

This discussion paper is/has been under review for the journal Geoscientific Model Development (GMD). Please refer to the corresponding final paper in GMD if available.

# Quality assessment concept of the World Data Center for Climate and its application to CMIP5 data

M. Stockhause, H. Höck, F. Toussaint, and M. Lautenschlager

German Climate Computing Center (DKRZ), World Data Center for Climate (WDCC),  
20146 Hamburg, Germany

Received: 22 March 2012 – Accepted: 23 March 2012 – Published: 13 April 2012

Correspondence to: M. Stockhause (stockhause@dkrz.de)

Published by Copernicus Publications on behalf of the European Geosciences Union.

781

## Abstract

The preservation of data in a high state of quality and suitable for interdisciplinary use is one of the most pressing and challenging current issues in long-term archiving. For high volume data such as climate model data, the data and data replica are no longer stored centrally but distributed over several local data repositories, e.g. the data of the Climate Model Intercomparison Project No. 5 (CMIP5). The most important part of the data is to be published as DOI according to the World Data Center for Climate's (WDCC) application of the DataCite regulations. The integrated part of WDCC's data publication process, the data quality assessment, was adapted to the requirements of a federated data infrastructure. A concept of a distributed and federated quality assessment procedure was developed, in which the work load and responsibility for quality control is shared between the three primary CMIP5 data centers: Program for Climate Model Diagnosis and Intercomparison (PCMDI), British Atmospheric Data Centre (BADC), and WDCC. This distributed quality control concept, its pilot implementation for CMIP5, and first experiences are presented.

## 1 Introduction

The International Panel on Climate Change (IPCC) aims to establish one common climate model data archive to advance the knowledge of climate change and variability. The results collected within the Climate Model Intercomparison Project No. 5 (CMIP5) are intended to underlie the coming fifth assessment report (IPCC-AR5). CMIP3 data for the last report IPCC-AR4 were collected and provided centrally by the Program for Climate Model Diagnosis and Intercomparison (PCMDI) without version control and with compact unformalized metadata information, which was imprecise in respect of model and simulation descriptions. The data volume for CMIP5 is expected to reach nearly 100 times that of CMIP3 (Taylor et al., 2012).

782

These experiences from CMIP3 together with the expected data volume led to three main improvements for the CMIP5 data infrastructure:

- Data is stored in several decentralized data nodes connected by the Earth System Grid (ESG; Williams et al., 2009, 2011). Three of them located at major data centers have built a federated system of data archives (also called primary CMIP5 data portals, Taylor et al., 2012): PCMDI, the British Atmospheric Data Centre (BADC), and the World Data Center for Climate (WDCC). These centers committed to hold replica of the most important part of the CMIP5 data, i.e. IPCC relevant data, on hard disks for quick access and data security.
- Information on models and simulations is enlarged significantly. The metadata schema used is the Common Information Model (CIM) developed by METAFOR and collected via a web-based questionnaire (Guilyardi et al., 2011).
- Data curation was improved by introducing a versioning concept and a quality assessment process providing a uniform identification of datasets as well as a persistent identifier DOI (Digital Object Identifier) for data citation in scientific publications. The data DOI, like a DOI for printed papers, gives scientific credits to data creators for their work and allows for persistent and direct data access.

The quality assessment procedure for CMIP5 has to support the federated data infrastructure and incorporate all available metadata resources, especially CIM metadata and those stored in the self-describing data headers of the netCDF files. A general concept for a distributed and coordinated quality assessment procedure suitable to use in a distributed data infrastructure was developed (Sect. 2). This concept was altered and adapted for its pilot application within CMIP5 (Sect. 3).

## 2 Concept of a distributed quality assessment of high volume data

Quality control and description of data in repositories and especially in long-term archives are generally viewed as essential. Moreover, a demand for more efficient evaluation services to convert data into information and information into knowledge is detected (Overpeck et al., 2011). This is of special importance for open-access data of interdisciplinary use, where no direct contacts between data users and original data creators exist any longer. However, contents of the quality checks as well as definitions of quality levels and the overall quality procedure vary significantly between data types and scientific disciplines.

