

# Quality Assessment Concept for CMIP5

*Martina Stockhause, Michael Lautenschlager*

# CMIP5 / IPCC-AR5 in Numbers

## Coupled Model Intercomparison Project (CMIP5)

- **Participants:**

- ca. 20 participating modelling centres
  - with ca. 40 model configurations (different resolutions)

- **Experiments:**

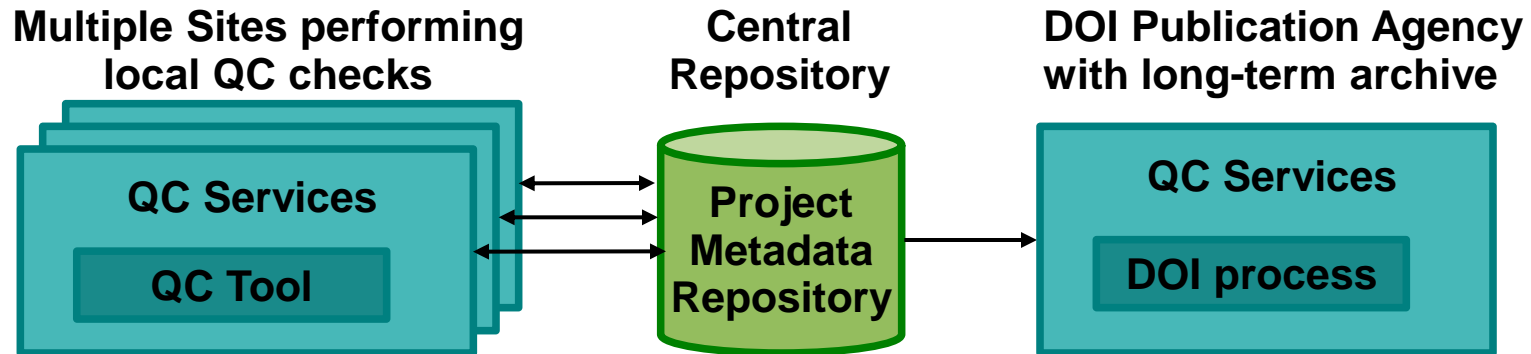
- 60 experiments with partly several realisations (ensemble members)
  - over 90 000 model years

- **ca. 2 Mio. atomic datasets of ca. 400 experiments**

- **Data Volume:**

- ca. 10 PB output: → **QC L1**
  - thereof ca. 2.5 PB requested *'output1'*, *'output2'* → **QC L2**
  - thereof ca. 1 PB replicated *'output1'*: IPCC-AR5 → **QC L3 / DOI**

# Distributed Quality Control Approach for High Data Volumes



## QC checks:

- QC Execution Service (qcWrapper.py): QC tool run and Repository ingests of configuration and results
- QC Analysis Services (qcDbselect.py) for data analyses and exception statistics
- QC Plotting Service (qcDbselect.py) and plot ingest in Repository
- QC level assignment (qcAssignL2.py)

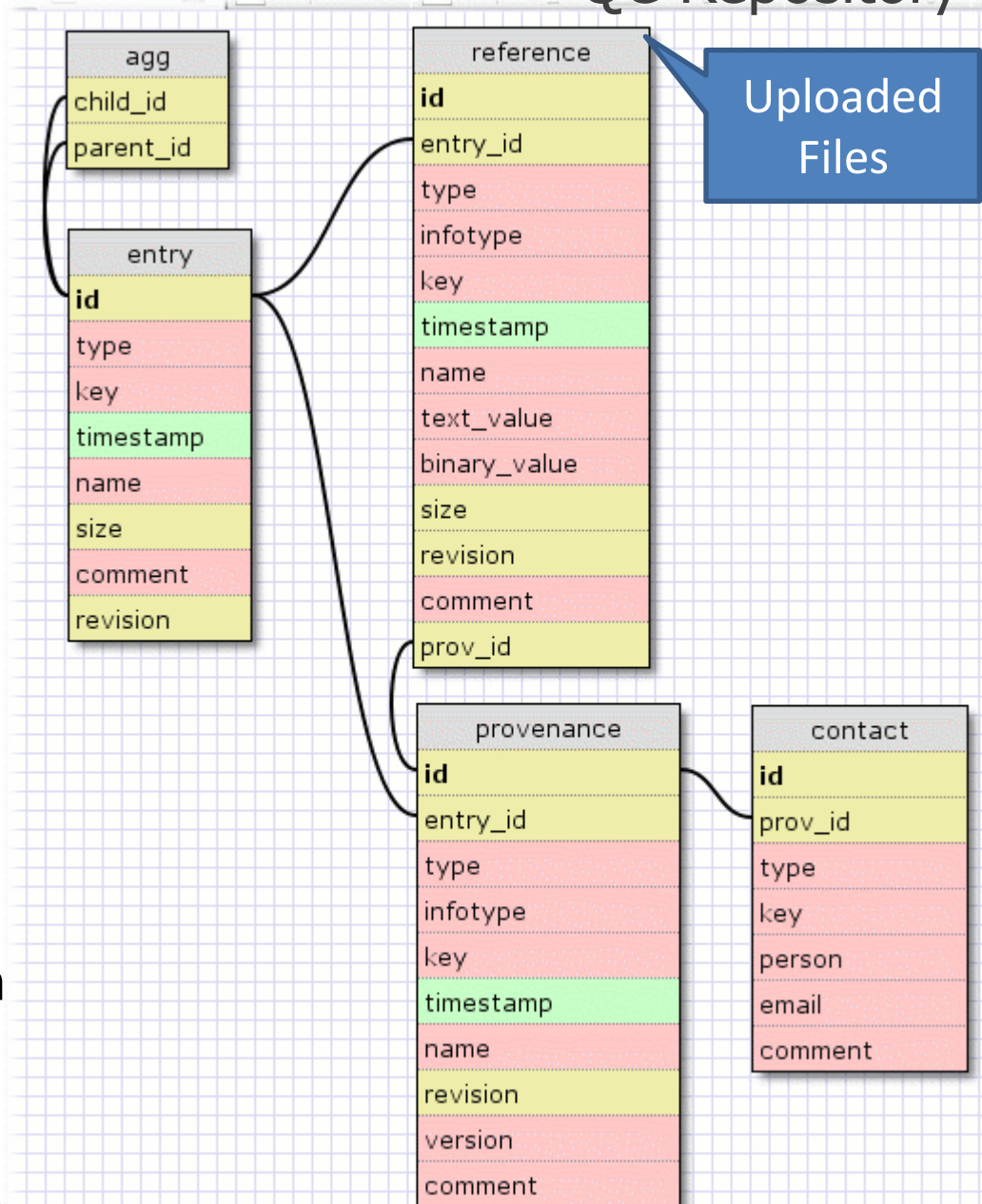
- 
- ## DOI publication:
- Export QC results (qcDbselect.py) for DOI publication process

