



# Development and exploitation of a controlled vocabulary in support of climate modelling

M.-P. Moine<sup>1</sup>, S. Valcke<sup>1</sup>, B. N. Lawrence<sup>2,3,4</sup>, C. Pascoe<sup>4</sup>, R. W. Ford<sup>5</sup>, A. Alias<sup>6</sup>, V. Balaji<sup>7</sup>, P. Bentley<sup>8</sup>, G. Devine<sup>9</sup>, S. A. Callaghan<sup>4</sup>, and E. Guilyardi<sup>9,10</sup>

<sup>1</sup>Centre Européen de Recherche et de Formation Avancée en Calcul Scientifique (CERFACS), CERFACS/CNRS SUC URA1875, Toulouse, France

<sup>2</sup>Department of Meteorology, University of Reading, Reading, UK

<sup>3</sup>Centre for Environmental Data Archival, STFC Rutherford Appleton Laboratory, Didcot, UK

<sup>4</sup>National Centre for Atmospheric Science (NCAS), Natural Environment Research Council, UK

<sup>5</sup>STFC Daresbury Laboratory, Warrington, UK

<sup>6</sup>Centre National de Recherches Météorologiques (CNRM), Meteo-France/CNRS, Toulouse, France

<sup>7</sup>NOAA Geophysical Fluid Dynamics Laboratory (GFDL) and University of Princeton, USA

<sup>8</sup>Met Office Hadley Centre, Exeter, UK

<sup>9</sup>National Center for Atmospheric Science (NCAS), University of Reading, Reading, UK

<sup>10</sup>Institut Pierre Simon Laplace (IPSL), CNRS, Paris, France

*Correspondence to:* M.-P. Moine (marie-pierre.moine@cerfacs.fr)

Received: 12 April 2013 – Published in Geosci. Model Dev. Discuss.: 23 May 2013

Revised: 20 January 2014 – Accepted: 26 January 2014 – Published: 21 March 2014

**Abstract.** There are three key components for developing a metadata system: a container structure laying out the key semantic issues of interest and their relationships; an extensible controlled vocabulary providing possible content; and tools to create and manipulate that content. While metadata systems must allow users to enter their own information, the use of a controlled vocabulary both imposes consistency of definition and ensures comparability of the objects described. Here we describe the controlled vocabulary (CV) and metadata creation tool built by the METAFOR project for use in the context of describing the climate models, simulations and experiments of the fifth Coupled Model Intercomparison Project (CMIP5). The CV and resulting tool chain introduced here is designed for extensibility and reuse and should find applicability in many more projects.

the earth system are being represented with an increasing number of physical processes taken into account. Higher spatial resolution is supported thanks to the emergence of high performance computing platforms. In addition, more and more research centres have been engaging in climate modelling, which increases the number of models involved. One important consequence is the growth of the volume of data produced. Climate Model Intercomparison Projects (CMIP) initiated and supervised by the World Climate Research Programme (WCRP) are an academic exercise on which climate projection assessment is based. Higher complexity of numerical models, explosion in the volume of data produced and the growing number of contributing modelling groups require a dedicated and expert infrastructure for data quality control, data documentation, data storage and access. Indeed, the ever-growing number of scientific groups producing and using climate model data requires more sophisticated data management systems, including good quality, understandable and shareable data documentation.

The technological part of this infrastructure in CMIP5 is ensured by ESGF (Earth System Grid Federation): several distributed data centres host the data produced by the

## 1 Introduction

Climate models have experienced outstanding evolution in the last 20 years, driven by scientific improvements and increases in computing capabilities. Additional components of