The ESIP (Federation of Earth Science Information Partners), a consortium of 120 organizations, formulated some principles on data stewardship and recommended practices (ESIP, 2011): Quality assessment and its documentation are tasks of the data creator. Data intermediaries like repositories should set time limits for quality control procedures in order to prevent it from delaying data accessibility. Data intermediaries additionally function as communicators between data creators and data users. ESIP (2011) focuses on the scientific content of the data in its principles for quality assessment. For scientific data distributed over several repositories this scientific quality assurance (SQA) has to be accompanied by a technical quality assurance (TQA). The TQA checks data and metadata consistency among the distributed data and metadata repositories and might include a check against data and metadata standards. This TQA can only be applied by the data intermediaries at the data repositories adding the TQA task including its documentation to their communicator role.

Quality control procedures of high volume data have to be carried out at the storage location before opening the repository for interdisciplinary data access and use. Together with the trend towards decentralized data repositories, quality control procedures have to become distributed/federated themselves and need to be coordinated and standardized (Sect. 2.2).



- the *analysis* of QC results by exception statistics, provenance information, and plotting,
- the *insert* of QC tool application information and results into the Project Metadata Repository,
- 5 – the *assignment* of QC levels (including a possibility to exclude certain data from the assignment to a data collection), and
- the final data checks for DOI data publication by *providing information* on project metadata.

### 2.3 Embedding the distributed quality control into a federated data infrastructure

Quality Control procedures rely on data and metadata accessibility. Data is stored in different local data repositories or Data Nodes (DN). Metadata (MD) is created during the whole project life time, starting with the description of model and model application (MD on model/simulation) provided by the data creator and inserted via a Metadata GUI (Fig. 2). Later metadata on the data in the data nodes and metadata on quality 15 of the different QC level checks are added. These metadata is collected and stored centrally in a project metadata repository.

The available metadata information in the project metadata repository is used within the cross- and double checks (TQA) of the DOI publication process (QC Level 3). At the end of the project and the QC procedure, data and metadata are long-term archived, 20 i.e. a data copy is stored at a data center along with all available metadata out of the project metadata repository. DOI published data is long-term archived at the long-term archive (LTA) of the DOI publication agency. The DOI is assigned via the registration agency DataCite to the IDF and is integrated into the global handle system. The DOI 25 resolves to an entry metadata page hosted by the DOI publication agency.

787

Data catalogs support data discovery by harvesting the metadata information of the project metadata repository or the LTA, and the DOI handle system supports data discovery of the long-term archived data after the end of the project.

## 3 Application of the distributed quality control in CMIP5

- 5 The distributed quality assessment procedure was adapted for the pilot implementation in the CMIP5 project. Detailed information on the quality control procedure within CMIP5 is available at <http://cmip5qc.wdc-climate.de>.

### 3.1 Quality control procedure within CMIP5

The definition of quality control levels and its implications are summarized in Table 1. 10 The quality procedure workflow with its actors is sketched in Fig. 3. The data collection for QC level assignment within CMIP5 is a CMIP5 simulation, i.e. all data of all realizations of one experiment carried out with a certain model. CMIP5 data is ESG published at decentralized local data centers. Most of the modeling centers decided to either host their own data node or ESG publishes their data at a national data node. 15 The data submission step is performed by the integration of the data node's metadata in the ESG gateway catalogs. During ESG publication at the data nodes, QC level 1 checks are performed, i.e. CMOR2 and ESG publisher conformance. Access of data of level 1 is restricted to selected users, who contribute to the QC process by reporting problems and errors to the data nodes or the ESG gateways.

20 In a second submission step the data is copied from the local data centers to one of the primary CMIP5 data centers (PCMDI, BADC, or DKRZ) for quality checks of level 2, followed by an ESG data publication at the CMIP5 data center. The QC Manager at the CMIP5 Center can alternatively decide to carry out the QC L2 checks at the local data center in parallel or prior to data replication. An example for a QC L2 check criteria is 25 the continuity of the time axis and the usage of the accurate CF (Climate and Forecast)

788



All QC L2 results stored in the QC repository become immediately accessible to the community via the QC result service. CIM quality documents are created and published via atom feed at QC level 2 and QC level 3 assignments. The data citation service provides a central entry to search by tracking\_id for the current preliminary citation recommendation in case of data of QC level 1 or 2 and for the persistent citation of DOI published data of QC level 3.

