

QA-DKRZ: The Annotation Model

H.-D. Hollweg

DKRZ, hollweg@dkrz.de

Overview

- **QA-DKRZ Tool**
 - Work-flow
 - Dependencies

- **Annotation Model**
 - Specification of actions tagged to checks
 - Structure of results: Files and directories
 - YAML formatted log-file output
 - JSON formatted summary

- **QA-DKRZ: status**

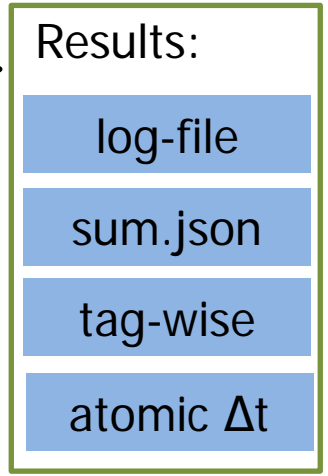
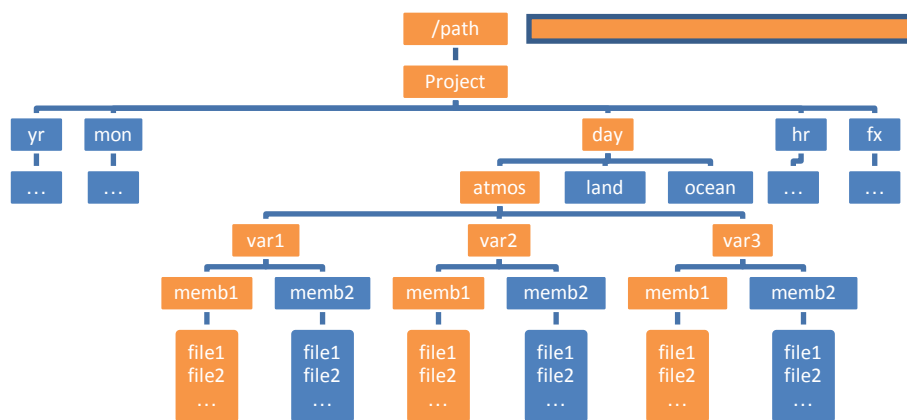
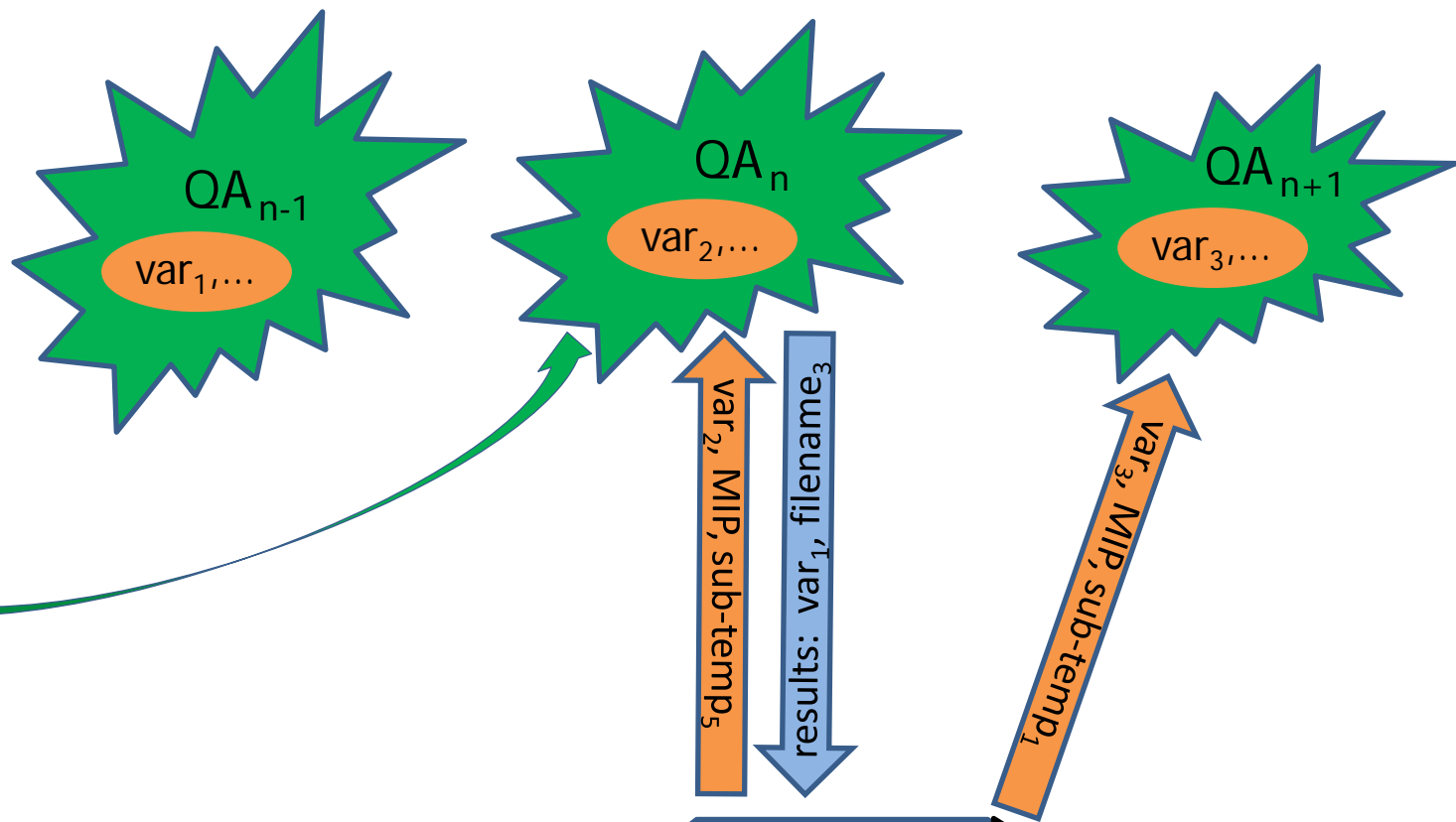
Purpose:

