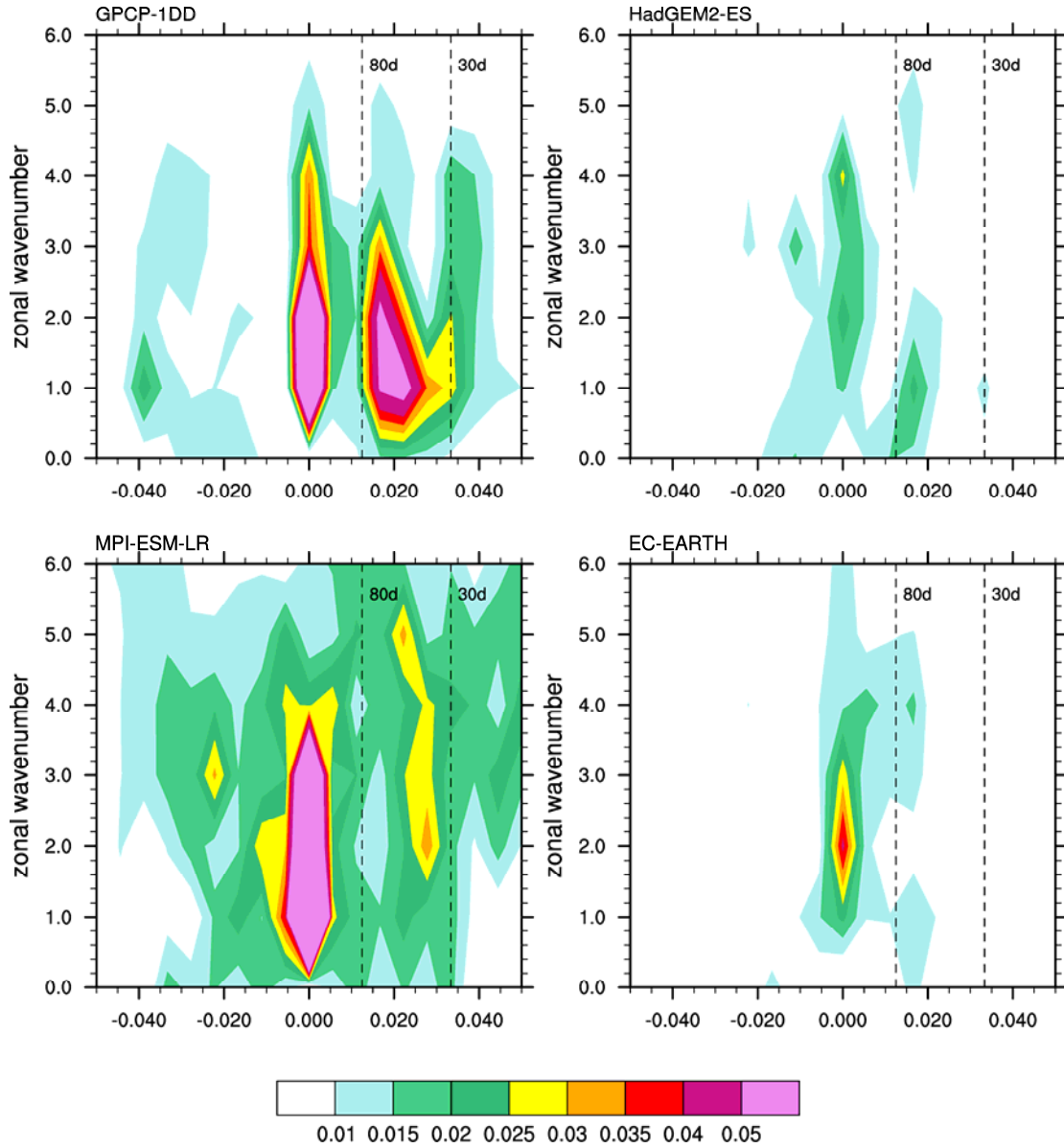
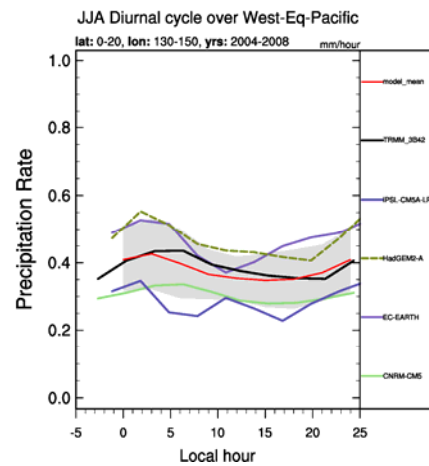
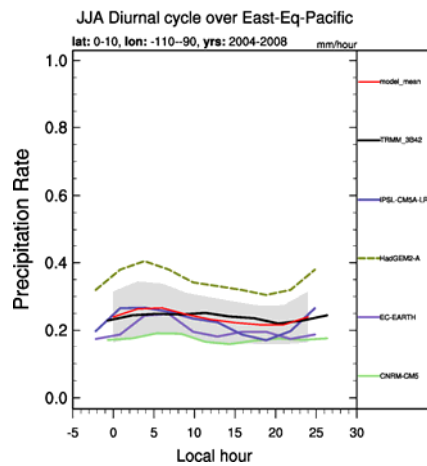
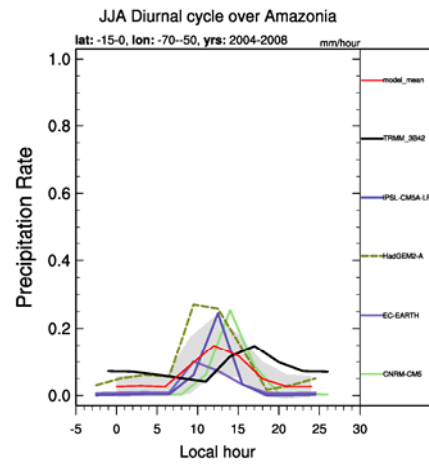
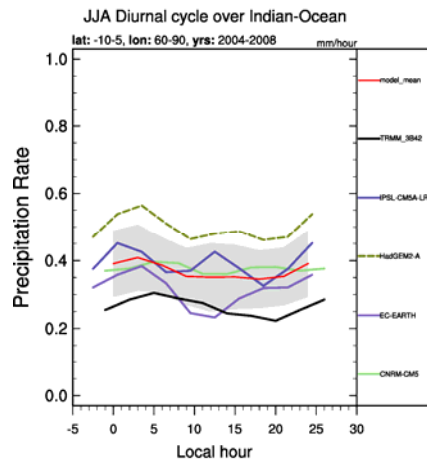
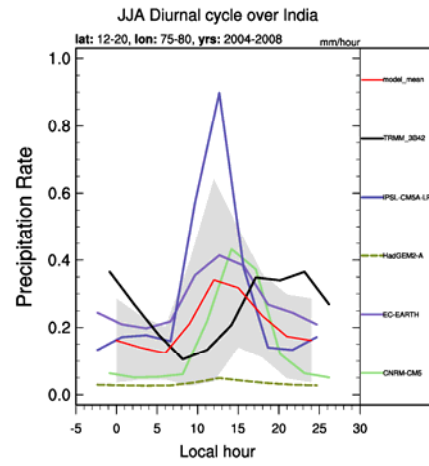
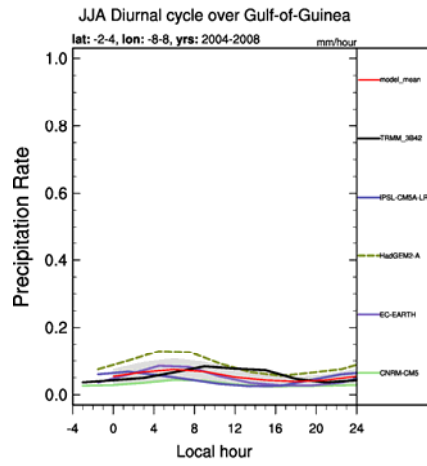
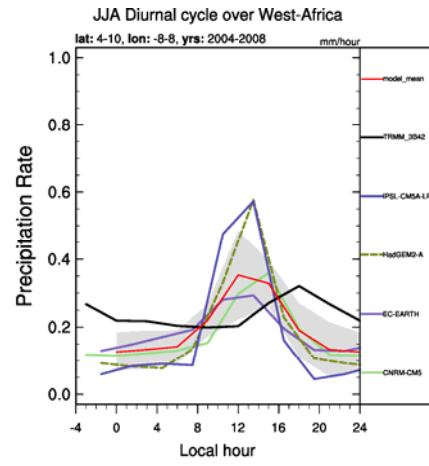
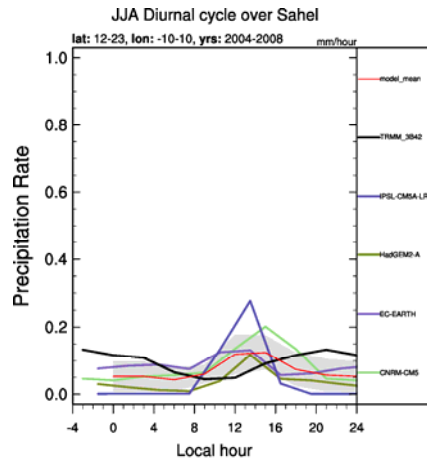
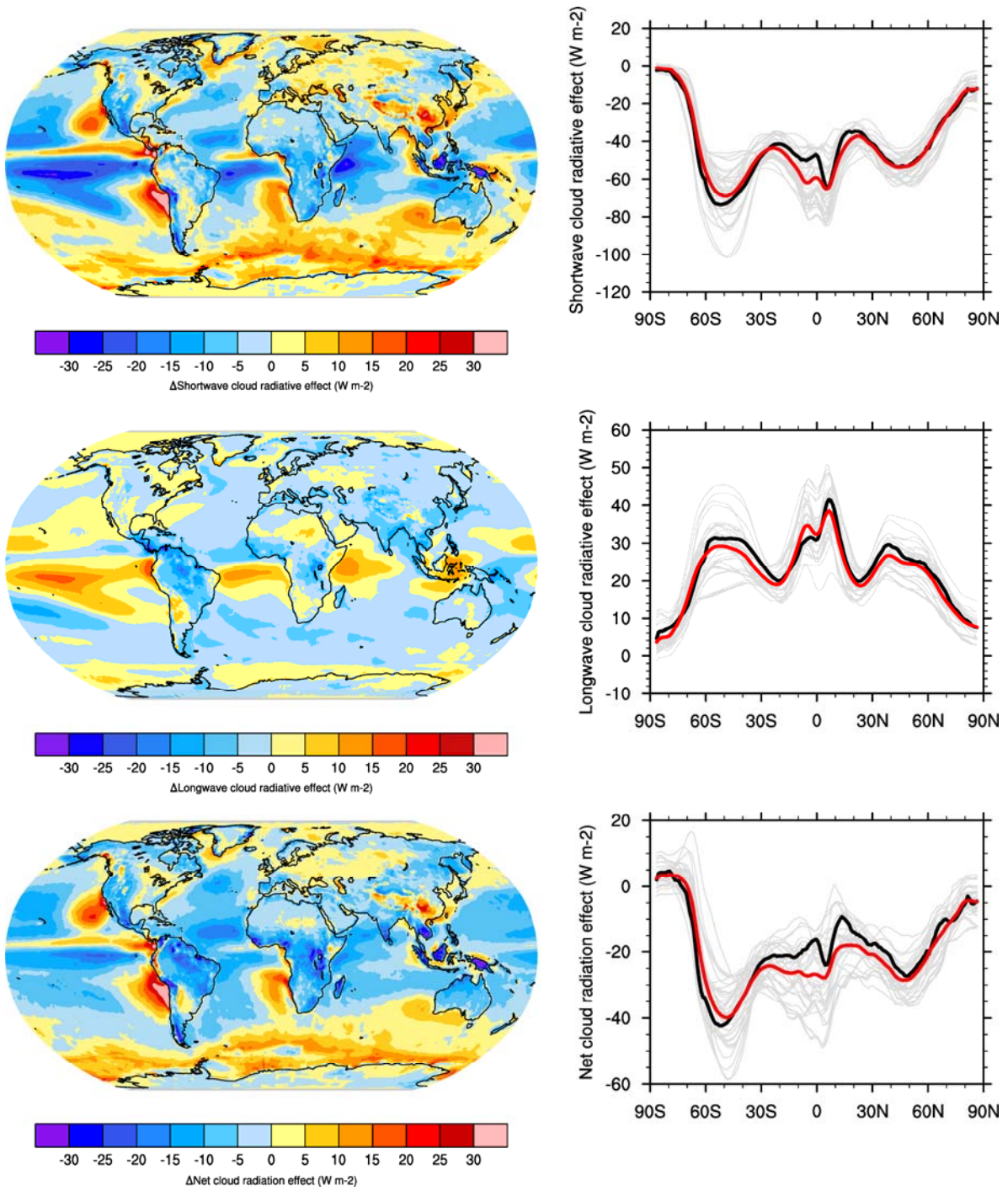