Thus, the QC repository serves within CMIP5 not only as intermediate QC result storage facility to support the QC process but additionally as long-term source of quality results and quality-related information for the climate community.

WDCC plays a double role in the federated quality procedure of CMIP5 by performing QC L2 on their share of the CMIP5 data and functioning as DOI publication agency for all replicated long-term archived data (QC L3).

### 3.3 Experiences of the CMIP5 quality approach

First experiences of the federated quality assessment procedure in CMIP5 are encouraging. The federated QC approach is capable to serve as QC procedure for CMIP5. The QC has helped to find data inconsistencies in order to improve the CMIP5 data quality. First DOIs on data are assigned, e.g. doi:10.1594/WDCC/CMIP5.MXELAM. Though the QC concept and its implementation worked out fine, several problems occurred during the QC application:

1. Most modeling centers are still in the process of ESG publishing new data versions or additional data, e.g. cfMIP data. Since no deadline exists for the creation of data for the DOI data publication process, the QC L2 process for the CMIP5 simulations cannot be finished but has to be continued several times. Among other reasons QC L2 findings contribute to these data revisions. The contact between QC L2 manager and data creator is more intensive than expected. Several findings during QC L2 require the interpretation of the data creator to distinguish a real error from a minor model-specific issue.

791

2. The time necessary for data replication was underestimated. Reasons are narrow bandwidths together with the ongoing data changes. As the long-term data archiving at the DOI publication agency is a precondition for DOI assignment, the DOI data publication is significantly delayed. The former aim to provide DOI data citation references for scientists of IPCC Working Group I (WG I) was altered in order to provide those citations for WGs II and III in fall 2012. Additionally, the data aggregation for a data DOI had to be changed from including all data of a CMIP5 simulation into including at least the monthly and yearly data of a CMIP5 simulation. The data of higher temporal frequency will be published as a new DOI, related to the first DOI via DataCite relation *isSupplementTo* (DataCite, 2011). The DOI data publication decision is a compromise between data completeness and providing data citation regulations for scientists, especially for those contributing to the 5th IPCC assessment report.

The data replication problem had three implications for the QC procedure: First, the QC L2 checks had to be distributed even more to enable QC L2 applications directly at the data nodes. The QC manager at PCMDI, BADC or DKRZ remains responsible for the QC L2 assessment and thus the QC level 2 assignment. However, he could either delegate the QC L2 checker tool run at the data node to the data node manager or run it himself. Secondly, the QC procedure for level 3 had to be altered to enable the exclusion of non-replicated data before starting the QC L3 process. And, thirdly, the scientist of WG I had to be supported in the citing of CMIP5 data, esp. of data of QC levels 1 and 2, i.e. data without DOIs. WDCC set up its citation service for that purpose.

3. In the data infrastructure data, replication as well as the inclusion of data replica and multiple data versions are not fully supported by the current ESG gateways. Thus, QC status or DOI data are not integrated in their data discovery functionality, but remain separated pieces of information. Even the available CIM quality documents are not harvested by the ESG gateways, so far, due to other more urgent development issues. To bridge this intermediate invisibility of QC status and