modelling groups around the world, some of them (PCMDI, BADC, WDCC) being gateways for data publication and download (Williams et al., 2011). It became clear during the set-up of this infrastructure that the definition and adoption of standard metadata (that is data describing the data), is crucial to guide end-users through data mining, data interpretation or data comparison tasks (Guilyardi et al., 2011) – even outside the climate modelling community itself, for example by the environment and health impact community. Furthermore, climate metadata must describe both the data content and the model and simulations that produced this data. The CMIP5 metadata standardization effort exploited work conducted jointly by the CURATOR project in the US (Dunlap et al., 2008) and the METAFOR project funded by the European Commission (Callaghan et al., 2010). The approach we followed in METAFOR was to define three key metadata components: a conceptual container to store and organize the information (the CIM, Common Information Model); the possible content (the controlled vocabulary) and a methodology to harvest a specific content, i.e. an instance of metadata (the so-called “CMIP5 Questionnaire”). The CIM is introduced in Lawrence et al. (2012). Here we concentrate on the controlled vocabulary (CV) and the specific harvesting tool developed for CMIP5. We begin by setting the context of earth system models and simulations so as to appreciate the challenge raised by climate metadata. We then present a brief inventory of existing metadata systems in the climate area, pointing out gaps and incompleteness and advocating for a unique and encompassing standard. We describe the methodology applied to build the METAFOR CV and its resulting structure based on key elements. Finally, we explain how this CV was used to construct the “CMIP5 Questionnaire” and how it was ingested by other metadata systems like ESGF.

## 2 Picture of a climate model and climate experiments

Climate study is a highly interdisciplinary science that historically emerged with the convergence of scientific expertise in the research areas related to the earth system, such as oceanography, atmospheric physics, sea ice dynamics, hydrology, etc. As a result, a climate model is a composition of models (hereafter referred to as “components”, some of which map onto “realms” using the nomenclature of Taylor et al. (2011a)), each one being devoted to a specific domain of the climate system. These models are generally assembled by coupling software (see Valcke et al., 2012, for a review). The role of the coupler is to exchange coupling fields at the interface of the component domains (for example, wind stress and radiative fluxes are transmitted from the atmosphere to the ocean, sea surface temperature and currents from the ocean to the atmosphere), performing the spatial remapping from the grid of one component to the other. The resulting global

model, including components and the coupler, is therefore referred to as a “coupled model”.

A given model can be run and integrated in time (i.e. a climate simulation can be performed) in a large number of different ways, depending on the temporal and dynamical schemes used, and according to the physical parameterizations selected to model subgrid phenomena within each physical scheme of each component. Initial conditions and external forcing that influence the climate system must be prescribed, e.g. green house gases, volcanoes, aerosol types and concentrations, and land-use changes. By adjusting model parameters such as orbital parameters or solar irradiance and by applying appropriate forcing and initial conditions, climate models can be run for various time durations (seasonal, decadal, centennial, millennial) and reproduce different climatic periods (paleo, present and future).

One particular model configuration is usually targeted at a specific scientific question: for example, to understand the sensitivity of a climate process to horizontal resolution or to provide a projection of future climate under a specific emission scenario. Hence, it is important to document not only the particular configuration, but also why that configuration was chosen. The purpose of an experimental protocol like CMIP5 (Taylor et al., 2011a) is to provide guidance for the set-up of models and simulations, so that the different modelling groups address the same questions in a comparable way with their own model. It is clear that the way the model is scientifically configured (including model parameterizations, initial conditions and forcing) and how it conforms to the experimental requirements is crucial information to interpret and compare results. It is therefore vitally important to preserve this information along with the data.

## 3 Existing metadata for weather forecast and climate

To ensure interoperability of geo-referenced and weather forecast data products, international organizations like the Open Geospatial Consortium (OGC) and the World Meteorological Organization (WMO) promote adoption of standards. These standards are currently used by national meteorology institutes and production centres of remote sensing and in situ observations all over the world.

The CF convention provides a set of standard names for geophysical variables associated with a precise scientific definition and units. In the CMIP5 framework, CF-NetCDF is the compulsory format for the output data set. Furthermore, CMIP5 output metadata are constrained by the CMIP5 tables which impose, among other things, short names and units and ensures correspondence with the CF standard names, both for dimensional and physical variables (Taylor and Doutriaux, 2010). Additional low-level<sup>1</sup> metadata are included in the output files as global attributes, for example *experiment\_id*

<sup>1</sup>“low-level metadata” term refers to metadata that applies to individual data sets and describes their content (i.e. what the data is),

or *model\_id* that respectively identify the CMIP5 experiment and the coupled model that produced the data set, according to terms defined in the Data Reference Syntax document (DRS) (Taylor et al., 2011b).