Figure 10. May-October wavenumber-frequency spectra of 10°S-10°N averaged precipitation ($\text{mm}^2 \text{ day}^{-2}$) for GPCP-1DD, HadGEM2-ES, MPI-ESM-LR and EC-Earth. Individual May-October spectra were calculated for each year and then averaged over all years of data. Only the climatological seasonal cycle and time mean for each May-October segment were removed before calculation of the spectra. The bandwidth is $(180 \text{ days})^{-1}$. The observed precipitation shows the dominant MJO spatial scale is zonal wavenumber 1-3 at the 30-80-day frequency. According to the definition, the positive frequency represents eastward propagation of the MJO. Compared with observations, both HadGEM2-ES and EC-Earth models have difficulties simulating precipitation variability on MJO timescales. Produced with *namelist_mjo_daily.xml*.



1 Figure 11. Mean diurnal cycle of precipitation (mm/hour) averaged over five summers (2004-2008)
2 over specific regions in the tropics (Sahel, West-Africa, Gulf of Guinea, India, Indian Ocean,
3 Amazonia, East-Equatorial Pacific and West-Equatorial Pacific) as observed by TRMM 3B42 V7
4 and as simulated by four CMIP5 models: CNRM-CM5, EC-Earth, HadGEM2-A and IPSL-CM5A-
5 LR. ESMs produce a too strong peak of rainfall around noon over land while the observed
6 precipitation maximum is weaker and delayed to 6 pm. At the same time, most models
7 underestimate nocturnal precipitation. Over the ocean, the diurnal cycle of precipitation is more flat
8 but rainfall maximum usually occurs a few hours earlier than in observations during the night, and
9 the amplitude of oceanic precipitation shows large variations among models. Produced with
10 *namelist_DiurnalCycle_box_pr.xml*.

11



1

2 Figure 12. Climatological (1985-2005) annual-mean cloud radiative effects from the CMIP5 models
 3 against CERES EBAF (2001–2012) in W m^{-2} . Top row shows the shortwave effect; middle row the
 4 longwave effect, and bottom row the net effect. Multi-model-mean biases against CERES EBAF
 5 2.7 are shown on the left, whereas the right panels show zonal averages from CERES EBAF 2.7
 6 (black), the individual CMIP5 models (thin grey lines), and the multi-model mean (thick red line).
 7 The multi-model mean longwave CRE is overestimated in models, particularly in the Pacific and
 8 Atlantic south of the inter-tropical convergence zone (ITCZ) and in the South Pacific convergence

1 zone (SPCZ). The longwave CRE is underestimated over Central and South America as well as
2 parts of Central Africa and southern Asia. The most striking biases in the multi-model mean
3 shortwave CRE are found in the stratocumulus regions off the west coasts of North and South
4 America, southern Africa, and Australia. Despite biases in component cloud properties, simulated
5 CRE is in quite good agreement with observations. Reproducing Figure 9.5 of Flato et al. (2013)
6 and produced with *namelist_flato13ipcc.nml*.

7

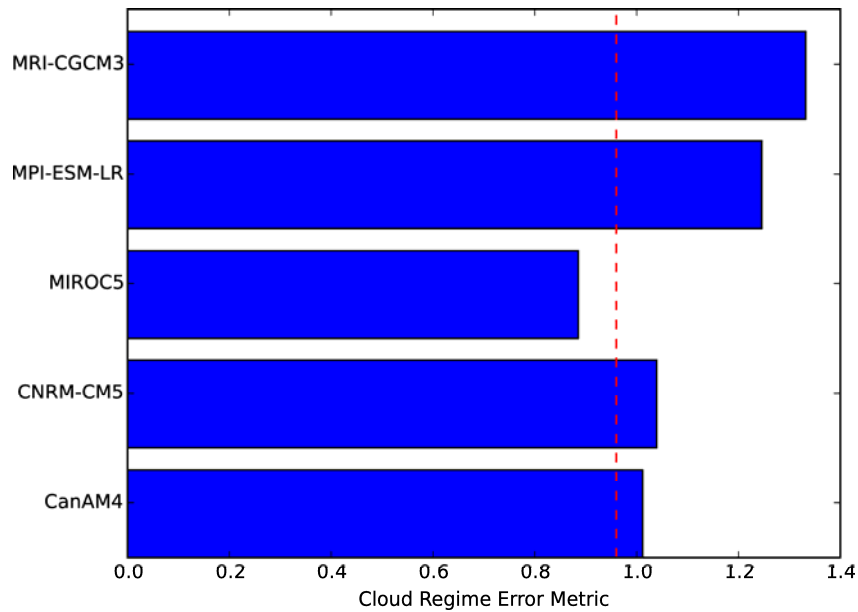


Figure 13. Cloud Regime Error Metric (CREM) from Williams and Webb (2009) applied to some CMIP5 AMIP simulations with the required data in the archive. The results show that MIROC5 is the best performing model on this metric, other models are slightly worse on this metric. The red dashed line shows the observational uncertainty estimated from applying this metric to independent data from MODIS. An advantage of the metric is that its components can be decomposed to investigate the reasons for poor performance. This requires extra print statements compared to the default code but might help to identify, for instance, cloud regimes that are too reflective or simulated too frequently at the expense of some of the other regimes. Produced with *namelist_williams09climdyn_CREM.xml*.