## Aggregation Levels:

- **NC:** chunk names  
(for completeness control)
- **DRS Atomic Dataset:**  
QC results
- **DRS Experiment / DOI:**  
QC I/O files; plots; status;  
contact

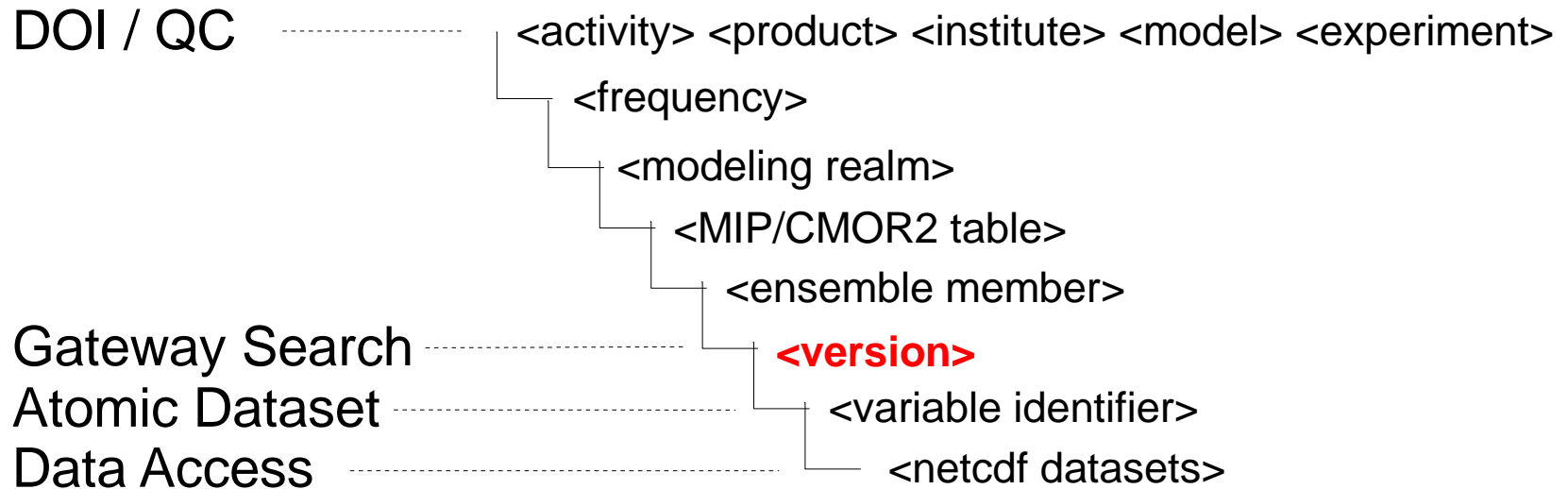
## Additional Services:

- View for version control
- Table with history information  
filled by a trigger