Several previous projects, such as NMM (Numerical Model Source Metadata, University of Reading) and NumSim (Numerical Simulation Discovery Metadata, BADC/NCAS, <http://proj.badc.rl.ac.uk/ndg/wiki/NumSim>) have tried to address the higher-level metadata issue, i.e. not only describing “what” are the data produced but also “how” they were produced (the model and simulation details). NMM and NumSim identified some key terms (e.g. *genealogy*, *boundary condition type*, *initial condition type*, *ensemble type*, *model component*, *model category*) and used ISO standards where relevant. Other specific metadata systems have addressed more technical aspects of climate modelling like the configuration of coupling exchanges between earth system components (BFG, Ford and Riley, 2011; OASIS4, Redler et al., 2010) or the grids on which climate model data is discretized (gridSpec, Balaji Institute, 2007). However, no one integrated high-level<sup>1</sup> metadata system able to encompass the whole “climate modelling” process emerged from these projects, leaving only pieces of metadata, often disconnected. In the previous CMIP phase 3, this resulted in asking scientists to provide additional information about models and simulations in unconstrained text-based documents (the CMIP3 questionnaire, see an extract in Appendix B).

#### 4 The METAFOR controlled vocabulary

Given that the metadata have to address all stages of the modelling process and given that they should serve data discovery and access tools, the prime objective of the METAFOR project was to design a conceptual metadata scheme and develop the associated hosting structure, the Common Information Model (CIM). The CIM defines objects, classes, and their relationships (Lawrence et al., 2012). Through specialized UML (Unified Modelling Language, [www.uml.org](http://www.uml.org)) packages, the CIM addresses the description of the constituent elements of climate modelling: the “activity” package includes the experimental context and simulations; the “software” package covers the climate model itself; the final data objects produced by simulations and their inputs are described by the “data” package and the numerical grids of the models by the “grid” package; finally, a “shared” package of reusable elements supports some “orphan” classes, such as quality control records and platform descriptions.

To be operational, each individual CIM package needs an associated controlled vocabulary (CV) that defines sets of allowed attributes (name/value pairs). For the “data”, “grid” and “shared” packages, the CV was mainly based on a list

<sup>1</sup>while “high-level metadata” refers to metadata that applies to whole data sets and addresses how the data were produced.

of already existing terms, respectively the CF standard, grid-Spec and some ISO standards. Vocabularies for the activity and software packages did not exist, and were developed from scratch. In the following we present the resulting “Model Controlled Vocabulary” and the “Simulations and Experiments Controlled Vocabulary”, used in support of CMIP5 to populate the software and activity packages respectively.

##### 4.1 The Model Controlled Vocabulary

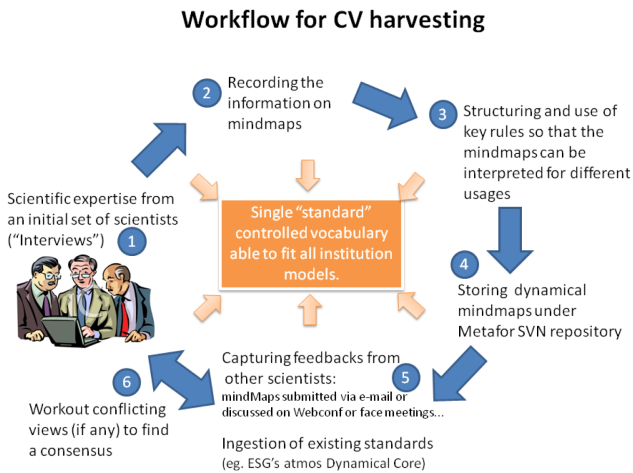
The Model Controlled Vocabulary describes the heart of the climate data production chain, that is, the numerical model itself. This work had to define the Model CV starting from scratch and had to go through the early steps of a classical CV building process:

1. identify the relevant and discriminating information (about the climate model components);
2. set an ensemble of appropriate terms (meaningful and non-ambiguous) to synthetically and faithfully express the information;
3. organize these terms hierarchically, with possible inter-dependencies;
4. attach a definition to each term;
5. identify allowed/possible values for each term.