Assure that every file entering ESGF
complies to conventions and project rules.

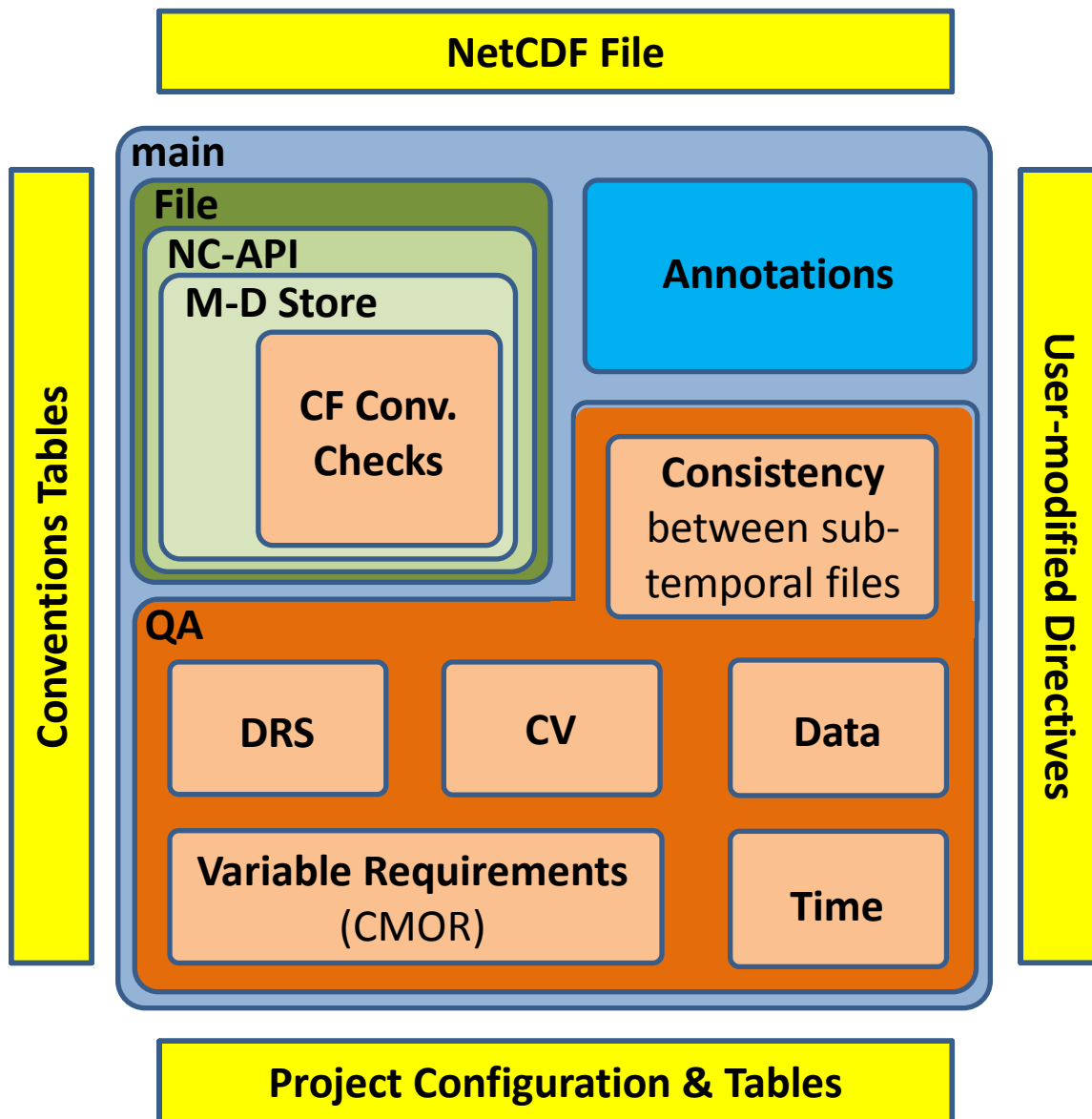
If not, then issue annotations.

Tables:

- Conventions
- Check-lists
- CV
- DRS
- Variable Requ.