792

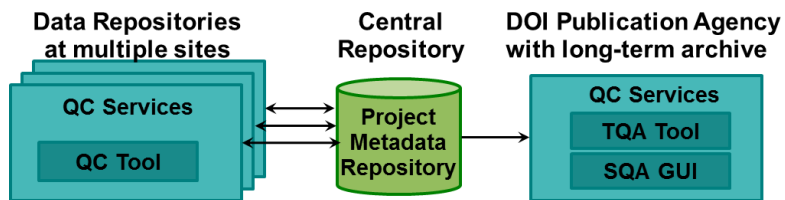




- Klump, J., Bertelmann, R., Brase, J., Diepenbroek, M., Grobe, H., Höck, H., Lautenschlager, M., Schindler, U., Sens, I., and Wächter, J.: Data Publication in the Open Access Initiative, *Data Science Journal*, 5, 15 June 2006. 785
- Overpeck, J. T., Meehl, G. A., Bony, S., and Easterling, D. R.: Climate Data Challenges in the 21st Century, *Science*, 331, 700, doi:10.1126/science.1197869, 11 February 2011. 784
- 5 Taylor, K. E., Stouffer, R. J., and Meehl, G. A.: An Overview of CMIP5 and the Experiment Design, submitted, *B. Am. Meteorol. Soc.*, doi:10.1175/BAMS-D-11-00094.1, 2012. 782, 783
- Williams, D.N., Ananthkrishnan, R., Bernholdt, D. E., Bharathi, S., Brown, D., Chen, M., Chervenak, A. L., Cinquini, L., Drach, R., Foster, I. T., Fox, P., Fraser, D., Garcia, J., Hankin, S., Jones, P., Middleton, D. E., Schwidder, J., Schweitzer, R., Schuler, R., Shoshani, A., Siebenlist, F., Sim, A., Strand, W. G., Su, M., and Wilhelmi, N.: The Earth System Grid: Enabling Access to Multimodel Climate Simulation Data, *B. Am. Meteor. Soc.*, 90, 195–205, doi:10.1175/2008BAMS2459.1, 2009. 783, 789
- 10 Williams, D. N., Taylor, K. E., Cinquini, L., Evans, B., Kawamiya, M., Lautenschlager, M., Lawrence, B. N., Middleton, D. E., and ESGF contributors: The Earth System Grid Federation: Software Supporting CMIP5 Data Analysis and Dissemination, *CLIVAR Exchanges*, 56, 16, 40 ff., 2 May 2011. 783, 789
- 15

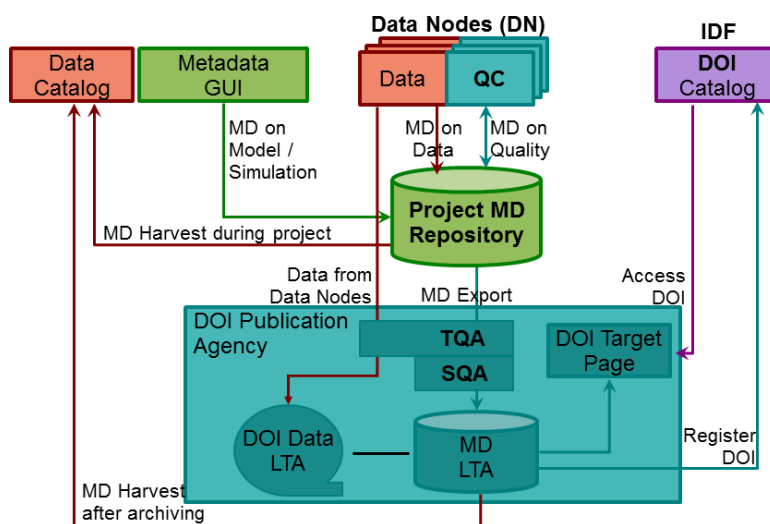
**Table 1.** CMIP5/IPCC-AR5 quality control levels and their implications.

	<b>QC Level 1:</b> CMOR2, ESG Conformance of Data and CIM Conformance of Metadata	<b>QC Level 2:</b> WDCC Conformance and subjective controls	<b>QC Level 3:</b> DOI Data Publication via DataCite
<b>Data</b>	Data preliminary; no user notification about changes; performed for all data; metadata may not be complete	Not finally agreed; no user notification about changes; performed for replicated data	published and persistent data with version and unique DOI as persistent identifier; performed for replicated data
<b>Access</b>	constrained to CMIP5 modeling centers	constrained to non-commercial research and educational purposes	constrained to non-commercial research and educational purposes or open for unrestricted use
<b>Citation</b>	no citation reference	informal citation reference	formal citation reference
<b>Quality Flag</b>	<i>automated conformance checks passed</i>	<i>subjective quality control passed</i>	<i>approved by author (in case of newer DOI available: approved by author, but suspended)</i>