Following the CMIP5 protocol (Taylor et al., 2011a), the first level decomposition of a coupled climate model was mapped onto eight identified realm components: *ocean*, *atmosphere*, *land surface*, *land ice*, *sea ice*, *atmospheric chemistry*, *aerosol* and *ocean biogeochemistry*. Each realm component is in its turn made of sub-components, one per main physical or dynamical process. Here the components are logical descriptions of the model, not descriptions of the actual software – it is important that users of these CVs understand the distinction, since with the version of the CIM used, there is not necessarily a direct mapping between the description of the components and the actual layout in software components.

The way of organizing the CV was both driven by typical structure of the numerical models themselves and by the scientific rationale for gathering ideas into main themes, the two being obviously closely related. The current CV granularity is a compromise driven by the requirements of model inter-comparison: to reach a level of details sufficient to be meaningful and discriminating across the various climate models but avoid overloading and too-specific information.

The CV could not be established ad hoc by exploring the model literature alone. The compromise reached is the result of a wide consultation with a number of climate modellers led by one dedicated person in METAFOR. The resulting collaboration of a significant number of scientists



**Fig. 1.** Consultation process with scientists to define the CV for climate model description.

from the international climate community, each working with different climate models, was a key part of the CV development. More than 35 experts from 13 research centres representing 6 countries contributed (see list of contributors in Appendix A), each bringing important scientific expertise to help in identifying the model characteristics important to capture and document for intercomparison. During face meetings or through audio screen-sharing sessions, modellers were asked to tell us about the science and algorithms of the climate model component they developed. The discussions were captured using mindmaps (Freemind software, [http://freemind.sourceforge.net/wiki/index.php/Main\\_Page](http://freemind.sourceforge.net/wiki/index.php/Main_Page)), one for each realm, which proved to be very appropriate for capturing structured information and feedback on the fly.

The interviewing and reviewing procedure is illustrated in Fig. 1: following a first-round interview with one realm expert (step 1), revision processes were launched with other scientists from other research centres (step 5). We integrated the feedback in a structured way (steps 2, 3), capturing their precise meaning, getting confirmation when necessary, working out possible conflicting views (step 6), and taking care not to introduce inconsistencies with previously collected CV. Following this consultation process, several iterations led to a consensus among the modellers interviewed. The resulting CV can be seen as the product of a converging process, giving ultimately both the content and the granularity of that content. For instance, the case of CV for atmospheric chemistry and aerosol modelling raised some debate within the scientific community since the CMIP5 steering committee had decided to separate them into two different realms. Intensive and rich scientific discussions and exchanges of views were necessary to raise a consensus.

The resulting scientific CV for climate models has three main categories:

1. the CV for the model realm components, including details of the numerical schemes deployed for dynamical processes (advection, diffusion, transport), for time integration and key information about the parameterizations used to model sub-grid-scale physical processes (e.g. precipitation and clouds in the atmosphere realm; soil hydrology in the land surface realm, gas phase processes in the atmospheric chemistry realm); this is the heart of the Model CV;
2. the CV associated with the numerical grids used by the models for spatial discretization;
3. the CV for describing the way components are coupled together for exchanging coupling fields, including selected terms for spatial regridding and time transformation of these fields; these latter have been derived from vocabulary used for standard configuration of couplers.