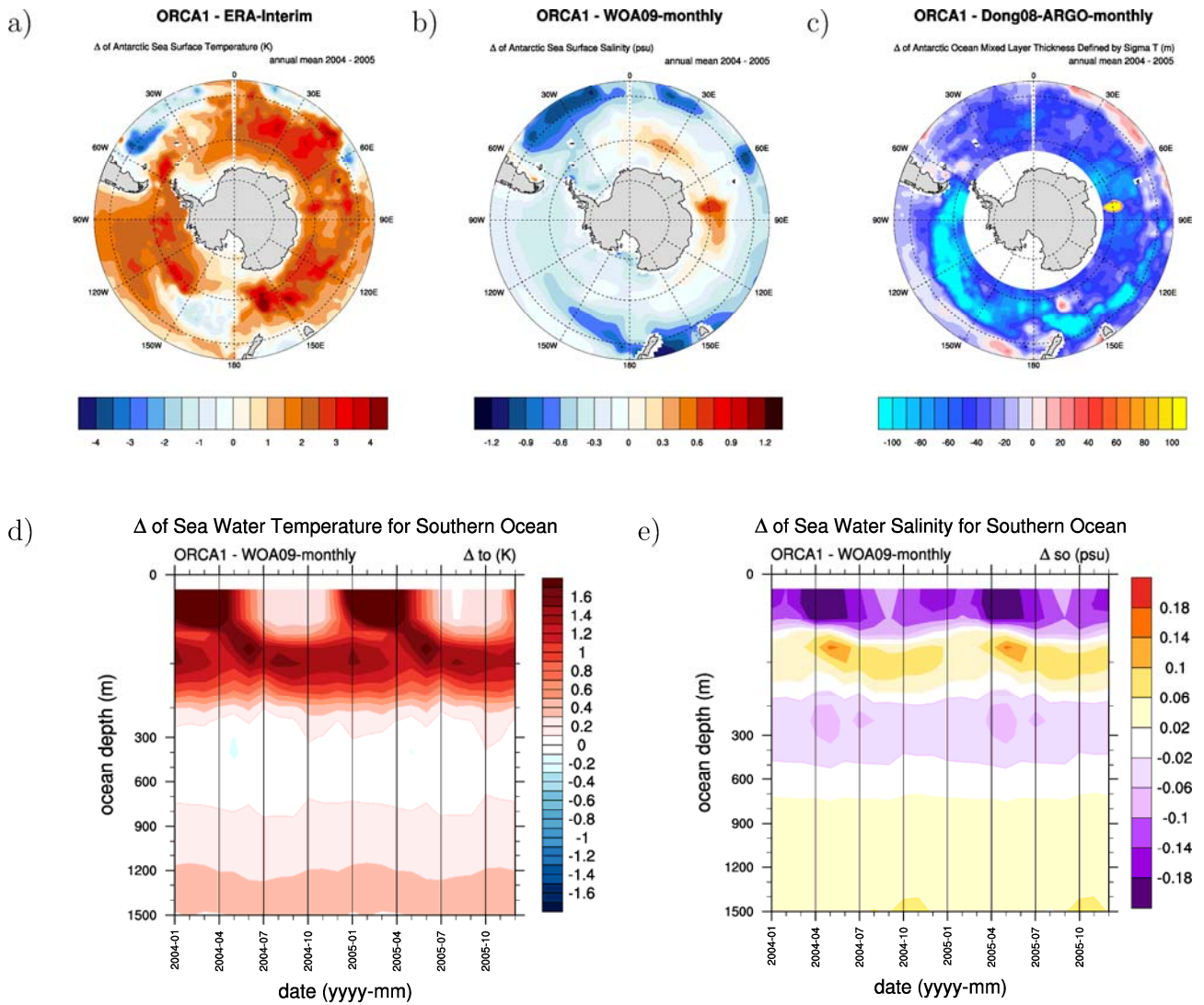
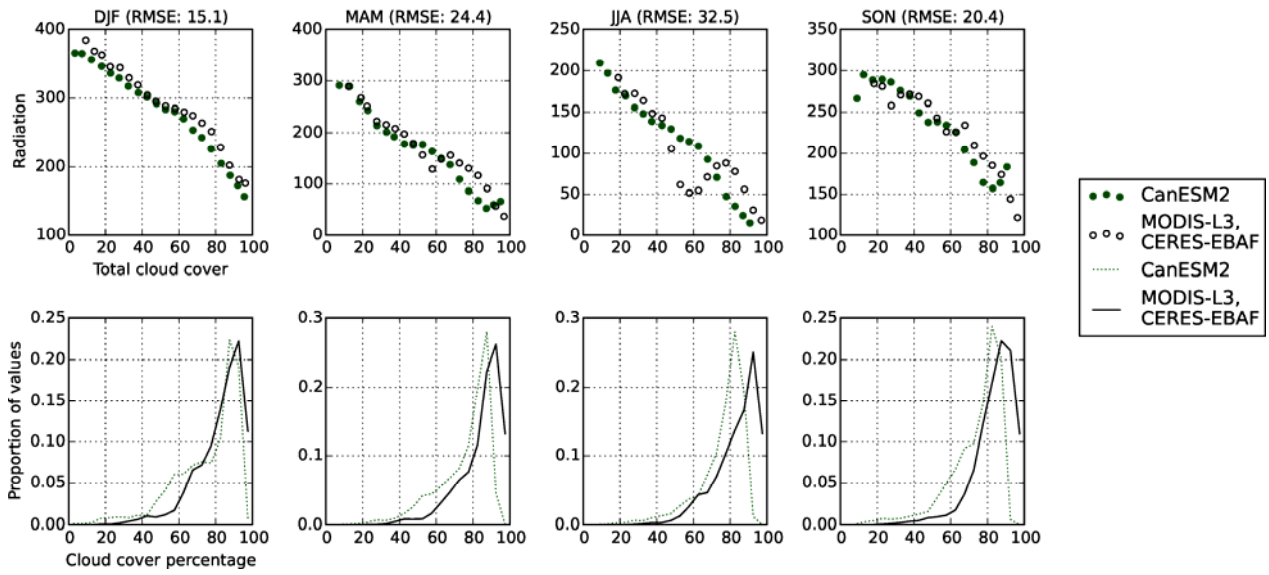


Figure 14. Annual-mean difference between EC-Earth/NEMO and ERA-Interim sea surface temperatures (a), the World Ocean Atlas sea surface salinity (b), and the Argo float observations for ocean mixed layer thickness (c), showing that in the Southern Ocean SSTs in EC-Earth are too high, sea surface salinity too fresh, and the mixed layer too shallow. The other available diagnostics of the *namelist_SouthernOcean.nml* help understanding these biases. Vertical sections of temperature (d) and salinity differences (e) reveal that the SST bias is mainly an austral summer problem, but also that vertical mixing is not able to penetrate a year-round existing warm layer below 80 m depth.

Surface incoming shortwave radiation sensitivity to Total cloud cover



1

2 Figure 15. Upper panel: covariability between incoming surface short wave radiation (rsds) and
 3 total cloud cover (clt). Lower panel: fraction occurrence histograms of binned cloud cover:
 4 observations are CERES-EBAF (radiation) and CloudSat (cloud cover). The CanESM2 model from
 5 the CMIP5 archive is shown as an example for comparison to observations (the namelists runs on
 6 all CMIP5 models). CanESM2 generally reproduces the observed slope of rsds as a function of clt,
 7 although there is a systematic positive bias in the amount of shortwave radiation reaching the
 8 surface for most cloud cover values. A positive bias is also seen in the CanESM2 histogram of
 9 cloud occurrence, with a strong peak in seasonal cloud fraction of 90% in most seasons. Produced
 10 with *namelist_SouthernHemisphere.xml*.