## CMIP5 Process

## DRS Name / Hierarchy Level





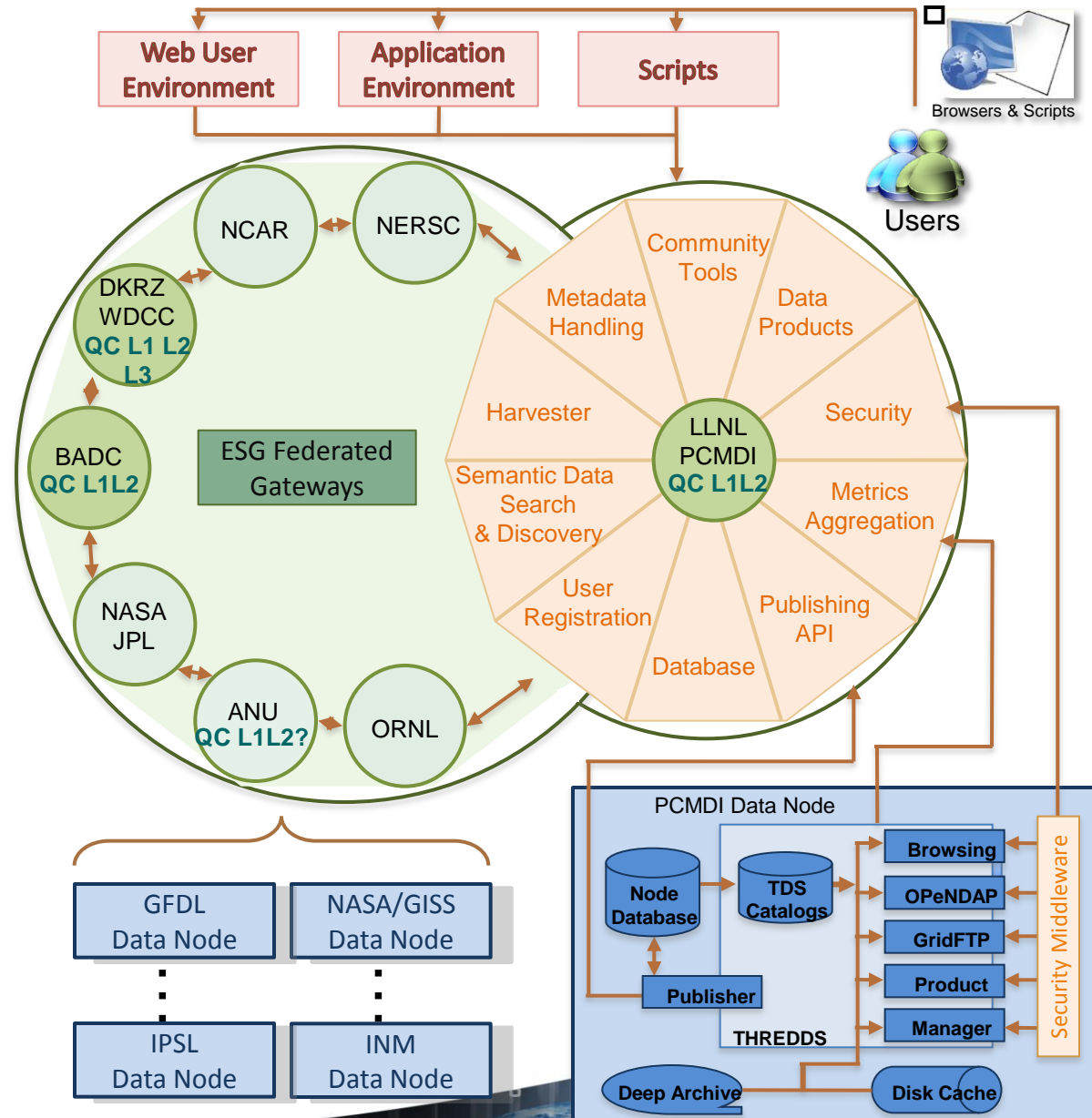
## ESG Federation

**PCMDI** (Data / Security Earth System Grid: ESG)

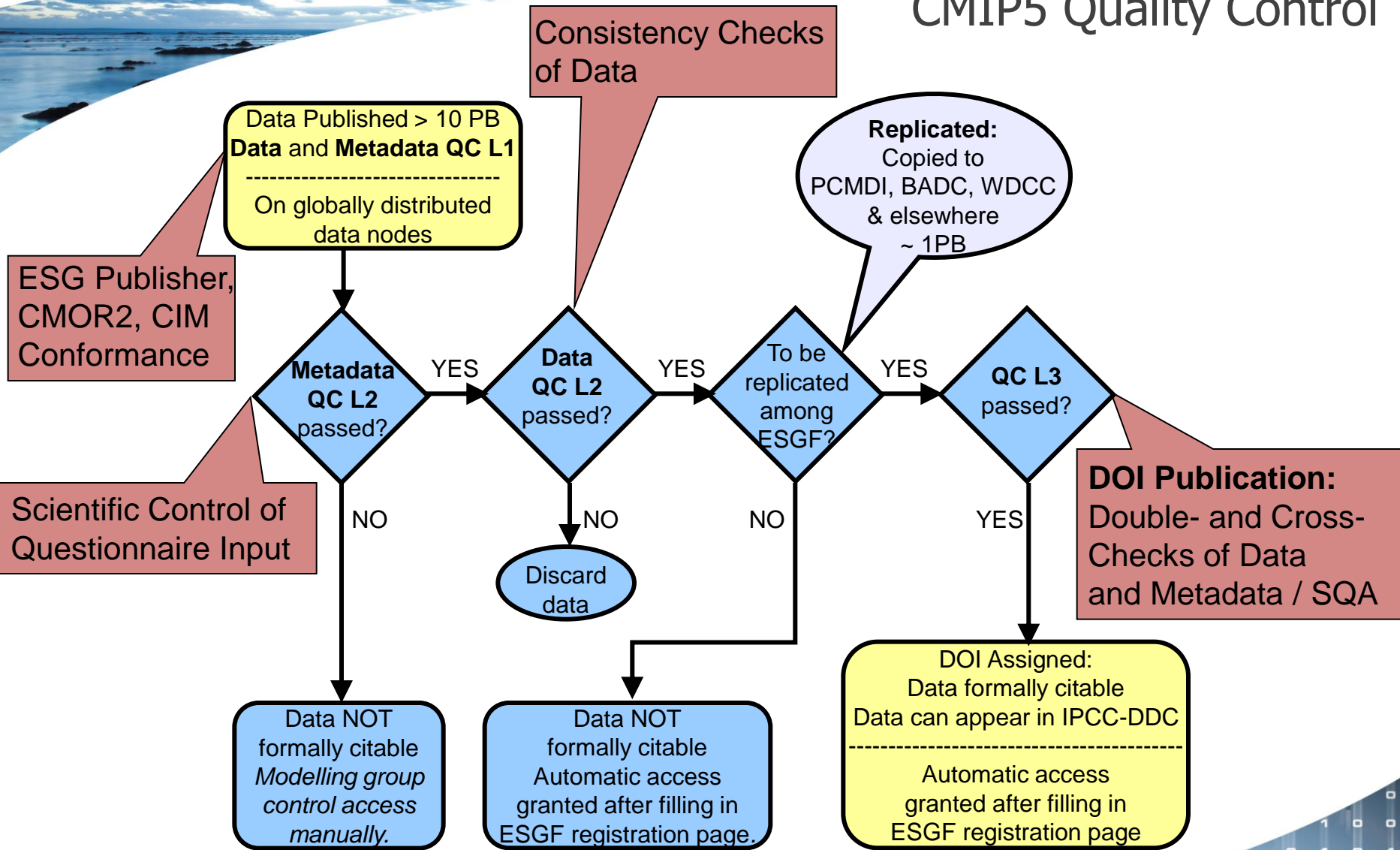
**BADC** (Metadata: CIM / Metafor)

**WDCC** (Quality Control, DOI data publication)

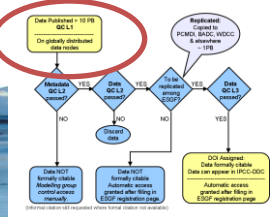
- Data Collection CMIP5
- ESG Gateway
- Data replication of CMIP5 data subset
- Quality Control (QC) on data for levels 1 and 2



# CMIP5 Quality Control



(Informal citation still requested where formal citation not available)