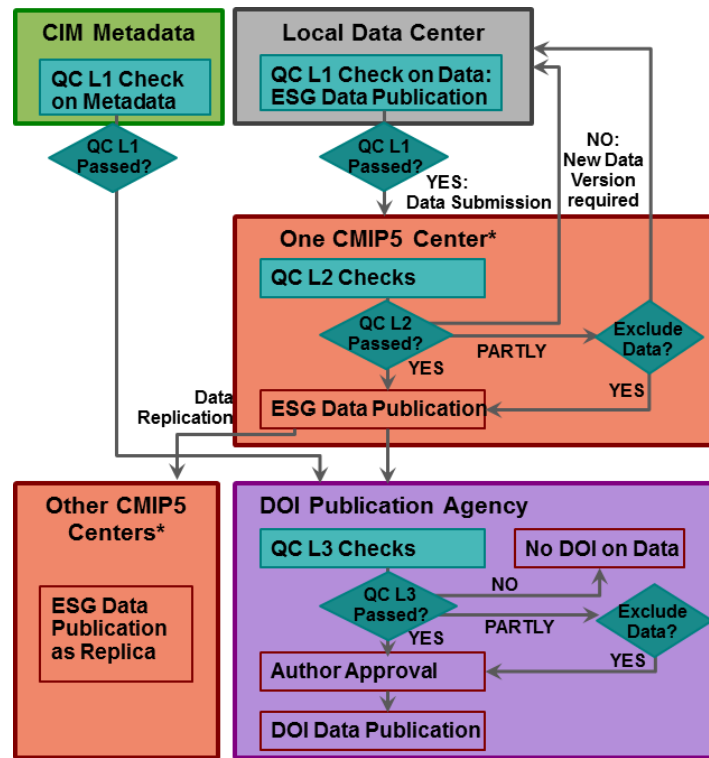
**Fig. 1.** Infrastructure components of a distributed quality assessment procedure.

799

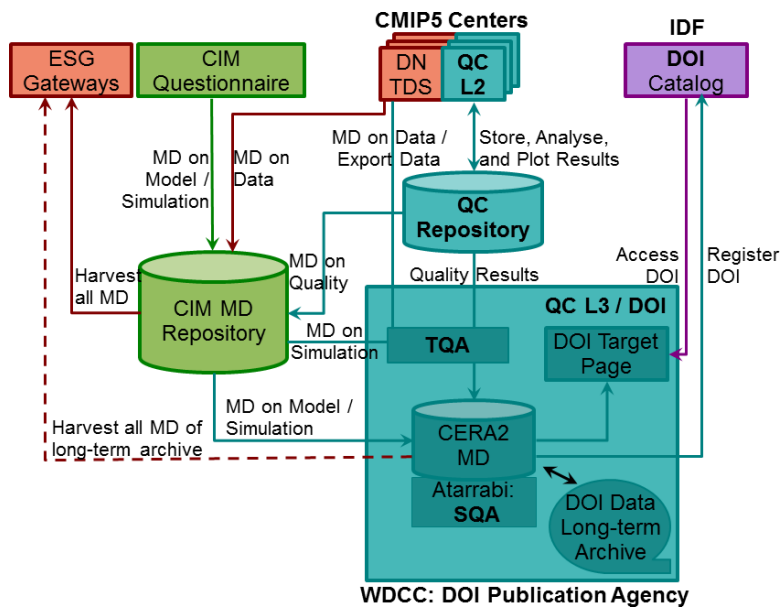


**Fig. 2.** Distributed Quality Control in a federated data infrastructure (MD: metadata, LTA: long-term archive).

800



**Fig. 3.** Workflow of the quality control procedure (\*: Current primary CMIP5 data centers are PCMDI, BADG, and DKRZ/WDCC).



**Fig. 4.** Implementation of the distributed quality control approach for CMIP5 (MD: metadata, DN: data node, TDS: THREDDS data server).