11

Pacific ocean [120E:100W] seasonal mean

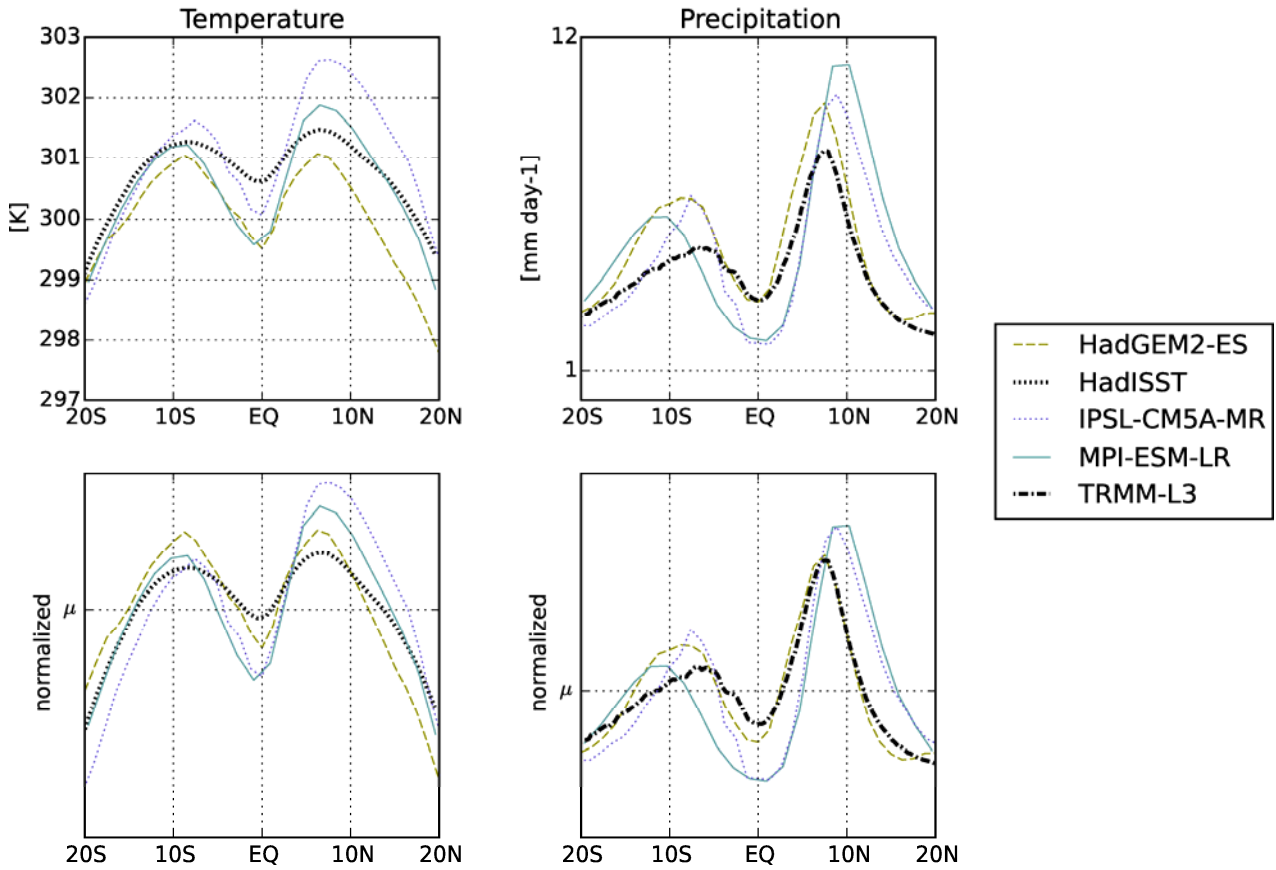
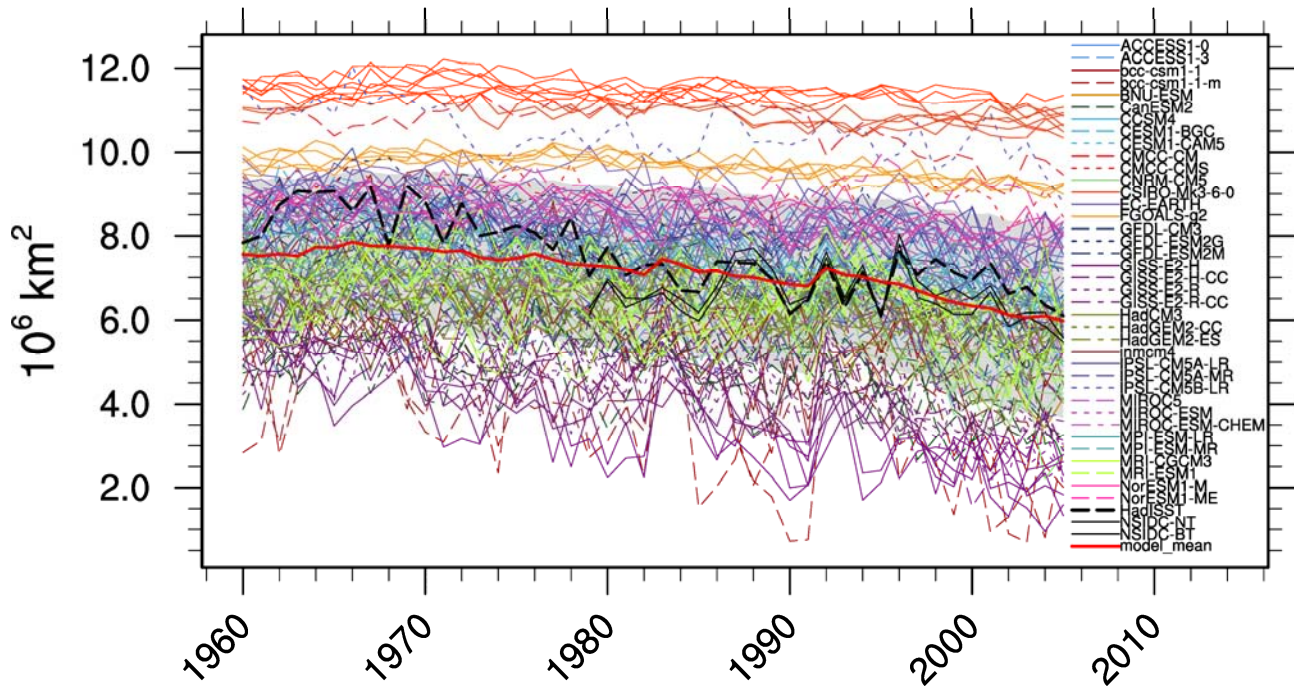


Figure 16. Latitude cross-section of seasonal and zonally averaged values of SSTs and precipitation for the tropical Pacific (zonal averages are made between 120°E and 100°W). Upper panel shows absolute values of SST and precipitation, lower panel shows values normalized by their respective tropical mean value (20°N to 20°S) The figure shows that HadGEM2-ES simulates a double ITCZ in the equatorial Pacific with excessive precipitation south of the equator. This bias is accompanied by off equatorial warm biases in normalized SST in both hemispheres and a relative cold bias along the equator. The IPSL-CM5A-MR and MPI-ESM-LR models better capture the SST and precipitation distributions in the tropical Pacific. Produced with *namelist_TropicalVariability.xml*.

September Arctic Sea Ice Extent



1
 2 Figure 17. Timeseries (1960-2005) of September mean Arctic sea ice extent from the CMIP5
 3 historical simulations. The CMIP5 ensemble mean is highlighted in dark red and the individual
 4 ensemble members of each model (coloured lines) are shown in different linestyles. The model
 5 results are compared to observations from the NSIDC (1978-2005, black solid line) and the Hadley
 6 Centre Sea ice and Sea Surface Temperature (HadISST, 1960-2005, black dashed line). Consistent
 7 with observations, most CMIP5 models show a downward trend in sea ice extent over the satellite
 8 era. The range in simulated sea ice is however quite large (between 3.2 and 12.1 x 10⁶ km² at the
 9 beginning of the timeseries). The multi-model-mean lies below the observations throughout the
 10 entire time period, especially after 1978, when satellite observation became available. Similar to
 11 upper left panel of Figure 9.24 of Flato et al. (2013) and produced with *namelist_Sealce.nml*.