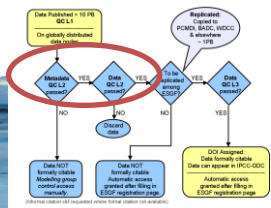


## QC Level 1 (CMOR2 and ESG publisher conformance checks): Performed at all ESGF partners during ESG publication

- **Data checks:**
  1. cmor2 compliance checks by the cmor checker `check_CMOR_compliant.py`
  2. esg publisher conformance
- **Metadata checks:** Completeness and technical validation of questionnaire input



# Quality Control Level 2

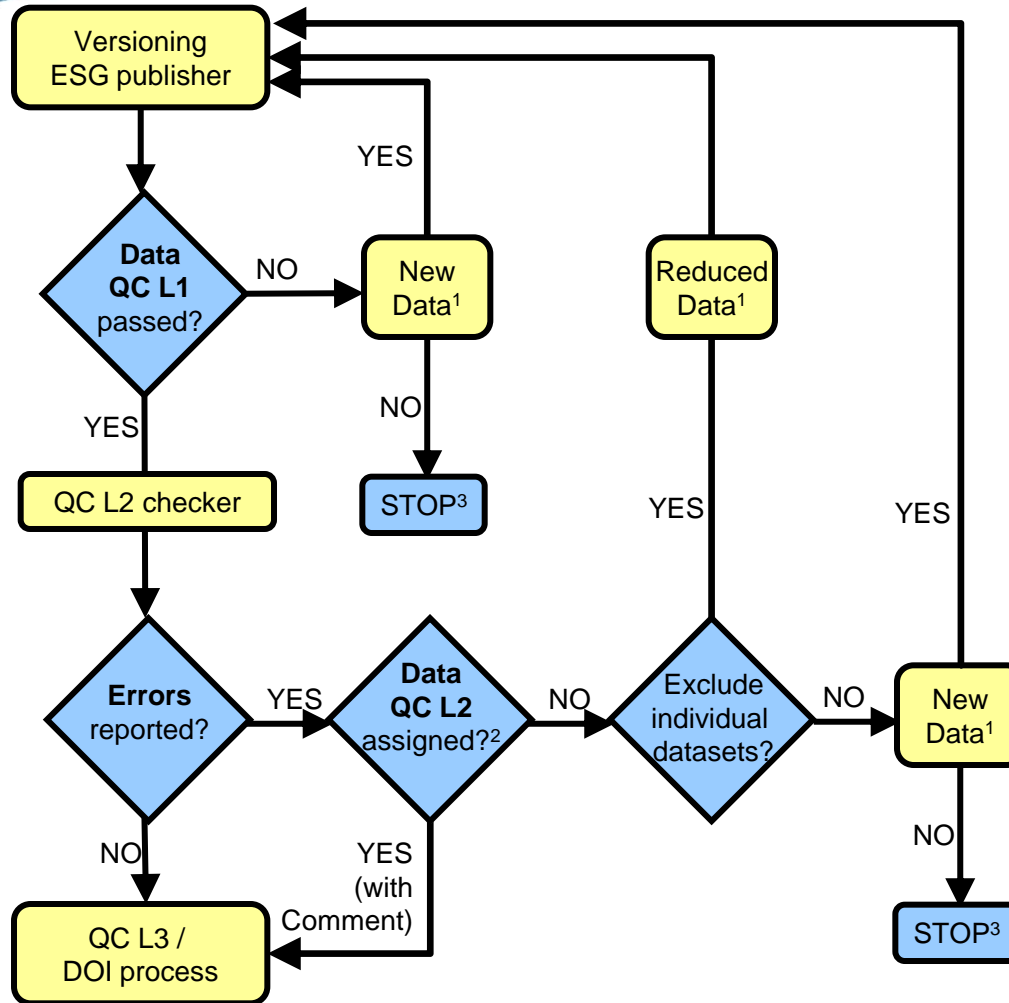
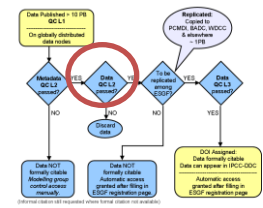