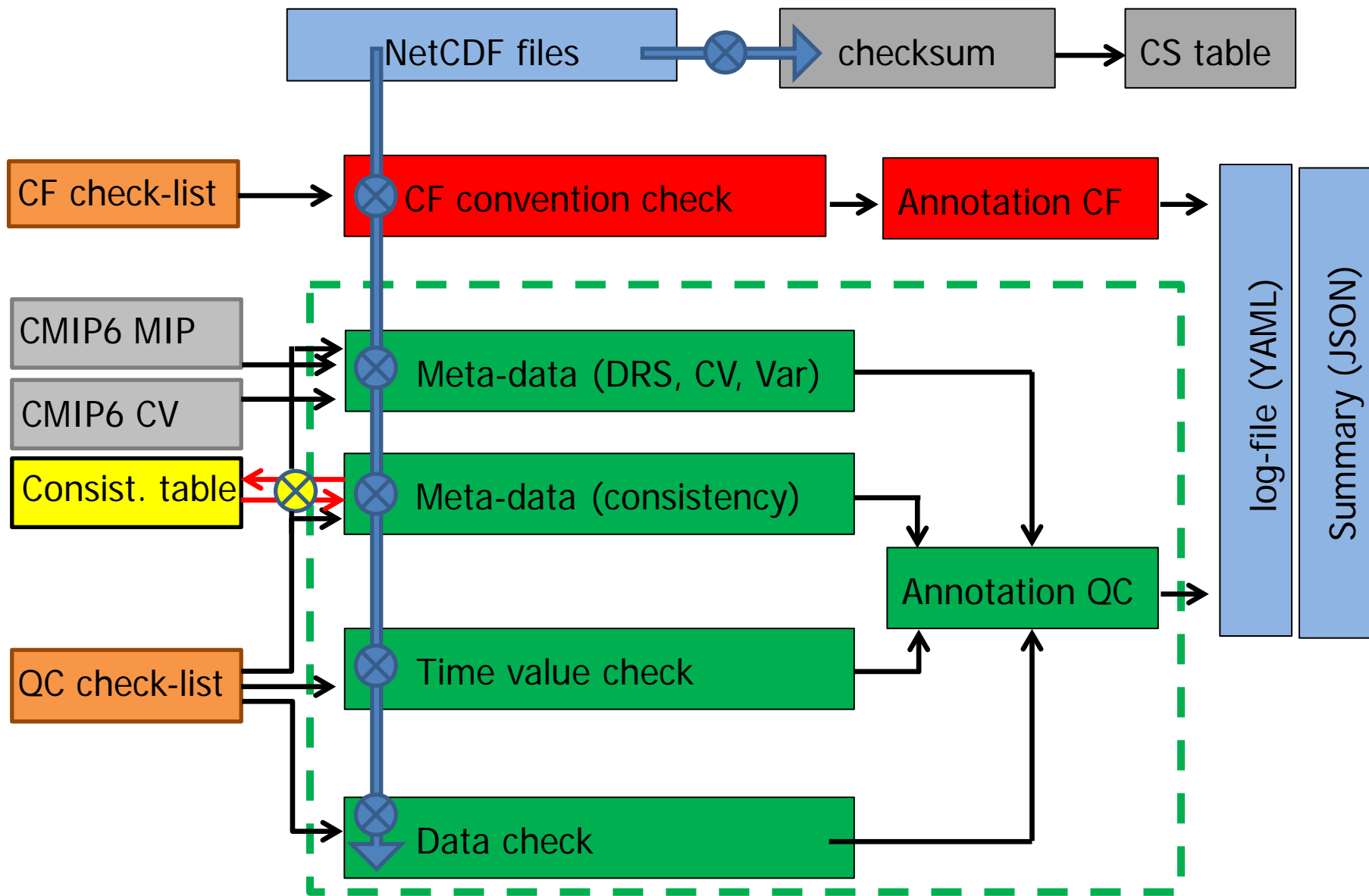


QA Program (C++)



Quality Assurance (QA)

- **Data Reference Syntax (DRS)**
 - **Controlled Vocabulary (CV)**
 - **Variable Requirements (CMIP Model Output Requir.)**
 - **Time Properties**
 - **Consistency** between parent - child files (atomic and experiments)
 - **Data Checks**
infinity and not-a-number
outlier tests
replicated record detection
- Note:**
every check may be disabled



Libraries

- zlib www.zlib.net
- hdf5 www.hdfgroup.org/HDF5
- netcdf www.unidata.ucar.edu/netcdf
- udunits2 www.unidata.ucar.edu/software/udunits