12

Jul-mean of Evapotranspiration

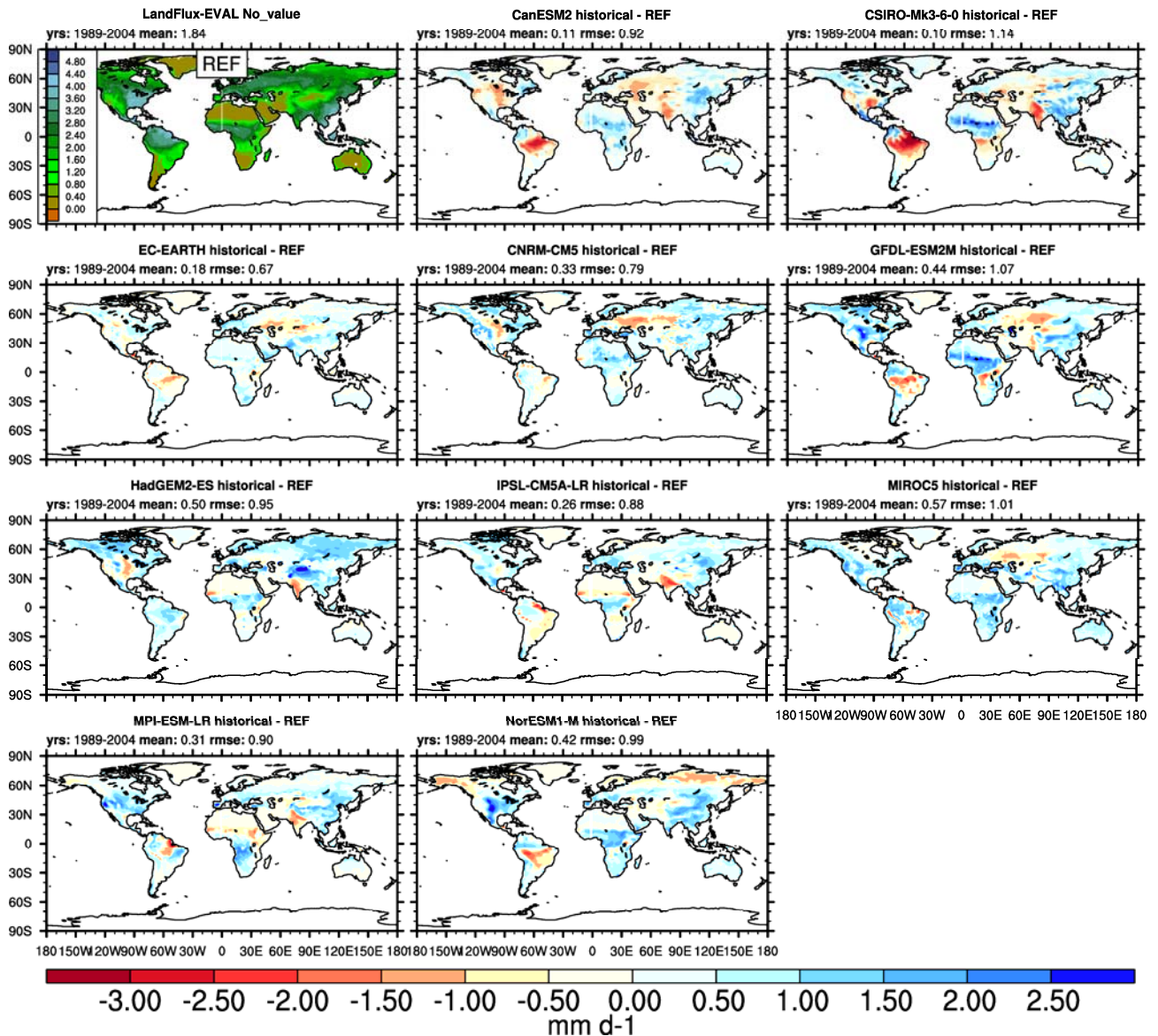
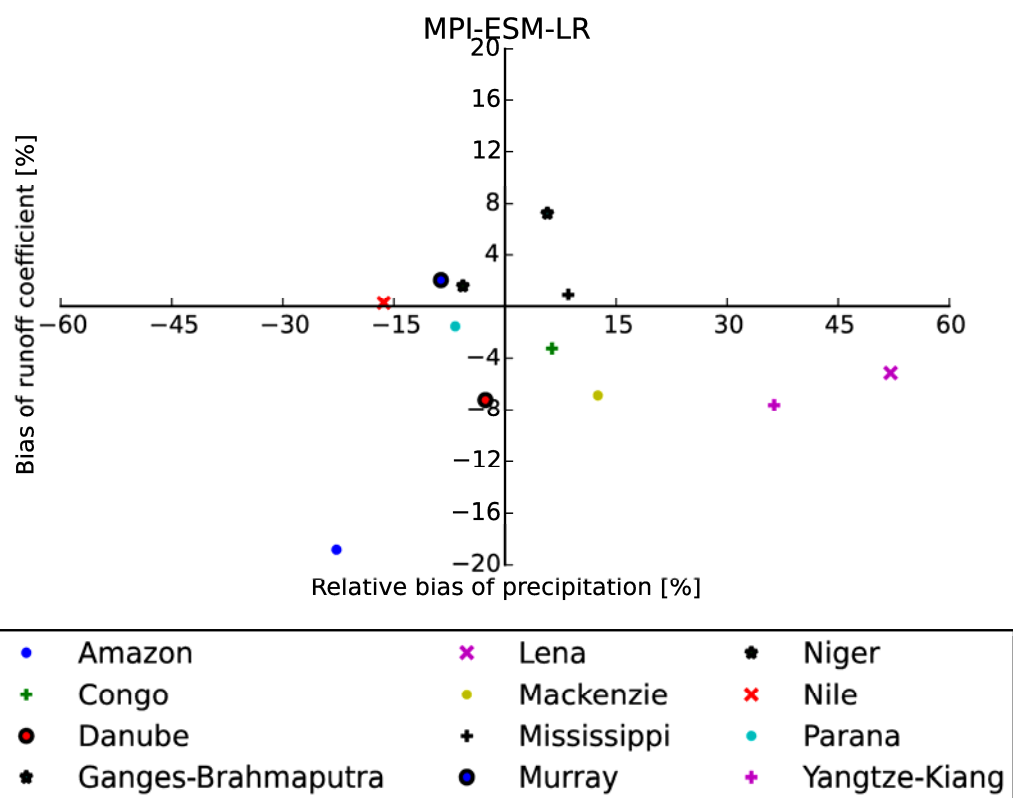


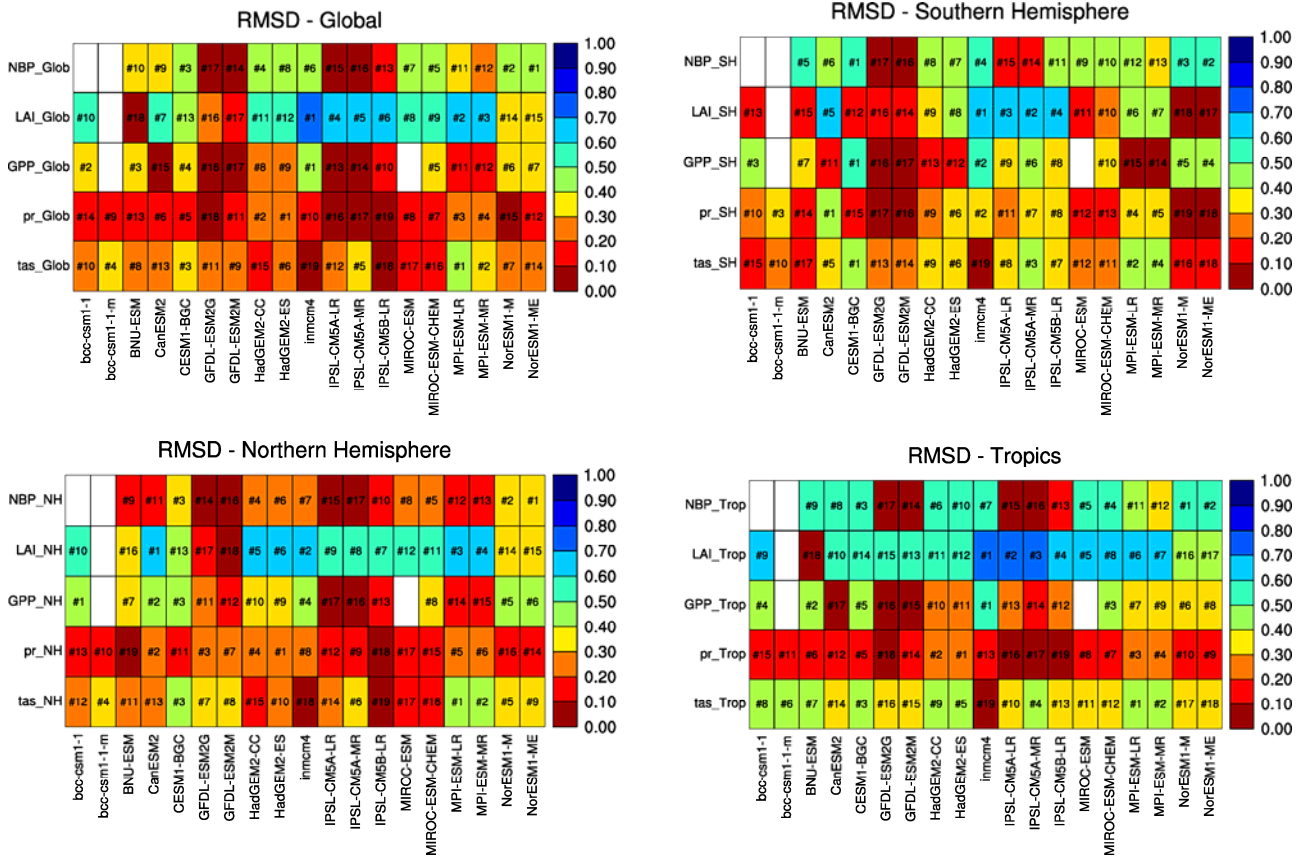
Figure 18. Bias in evapotranspiration (mm/day) for July in a subset of CMIP5 models in reference to the LandFlux-EVAL evapotranspiration product. The global mean bias is also indicated for each model as well as the RMSE. The comparison reveals the existence of biases in July evapotranspiration for a subset of CMIP5 models. All models overestimate evapotranspiration in summer, especially in Europe, Africa, China, Australia, Western North America, and parts of Amazonia. Biases of the opposite sign (underestimation in evapotranspiration) can be seen in some other regions of the world, notably over parts of the tropics. For most regions, there is a clear correlation between biases in evapotranspiration and precipitation (see precipitation bias in Fig. 4). Produced with *namelist_Evapotranspiration.xml*.