QC Level 2 (subjective quality control passed):

**Performed on requested subset of CMIP5 data at all 3 (4) sites**

- **Data checks:** Consistency checks: Check of statistical global values and additional DRS checks
- **Metadata checks:** Subjective metadata control by scientist metadata available via AtomFeed

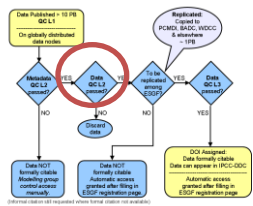
# Assignment of Quality Control Level 2



<sup>1</sup> New version assigned to ESG datasets

<sup>2</sup> Assignment in co-operation with data author according to criteria <http://www.leuchtturm-atlas.de/SCR/qc2list.html>

<sup>3</sup> Delete possible QC L2 results out of QCDB



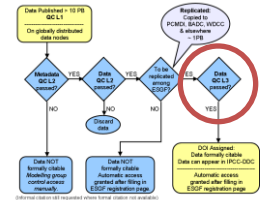
## QC2 Exception Codes (preliminary version)

Fatals: immediate action necessary,  
 errors: data unacceptable,  
 warnings: data possibly not ok,  
 informatory: fine - just for info,  
 unclear: title lines or open issues.

Thematic exception groups (see first column):  
**ACCESS** errors, **GENERAL** checks, **METADATA** and file name checks,  
**TABL**: inconsistencies in comparison to meta data tables,  
**TIME** axis checks, **VARIABLES**' checks.  
**OBSOLETE** messages (not used for CMIP5 project).

The following flags F<n> refer to general checks, mainly on the time axis.

key group	description	comment
<b>F-1</b> GENE	-- Not checked	
<b>F0</b> GENE	-- No error found	
<b>F1</b> TIME	testTimeStep() ^0, ^5 Error: negative time step	
<b>F2</b> TIME	testTimeStep() ^0, ^5 Error: missing time step	This, of course, is no error if the QC is run over several time slices with intentional gaps in between. You may want to set NON_REGULAR_TIME_STEP in the setup file to check only for positive increments (of perhaps different sizes).
<b>F4</b> TIME	testTimeStep() ^0, ^5 Error: identical time step	
<b>F8</b> TIME	testCalendarTimeBounds() Error: negative/zero time bounds range	
<b>F16</b> TIME	testCalendarTimeBounds() ^0 Error: overlapping time bounds ranges	
<b>F32</b> TIME	testCalendarTimeBounds() ^0 Warning: gap between time bounds ranges	



## QC Level 3 (approved by author):

- **Double and cross-checks of data and metadata**
- **Author check and approval of data and metadata (SQA)**
- **DOI publication of data (DataCite):**  
assignment of DOI as persistent identifier and citation direction  
resolving the DOI opens the DOI landing page.

	Permission: QC L2	Scientific Q. Assurance	Technical Q. Assurance	DOI- Publication
Scientist				
Publication Agency				
Registration Agency				

TIME →

E.g. doi:10.1594/WDCC/CMIP5.MXELr4

**DOI for Scientific Data**  
10.1595/WDC/TEST\_AMP\_TR

**Title**  
cmip5 output MPI-M ECHAM6-M

**Citation**  
Lautenschlager, Michael (2011)

**Publication Date**  
2011-02-17

**Contact for data entity**  
[Joerg Wegner](#)

**CMIP5 Metadata hoste**  
<http://cera-www.dkrz.de/WDC/>

**Summary**  
amip is an experiment of the CMIP5 experiments for the next five years.

3.3 amip (3.3 AMIP): AMIP (1971-2030)

Experiment design is described in [pcmdi.llnl.gov/cmip5/docs/standard\\_output](http://pcmdi.llnl.gov/cmip5/docs/standard_output). The output is stored in netCDF repository.