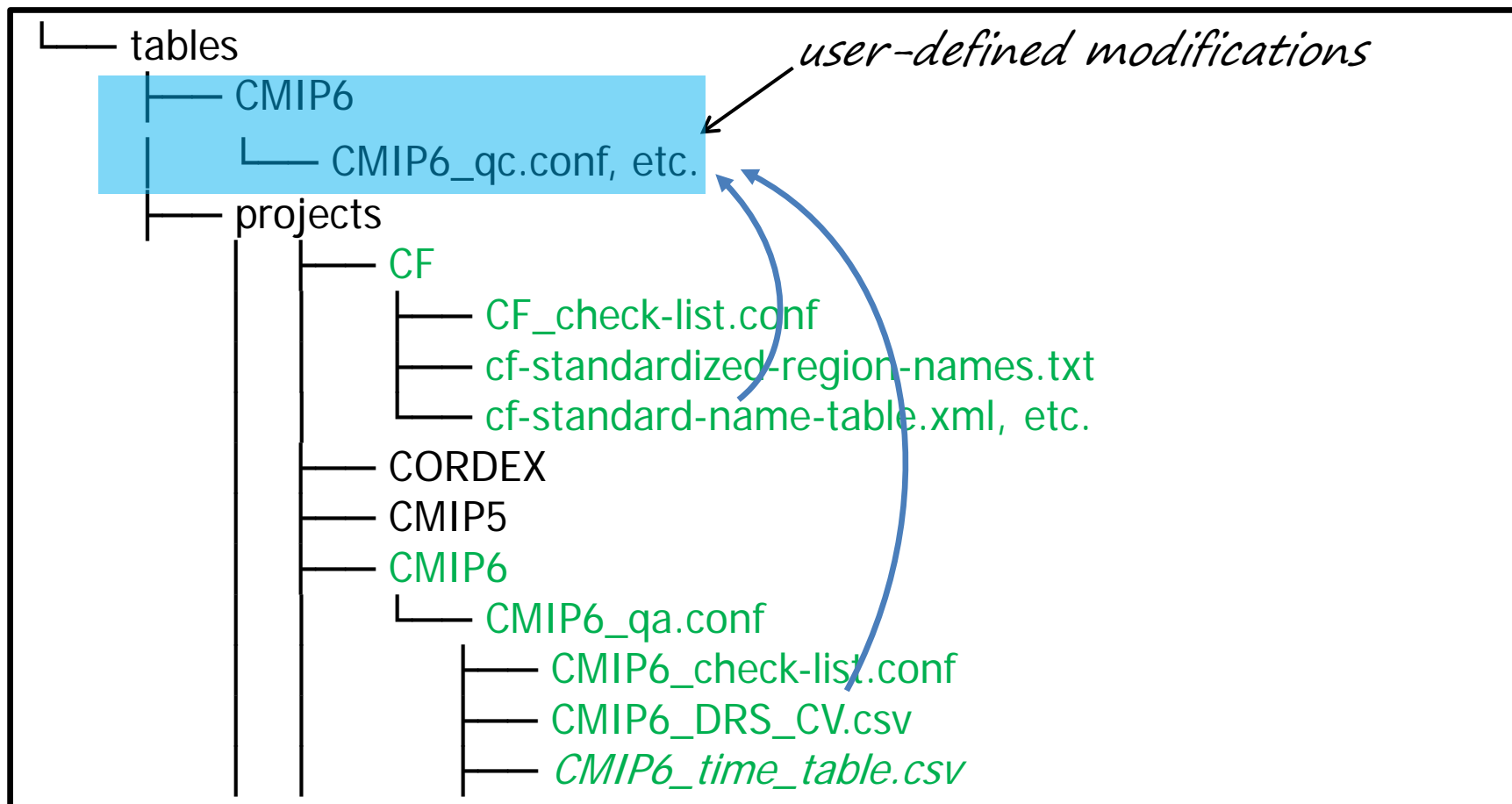
Tables

- CF Conv. <http://cfconventions.org>
- CMIP6_MIP http://proj.badc.rl.ac.uk/svn/exarch/CMIP6dreq/tags/latest/dreqPy/docs/CMIP6_MIP_tables.xlsx
- CMIP6_CV https://github.com/WCRP-CMIP/CMIP6_CVs

Externals

- xlsx2csv <http://github.com/dilshod/xlsx2csv>
- jsoncpp <https://github.com/open-source-parsers/jsoncpp>
- PrePARE http://cmor.llnl.gov/mydoc_cmip6_validator

Path: /home/user/.qa-dkrz



Precedence of Tables (Directives)

- Command-line
- Current Directory (Task File)
- User-defined Modifications
- Default (Projects Directory)

QA-DKRZ

- **Sources: GitHub**

<https://github.com/IS-ENES-Data/QA-DKRZ>

- **Binaries**

conda install -c birdhouse -c conda-forge qa-dkrz

ehbrecht@dkrz.de

- **Documentation: ReadTheDocs.org**

<http://qa-dkrz.readthedocs.io/en/latest>

Annotation Model

- Check-lists
- QA_Results
 - Log-file (YAML)
 - Summary
 - Annotations (JSON)
 - Atomic Time Interval