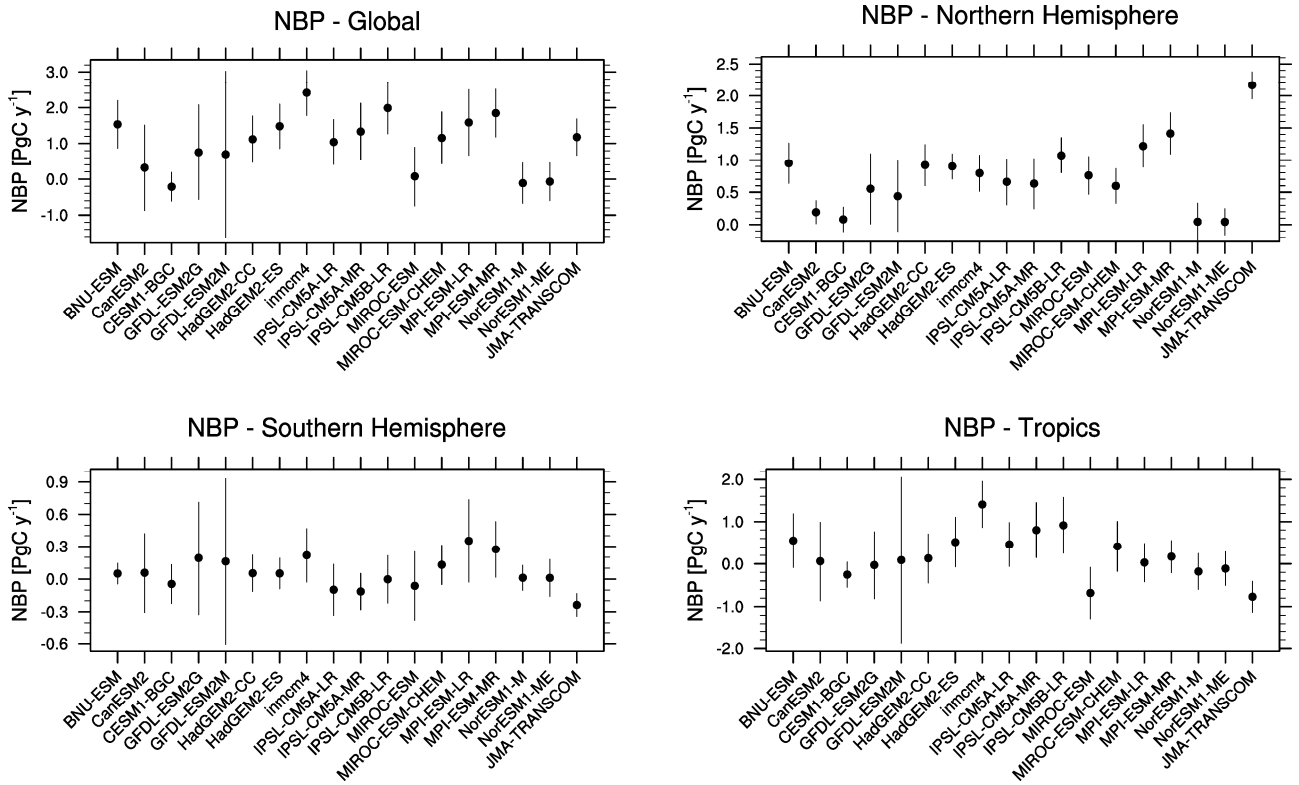
1

2 Figure 19. Biases in runoff coefficient (runoff/precipitation) and precipitation for major catchments
 3 of the globe. The MPI-ESM-LR historical simulation is used as an example. Even though positive
 4 and negative precipitation biases exist for MPI-ESM-LR in the various catchment areas, the bias in
 5 the runoff coefficient is usually negative. This implies that the fraction of evapotranspiration
 6 generally tends to be overestimated by the model independently of whether precipitation has a
 7 positive or negative bias. Produced with *namelist_runoff_et.xml*.

8



1

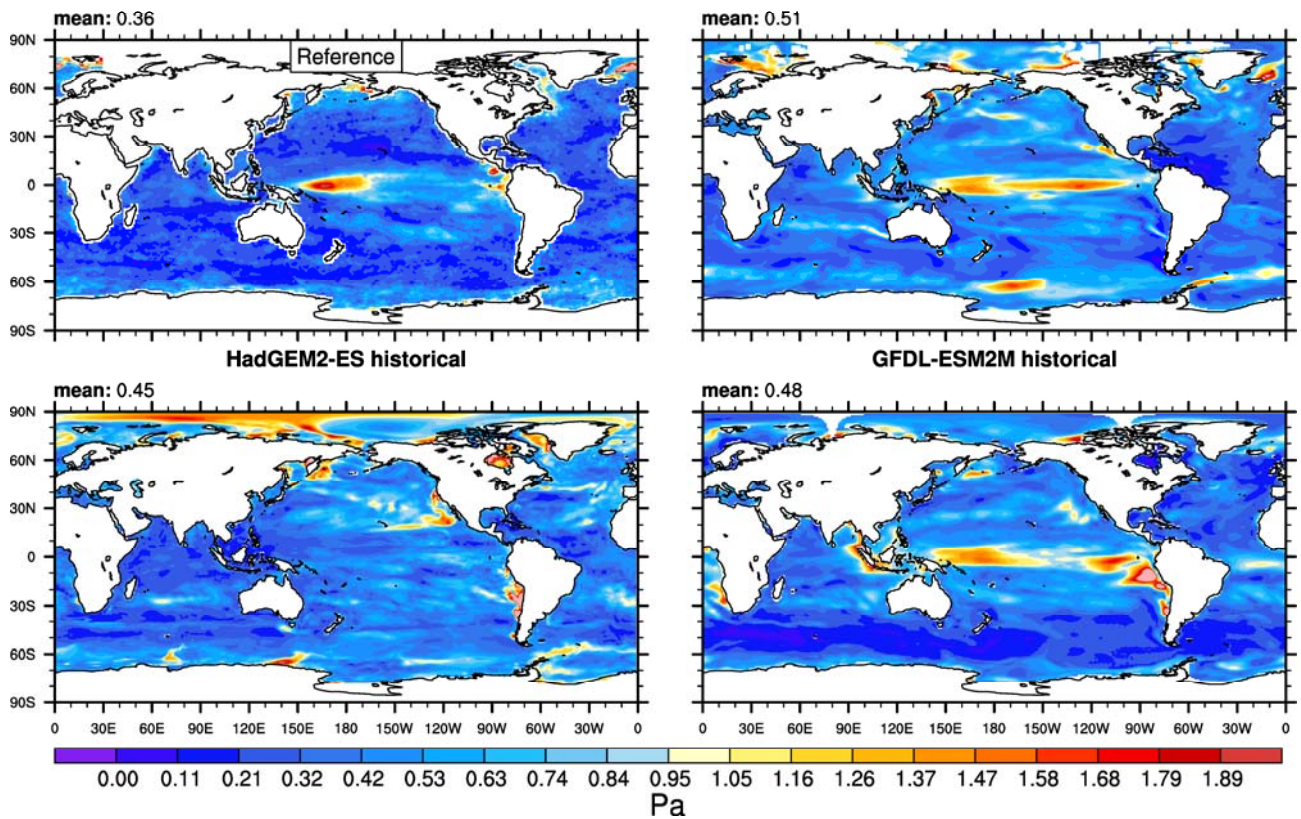


2

3 Figure 21. Error-bar plot showing the 1986-2005 CMIP5 integrated NBP for different land
 4 subdomains. Positive values of NBP correspond to land uptake, vertical bars are computed
 5 considering the interannual variation. The models are compared to JMA inversion estimates. The
 6 models' range is very large and results show that ESMs fail to accurately reproduce the global net
 7 land CO_2 flux. At the hemispheric scale, there is no clear bias common in most ESMs, except in the
 8 tropics where models simulate a lower CO_2 source than that estimated by the inversion.
 9 Reproducing Figure 6 of Anav et al. (2013) and produced with *namelist_anav13jclim.xml*.