**Quality**  
**Accuracy:** not filled

**Consistency:** Quality Control Level 0: Spot checks on selected variables  
\* Level 1: CMOR2 and ESG put  
\* Level 2: Technical checks on files  
\* Level 3: Data approved by authors

**Completeness:** not filled

**Specification:** [CMIP5:QualityLevel=0]  
[CMIP5:QualityControl2Comments=] sub-T=6hrPlev, var=va, dim=time, metadata not found or not accessible

## Data for CMIP5 experiment MXETam

Please choose your desired data destination. The displayed list of links refer to ESG datasets, i.e. collections of data belonging to a model realm / ensemble member (see DRS syntax definition [\[http://cmip-pcmdi.llnl.gov/cmip5/docs/cmip5\\_data\\_reference\\_syntax.pdf\]](http://cmip-pcmdi.llnl.gov/cmip5/docs/cmip5_data_reference_syntax.pdf)). Links resolve to the corresponding Thredds Data Server entries.

**Location**

### Dataset

- [cmip5\\_output MPI-M ECHAM6-MPIOM-TR amip\\_mon\\_atmos Amon r2i1p1 v20100928](#)
- [cmip5\\_output MPI-M ECHAM6-MPIOM-TR amip\\_fx\\_atmos fx r0i0p0 v20100928](#)
- [cmip5\\_output MPI-M ECHAM6-MPIOM-TR amip\\_mon\\_atmos Amon r1i1p1 v20100928](#)
- [cmip5\\_output MPI-M ECHAM6-MPIOM-TR amip\\_mon\\_landlce Lmon r1i1p1 v20100928](#)
- [cmip5\\_output MPI-M ECHAM6-MPIOM-TR amip\\_mon\\_land Lmon r1i1p1 v20100928](#)
- [cmip5\\_output MPI-M ECHAM6-MPIOM-TR amip\\_6hr\\_atmos\\_6hrPlev r2i1p1 v20100928](#)
- [cmip5\\_output MPI-M ECHAM6-MPIOM-TR amip\\_day\\_atmos\\_day r1i1p1 v20100928](#)
- [cmip5\\_output MPI-M ECHAM6-MPIOM-TR amip\\_day\\_land\\_day r1i1p1 v20100928](#)
- [cmip5\\_output MPI-M ECHAM6-MPIOM-TR amip\\_day\\_atmos\\_day r2i1p1 v20100928](#)
- [cmip5\\_output MPI-M ECHAM6-MPIOM-TR amip\\_mon\\_land Lmon r2i1p1 v20100928](#)
- [cmip5\\_output MPI-M ECHAM6-MPIOM-TR amip\\_6hr\\_atmos\\_6hrPlev r1i1p1 v20100928](#)
- [cmip5\\_output MPI-M ECHAM6-MPIOM-TR amip\\_fx\\_land\\_fx r0i0p0 v20100928](#)
- [cmip5\\_output MPI-M ECHAM6-MPIOM-TR amip\\_mon\\_landlce Lmon r2i1p1 v20100928](#)

Back to [the metadata page](#).

This page is hosted at [WDC](#), please send technical inquiries to [data@dkrz.de](mailto:data@dkrz.de).

### Link to primary data

[CMIP5Links.jsp?acronym=MXETam](#) (to be replaced by link into PCMDI gateway)