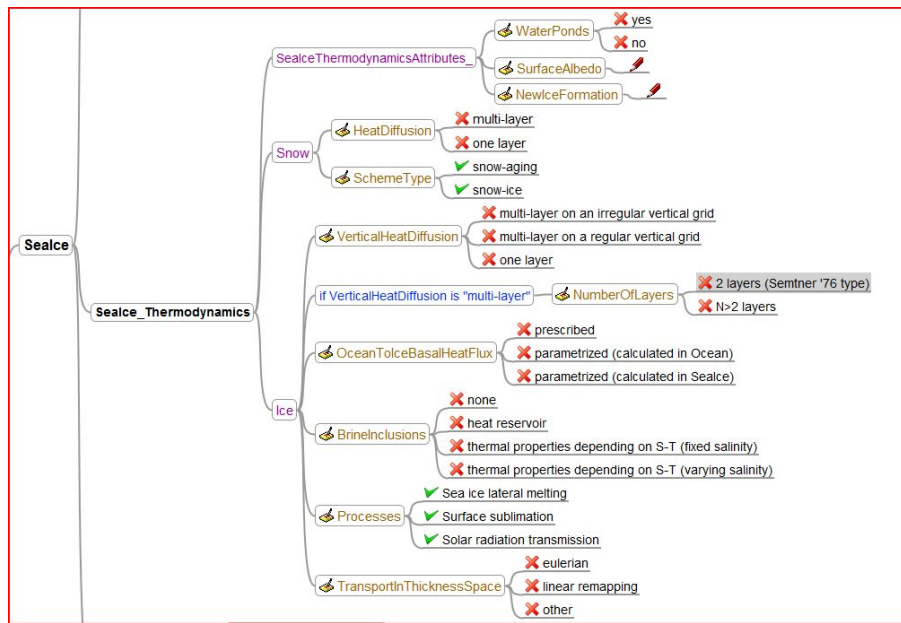
#### 4.1.1 Model realm component CV

The complete set of CV for realm components addresses more than 570 leaf parameters over 8 realms<sup>2</sup>.

The CV schema adopted for describing the model components has a hierarchical structure we illustrate with the *SeaIce* realm component (Fig. 2). The CV is made of possibly embedded elements: single "leaf parameters" (name/value pairs; e.g. *SchemeType/snow-aging* in Fig. 2) are gathered within "parameter groups" containers (e.g. *Snow*, to follow the same example in Fig. 2), themselves gathered within "components" (e.g. *SeaIce\_Thermodynamics*). Some groups of parameters are "conditional parameter groups" (e.g. *if VerticalDiffusion is multi-layer* in Fig. 2) depending on the value taken by another parameter (here *VerticalDiffusion*). The tree structure of these different container families define the allowable embedding of the controlled vocabularies and their relationships.

The CV forms a semantic database (the possible content) for building a metadata instance (an actual content recorded as a CIM document) for a given model and related simulation. A suite of tools were developed to exploit this semantic database in an automatic way so as to feed downstream tools such as the CMIP5 Questionnaire (see Sect. 5.1). To that end, coding rules were added to the mindmaps. We defined a set of formal typographic rules (e.g. different font formats and icons) to distinguish the different types of CV containers and the different types of choice (exclusive or not) among possible values for parameters or to define the type of expected value (numeric or string). These rules are illustrated in Fig. 2 and detailed in the legend of this figure. A definition of parameters is provided (as attached note, not shown in Fig. 2) and units are prescribed where numeric values are expected.

<sup>2</sup>See [http://METAFORclimate.eu/trac/browser/controlled\\_vocabularies/branches/cmip5/Software](http://METAFORclimate.eu/trac/browser/controlled_vocabularies/branches/cmip5/Software).



**Fig. 2.** A portion of the sea ice CV, showing the *SealceThermodynamics* sub-component. Black bold font denotes model components; purple is for parameter groups; blue is for conditional parameter groups; brown is for leaf parameters expecting values; black is for possible values for the leaf parameters; red cross icons mark single choice (XOR), green tick mark icons symbolize multiple choice (OR); pencil icons are for free text entry (numeric entries are also possible (not shown)); notebook icons ahead of a leaf parameter indicate that a definition is attached as a footnote.

#### 4.1.2 Model grid CV

With the model grid CV, METAFOR describes the computational grids of the model components. These grids may differ from the grid the data is expressed on, which, according to the CMIP5 guidance, should be described following the grid-Spec standard (Balaji Institute, 2007). The model numerical grid CV has to provide information about the horizontal and vertical coordinate system, the vertical coordinate used, the number of levels in the mixed layer and boundary layer, for ocean and atmosphere respectively, etc. A systematic comparison with gridSpec vocabulary was conducted prior to establishing the numerical grid CV so as to reuse terms when possible. A part of this model grid CV, dealing with the vertical coordinate system, is shown in Fig. 3: according to the value of the *VerticalCoordinateType* leaf parameter, different values for vertical coordinate are proposed (e.g. *sigma* coordinate is proposed only if the type of vertical coordinate is *terrain following*).