10

JFMAMJJASOND-mean of stddev of Surface ocean $p\text{CO}_2$



1

2 Figure 22. Inter-annual variability in de-trended annual mean surface $p\text{CO}_2$ (Pa) for the period
 3 1998–2011 from an observation-based reference product (ETH-SOM-FFN; upper left) and three
 4 CMIP5 models (1992-2005). The spatial structure of inter-annual variability differs between
 5 individual CMIP5 ESMs, however both BNU-ESM and GFDL-ESM2M are able to reproduce
 6 pronounced variability in surface ocean $p\text{CO}_2$ within the Equatorial Pacific, primarily associated
 7 with ENSO variability (Rodenbeck et al., 2014). Produced with *namelist_GlobalOcean.xml*.

8

Ambient Aerosol Optical Thickness at 550 nm

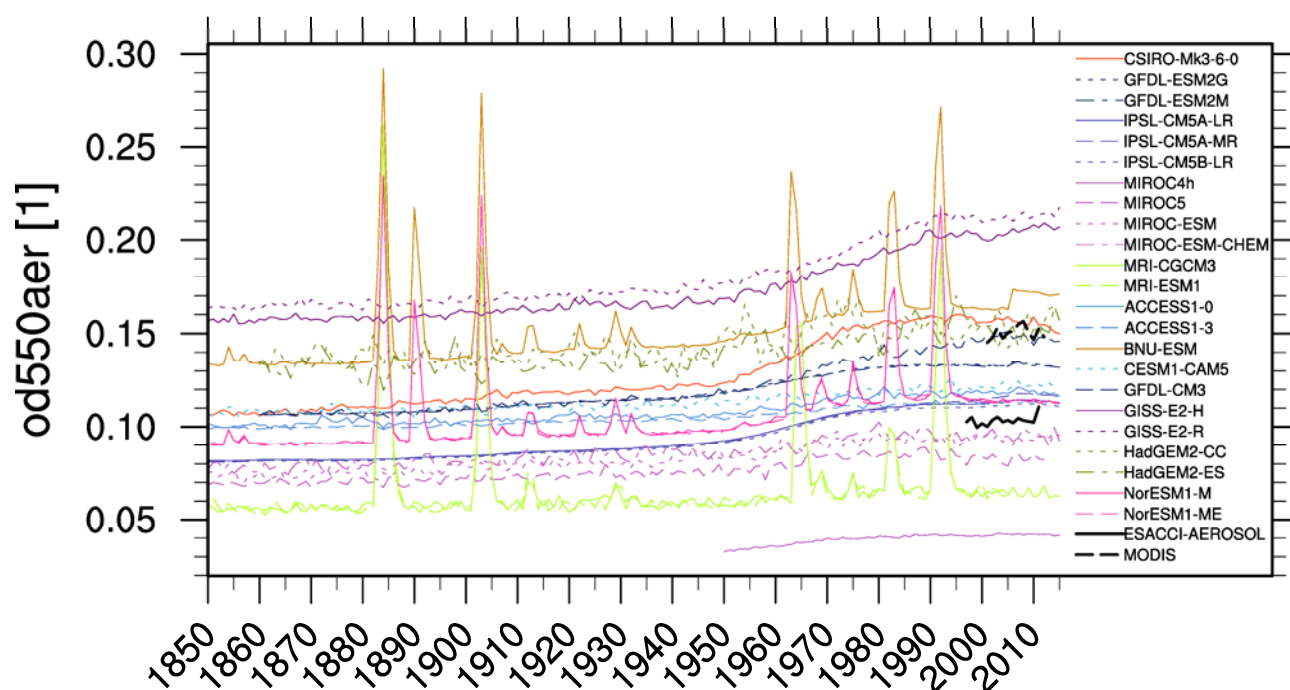


Figure 23. Timeseries of global oceanic mean aerosol optical depth (AOD) from individual CMIP5 models' historical (1850–2005) and RCP 4.5 (2006–2010) simulations, compared with MODIS and ESACCI-AEROSOL satellite data. All models simulate a positive trend in AOD starting around 1950. Some models also show distinct AOD peaks in response to major volcanic eruptions, e.g. El Chichon (1982) and Pinatubo (1991). The models simulate quite a wide range of AODs, between 0.05 and 0.20 in 2010, which largely deviates from the observed values from MODIS and ESACCI-AEROSOL. A significant difference, however, exists also between the two satellite data sets (about 0.05), indicating an observational uncertainty. Similar to Figure 9.29 of Flato et al. (2013) and produced with *namelist_aerosol_CMIP5.xml*.

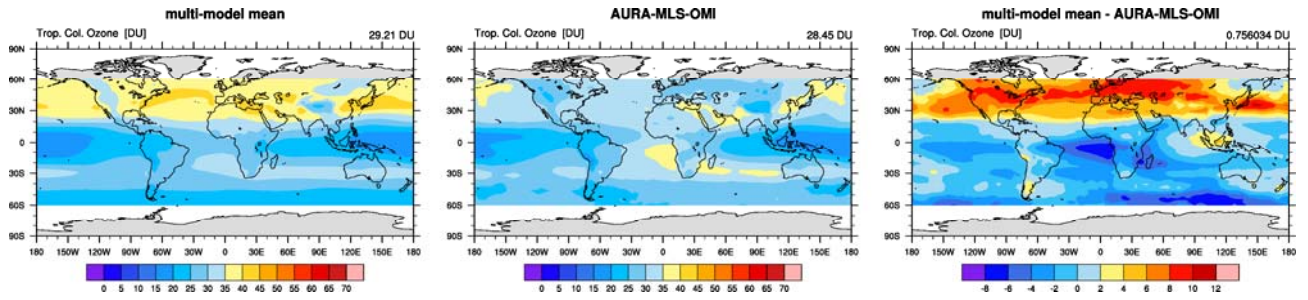


Figure 24. Climatological mean annual mean tropospheric column ozone averaged between 2000 and 2005 from the CMIP5 historical simulations compared to MLS/OMI observations (2005-2012). The values on top of each panel show the global (area-weighted) average, calculated after regridding the data to the horizontal grid of the model and ignoring the grid cells without available observational data. The comparison shows a high bias in tropospheric column ozone in the Northern Hemisphere and a low bias in the Southern Hemisphere in the CMIP5 multi-model mean. Similar to Figure 13 of Righi et al. (2015) and produced with *namelist_righi15gmd_tropo3.xml*.