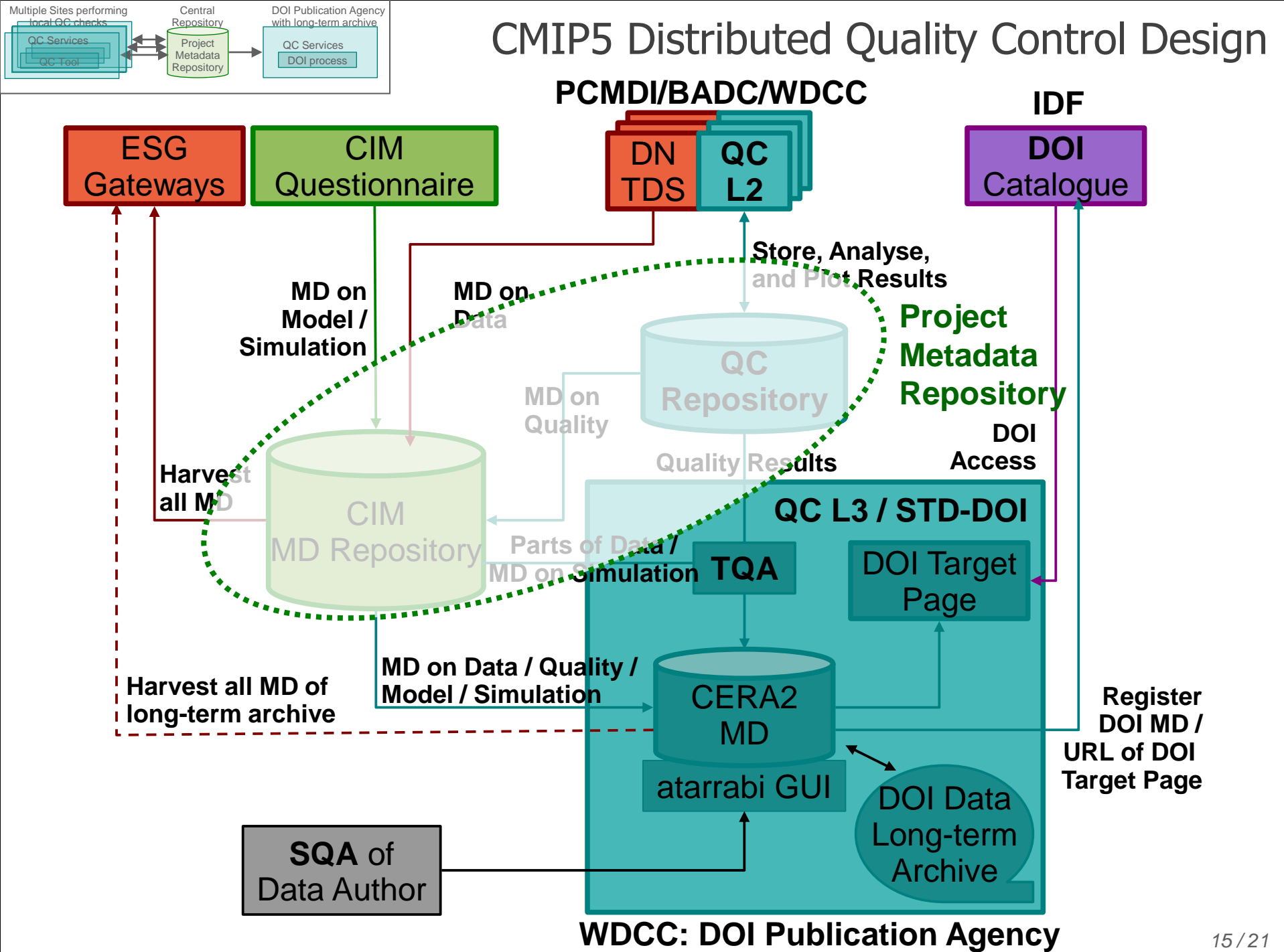
### Please note

WDC Climate as DOI publishing agency grants the validity of data accessed at WDC Climate. It recommends to check all downloaded data files by their tracking\_id with the CMIP5 Data Validation Service at [http://cera-www.dkrz.de/CMIP5Tracking.jsp?tracking\\_id=tracking\\_id](http://cera-www.dkrz.de/CMIP5Tracking.jsp?tracking_id=tracking_id). The tracking\_id can be found in the file headers (e.g. using `ncdump -h`).





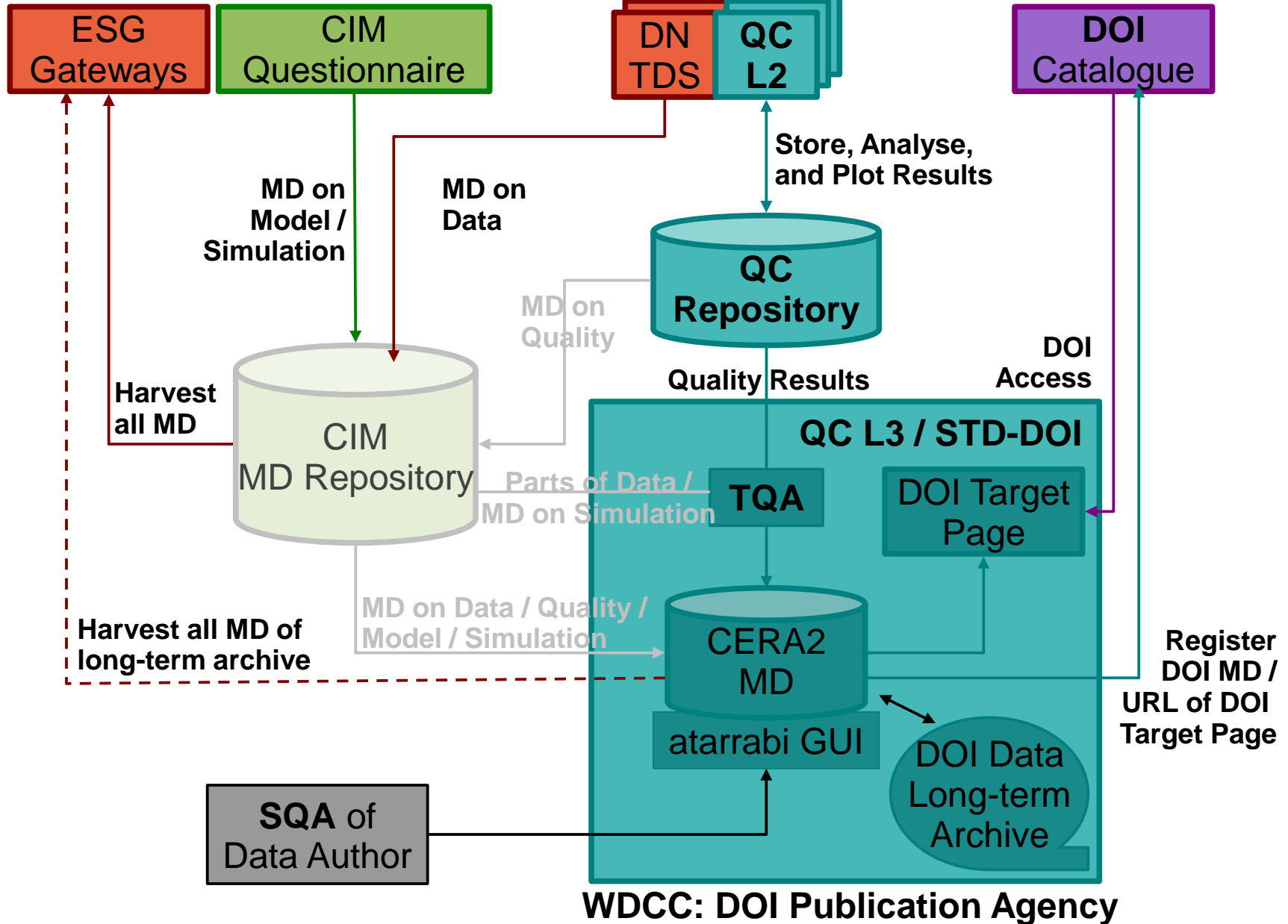
# CMIP5 Distributed Quality Control Design