Check-list File

Format: [text] & tag [,level] [,task] [,variable] [,constraint]

Brace grouping {}:

Example: given: a,b{v{D(z),x,b=2}},{u,v},w

result: 'a,b,w', 'a,v,x,b=2,w', 'a,b,u,v,w'

Key words of actions: {Ln, D, EM, tag, var, V=value, R=record}

- level: L1 – L4 (warning – emergency stop)
- D: Discard
- tag: Identifier.
- EM: Email notification (EM)
- var: Comma-separated acronyms of variables; directive is only applied to these variable(s).
- value: Constraining value, *e.g {tag,D,V=0,var} discards test for variable var only if value=0*
- record: apply to time value(s) r_0 [- r_1]

Check-list File

Groups:

- (1) Directory and Filename Structure (DRS)
- (2) Required Attributes (CV)
- (3) Variables
- (4) Dimensions
- (5) Auxiliaries
- (6) Tables
- (7) Consistency Check
- (8) Miscellaneous
- (M) Pre-check failure detection
- (R) Data / Time (record/based)

Examples (from `CORDEX_check-list.conf`):

Height requires units=m

& `2_3,L1`

every height variable is checked for units [m]

Near-surface height must be 0 - 10m

& `5_6,L1,{D,rlut,rsdt,rsut}`

variables discarded from check: rlut, rsdt, rsut

Suspecting replicated records

& `R3200,L1{D,sund},{D,V=0,clivi,mrfso,prsn,sftgif}`

sund discarded,

clivi ... discarded for records

with constant value=0.

QA Results

Structure of QA-Results: Files and Directories

QA_Results/project/institute (root directory)

└─ check_logs

├─ AFR-44_CNRM_ECMWF-ERAINT_evaluation_r1i1p1_v1.log

├─ Annotation

├─ AFR-44_CNRM_ECMWF-ERAINT_evaluation_r1i1p1_v1.json

├─ AtomicTimeRange

├─ AFR-44_CNRM_ECMWF-ERAINT_evaluation_r1i1p1_v1.period

├─ AFR-44_CNRM_ECMWF-ERAINT_evaluation_r1i1p1_v1.range

└─ Tags

├─ AFR-44_CNRM_ECMWF-ERAINT_evaluation_r1i1p1_v1

├─ CF_73d

├─ L1_1_2

├─ L1_CMOR_xzy

├─ L2_SF

Log-file (YAML)

Log-file of a QA session started by qa-DKRZ

configuration:

command-line: -m -f task.CMIP6 -e_check_mode=-CNSTY -e_next

options:

APPLY_MAXIMUM_DATE_RANGE:

...

SELECT_VAR_LIST: .*

start:

date: 2016-12-02T11:23:38

qa-revision: master-66ca331

items:

- **date:** 2016-12-02T11:23:40

file: tas_Amon_1pctCO2_MPI-ESM-LR_r1i1p1f2_gn_200601-210012.nc

data_path: /path/CMIP6/CMIP/MPI-M/.../r1i1p1f2/Amon/tas/gn/v20161130

conclusion: 'CF: FAIL, CV: FAIL, DATA: PASS, DRS(F): PASS, DRS(P): FAIL, TIME: PASS

checksum: ce5e24ffeb5c38665a17570f4a564f0e.md5

creation_date: 2016-12-02T12:40:29Z

tracking_id: 06cfd581-917a-4888-9b92-a07a726469d0

events:

- event:

caption: 'DRS path: path component member_id=<r1i1p1f2> does not match global attribute value <r1i1p1f1>.'

impact: L1

tag: '1_2'

- event:

caption: 'Attribute institution:
found <Max Planck Institute for Meteorology>,
expected from CMIP6_institution_id.json
<Max Planck Institute for Meteorology, Hamburg 20146,
Germany>.'