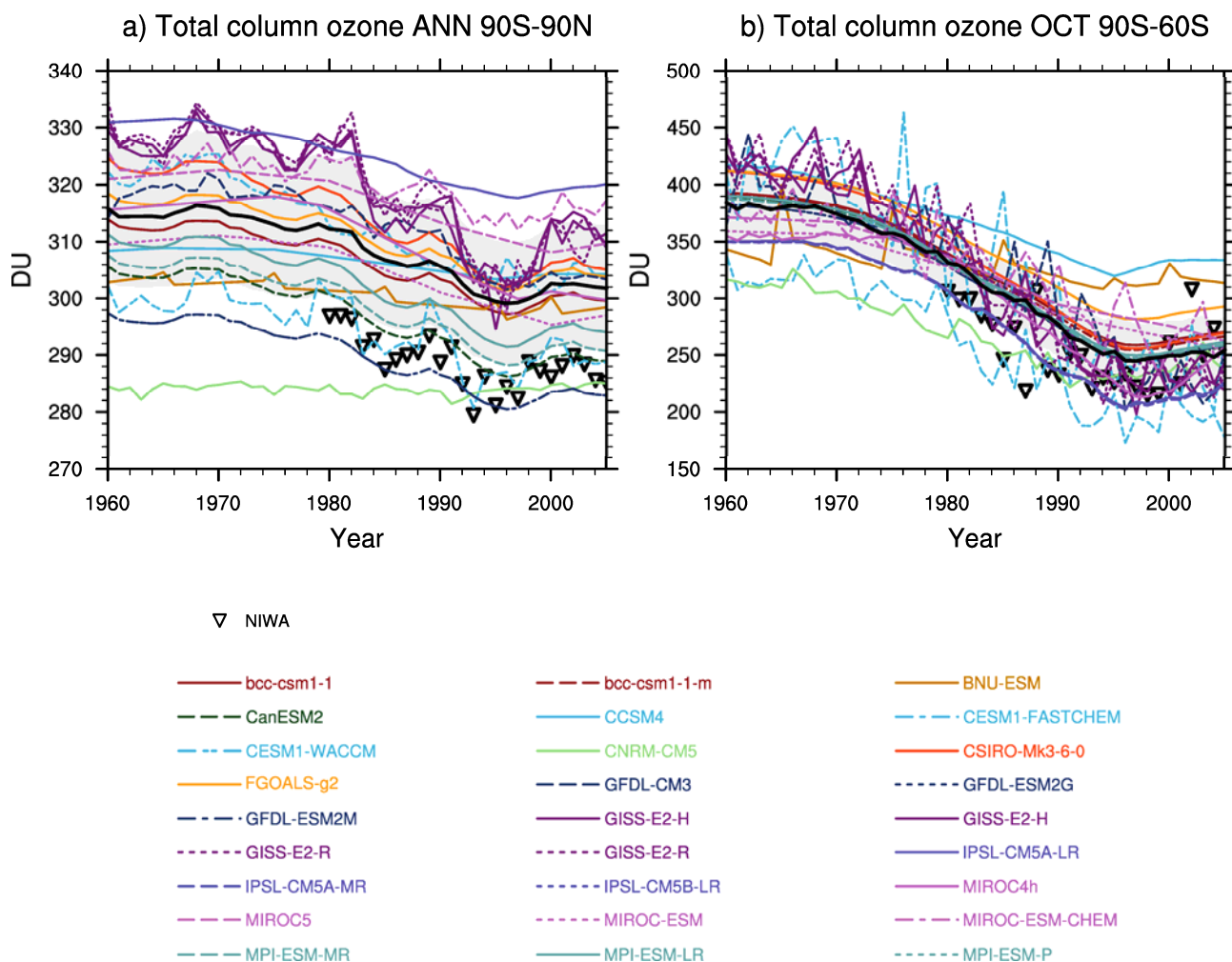


Figure 25. Total column ozone time series for (a) annual global and (b) Antarctic October mean. CMIP5 models are shown in coloured lines and the multi-model mean in thick black, their standard deviation as grey shaded area, and observations from NIWA (black triangles). The CMIP5 multi-model mean is in good agreement with observations, but significant deviations exist for individual models with interactive chemistry. Based on Figure 2 of Eyring et al. (2013) and reproducing Figure 9.10 of Flato et al. (2013), with *namelist_eyring13jgr.xml*.

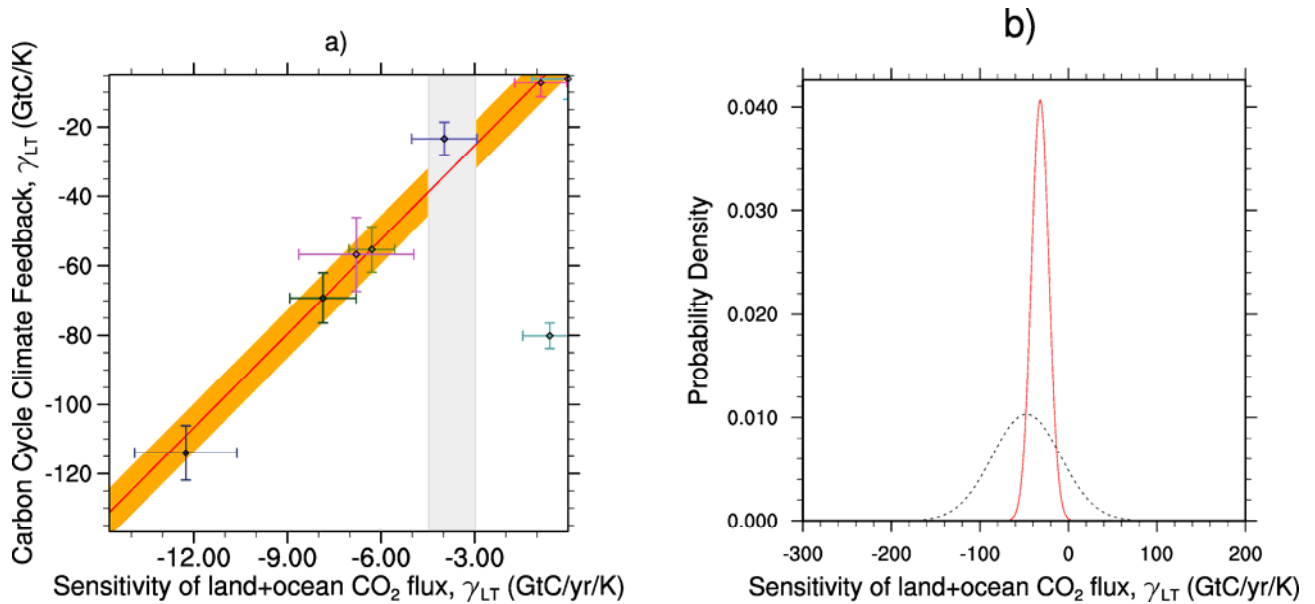
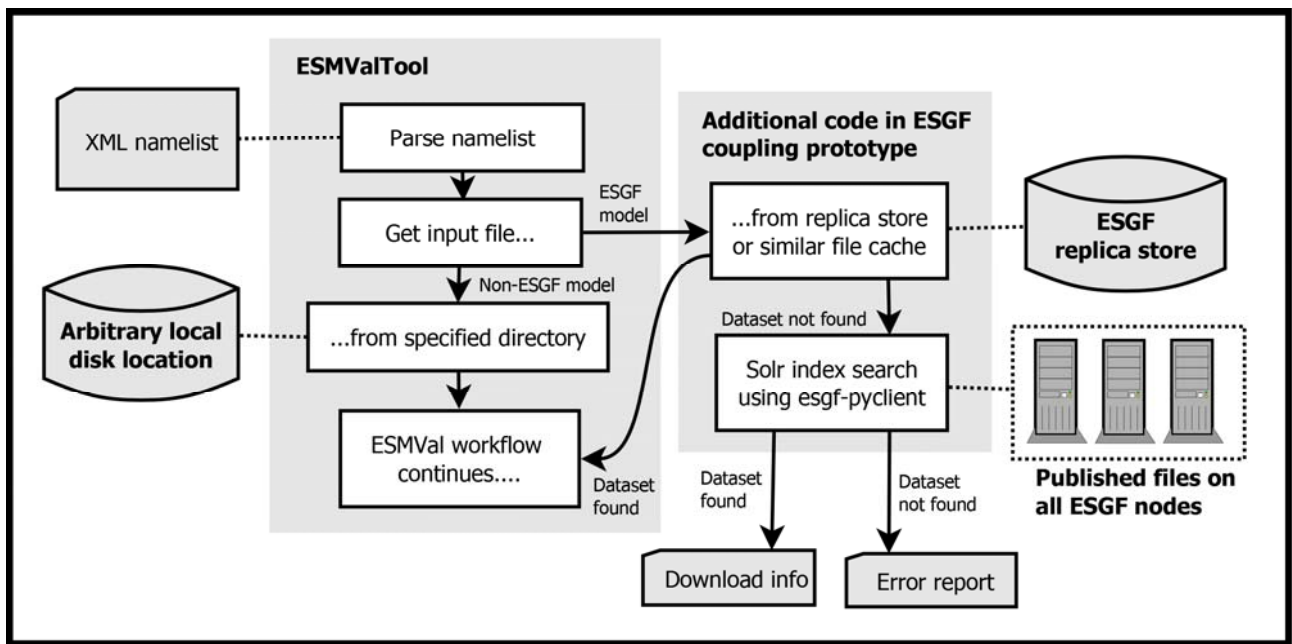


Figure 26. (a) The carbon cycle-climate feedback (γ_{LT}) versus the short-term sensitivity of atmospheric CO_2 to interannual temperature variability (γ_{IAV}) in the tropics for CMIP5 models. The red line shows the best fit line across the CMIP5 simulations and the vertical dashed lines show the observed range of γ_{IAV} . (b) probability distribution function (PDF) for γ_{LT} . The solid line is derived after applying the interannual variability (IAV) constraint to the models while the dashed line is the prior PDF derived purely from the models before applying the IAV constraint. The results show a tight correlation between γ_{LT} and γ_{IAV} that enables the projections to be constrained with observations. The conditional PDF sharpens the range of γ_{LT} to -44 ± 14 GtC/K compared to the unconditional PDF which is $(-49 \pm 40$ GtC/K). Similar to Figure 9.45 of Flato et al. (2013) and reproducing the CMIP5 model results from Figure 5 of (Wenzel et al. (2014)) with *namelist_wenzel14jgr.xml*.



1

2 Figure 27. Schematic overview of the coupling of the ESMValTool to the ESGF.

3

4