# CMIP5 Distributed Quality Control Status

PCMDI/BADC/WDC

IDF



## QC Development:

## Status QC in CMIP5 (1)

- **QC L2:**  
QCDB moved from postgres to oracle for production due to maintenance reasons (21st April 2011);  
bug-fixes in QC tool
- **QC L3:**  
DOI data publication moved from STD-DOI to DataCite (21st April 2011)  
SQA GUI atarrabi final release 1.5 (2nd May 2011):  
[cera-www.dkrz.de/atarrabi/](http://cera-www.dkrz.de/atarrabi/)
- **DOI landing page:**  
set up: [cera-www.dkrz.de/CMIP5/CMIP5Compact.jsp?acronym=<acronym>](http://cera-www.dkrz.de/CMIP5/CMIP5Compact.jsp?acronym=<acronym>)  
design under discussion:  
[cera-www.dkrz.de/CMIP5/CMIP5example.html](http://cera-www.dkrz.de/CMIP5/CMIP5example.html)  
*data set access realized by a list of ESG dataset TDS entry points because of the DOI granularity on DRS experiment level*
- **DOI services:**
  - tracking\_id service: [cera-www.dkrz.de/CMIP5/CMIP5Tracking.jsp](http://cera-www.dkrz.de/CMIP5/CMIP5Tracking.jsp)
  - view for QC status/DOI in place: `cera2.v_qc_status`
  - view for ESG datasets belonging to DOI in place: `cera2.v_ipcc_files`

## CIM access:

## Status QC in CMIP5 (2)

- **Get Contact/Authors/Title from simulationRun object:**  
*no access of AtomFeed entries by DRS experiment name*  
-> use of persons from email receiver list for QC L3 process
- **Ingest of QC L2:**  
QC questionnaire for QC check definition exists  
Ingest tool for QC result ingest into CIM missing  
*no QC flag communication to gateways, no QC results in CIM*  
-> use WDCC view cera2.v\_qc\_status (cera2.v\_ipcc\_files)
- **Update CIM after QC L3 / DOI publication:**  
unclear  
*no QC flag/DOI communication to gateways, no update in CIM*  
-> use WDCC view cera2.v\_qc\_status (cera2.v\_ipcc\_files)
- **DOI landing page:**  
*no link to CIM metadata available so far*  
-> use link to CERA metadata in the meantime

Technical QC Meeting planned in UK in May 24 + 25<sup>th</sup>, 2011



# Open Issues (1)

## 1. Will ANU perform QC L2 checks on the Australian data?

- a. Yes, Australian data. Contact person needed.
- b. No.

## 2. How will the gateway show QC flag and DOI?

- a. For the DOI / DRS experiment and the ESG datasets:  
Then we could use a single link out of the DOI landing page into the ESG gateway. Harvest the information from:
  - I. CIM ?
  - II. WDCC: Views `cera2.v_qc_status` und `cera2.ipcc_files`
- b. Only for DOI / DRS experiment:  
Then we do not have a connection between DOI and the datasets belonging to it in the gateways. Harvesting:
  - I. CIM ?
  - II. WDCC: View `cera2.v_qc_status`

## Open Issues (2)

### 3. How do we deal with missing or unfinished CIM metadata descriptions?

- a. Part of CIM metadata might be sufficient but even not ideal: Specification of mandatory parts or percentage of completeness of CIM needed.
  - I. QC on data is assigned QC L2 without CIM metadata. For DOI publication CIM metadata have to exist, i.e. published by AtomFeed.
  - II. QC stops before assignment of QC L2 if no or insufficient CIM metadata exists.
- b. CIM metadata is regarded as optional: QC L3 is finished by DOI assignment only with TDS and DOI metadata. Data Authors have to complete metadata for DOI data publication.

### 4. How do we ensure identical QC L2 data checker application at all three (four) sites? This is crucial for QC-L3 and the DOI publication process.

- a. Use of identical software version of QC tool and in same configuration.
- b. Otherwise during QC L3 the QC L2 checks have to be inferred and completed if necessary.