impact: L2

tag: '2_4'

- event:

caption: 'Coordinate variable <height>: No data.'

impact: L1

tag: 'CF_0d,

status: 2

Time Intervals of atomic Variables (*YAML*):

```

--- # Time intervals of atomic variables.
- frequency: day
  number_of_variables: 42
    - variable: evspsbl_AFR-44_ECMWF-ERAINT_evaluation_r1i1p1_v1_day
      begin: 1989-01-01T00:00:00
      end: 2009-01-01T00:00:00
      status: PASS | FAIL:B | FAIL:E
    - ...
- frequency: mon
  number_of_variables: 71
    - variable: evspsbl_AFR-44_ECMWF-ERAINT_evaluation_r1i1p1_v1_mon
      begin: 1989-01-01T00:00:00
      end: 2009-01-01T00:00:00
      status: PASS | FAIL:B | FAIL:E
    - ...
  
```

Note: FAIL:B | FAIL:E means that not all files begin|end with the same date.

Time Intervals of atomic Variables (*human read.*):

Frequency: day

Number of variables: 4

clh_EUR-11_ECMWF-ERAINT...day	1979-01-01T00:00:00 - 2013-01-01T00:00:00
clivi_EUR-11_ECMWF-ERAINT...day	--> 1980-01-01T00:00:00 - 2013-01-01T00:00:00
cli_EUR-11_ECMWF-ERAINT...day	1979-01-01T00:00:00 - 2010-01-01T00:00:00 <--
clm_EUR-11_ECMWF-ERAINT...day	1979-01-01T00:00:00 - 2013-01-01T00:00:00

Frequency: fx

Number of variables: 1

orog_EUR-11_ECMWF-ERAINT	-
--------------------------	---

Frequency: mon

Number of variables: 4

clt_EUR-11_ECMWF-ERAINT...mon	--> 1980-01-01T00:00:00 - 2013-01-01T00:00:00
evspsbl_EUR-11_ECMWF-ERAINT...mon	1979-01-01T00:00:00 - 2013-01-01T00:00:00
hfls_EUR-11_ECMWF-ERAINT...mon	--> 1980-01-01T00:00:00 - 2010-01-01T00:00:00 <--
hfss_EUR-11_ECMWF-ERAINT...mon	1979-01-01T00:00:00 - 2013-01-01T00:00:00

Summary (JSON)

```

{
  "QA_conclusion": [ PASS | FAIL ] ",
  "project": "CORDEX",
  "DRS_0": "cordex",
  "DRS_1": "output",
  "DRS_2": "AFR-44",
  ...
  "DRS_8": "v1",
  "DRS_9": "SHARED",
  "DRS_10": "SHARED",
  "annotation":
  [
    {
      "DRS_9": ["day", "mon"],
      "DRS_10": ["tauv", "zg500"],
      "caption": "DRS CV path: global attribute RCMModelName = <QWER> vs. <ASDF>.",
      "severity": "L1"
    },
    {...}
  ]
}

```

CMIP6 Files → ESGF

QA Procedure

- **Check (only) DRS of paths and filenames.**
- **Run PrePARE checker for CMIP6 CV.**

EXAMPLE: CMIP6 Test File with Faults

QA-DKRZ: DRS Check

- event:

capt: DRS path component

member_id=<r1i1p1f2> does not match
global attribute value <r1i1p1f1>.

impact: L1

tag: 1_2

QA-DKRZ: status

		CMIP5	CORDEX	CMIP6	Comment
Conv	CF	v1.4	v1.4	v1.7	www.cfconventions.org
	UGRID	-	-	v1.0	ugrid-conventions.github.io
DRS	(Path)				
	(File)				
CV		1)			1) CMOR guide → machine read.
Var. Requir.				2)	2) CMIP6_MIP_tables.xlsx
Consistency					files across atomic & exp. scope
Time					
Data					NaN, Inf, replications, outlier
CMOR		-	-	PrePARE	http://cmor.llnl.gov
WPS					
OpenDAP					