#### 4.1.3 Coupling exchanges CV

The CV defined in METAFOR to describe the coupling exchanges between the component models should be considered as an elementary first step. For each exchange, the source and target components are identified, and the coupling CV covers the coupling software used, the type of the spatial regridding and time transformation of the fields (if any). As

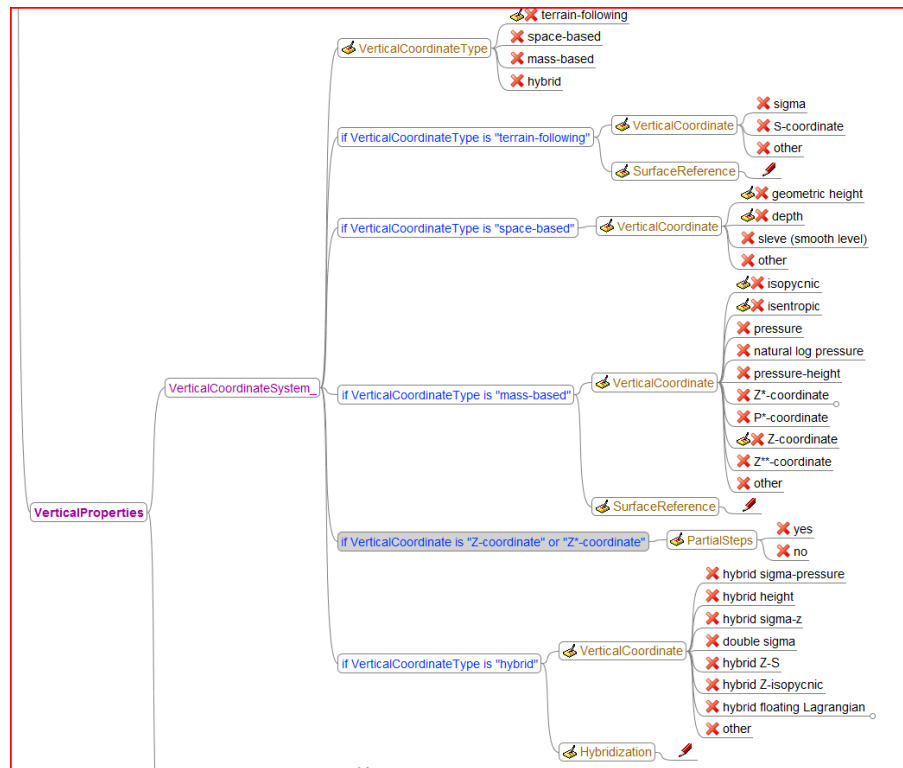
one can see, the coupling exchange CV is currently quite limited.

#### 4.1.4 Climate Model CV evolution and preservation

Even though frozen in the context of CMIP5, we expect that, with usage, this climate Model CV will evolve, improve and be reused in other scientific projects. Thus, we will have to manage the evolution and ensure the preservation of this CV, which is the first one encompassing all components of a coupled climate model. To that end, it is planned to set up an international governance committee under the auspices of IS-ENES2 (<https://verc.enes.org/ISENES2/>), the EU-FP7 (EU’s Seventh Framework Programme for Research) project that follows IS-ENES (InfraStructure for the European Network for Earth System Modelling).

#### 4.2 Controlled vocabulary for simulations and experiments

Although the Model CV discussed above is valid for any climate model, the vocabulary necessary to describe an experimental framework depends on the experiment context and aims. In contrast to the model description, METAFOR was not asked to define a specific vocabulary for experiments and simulations, the latter being extensively defined in the CMIP5 experiment design document (Taylor et al., 2011a). This document addresses two main sets of experiments, long-term and near-term, further subdivided according to



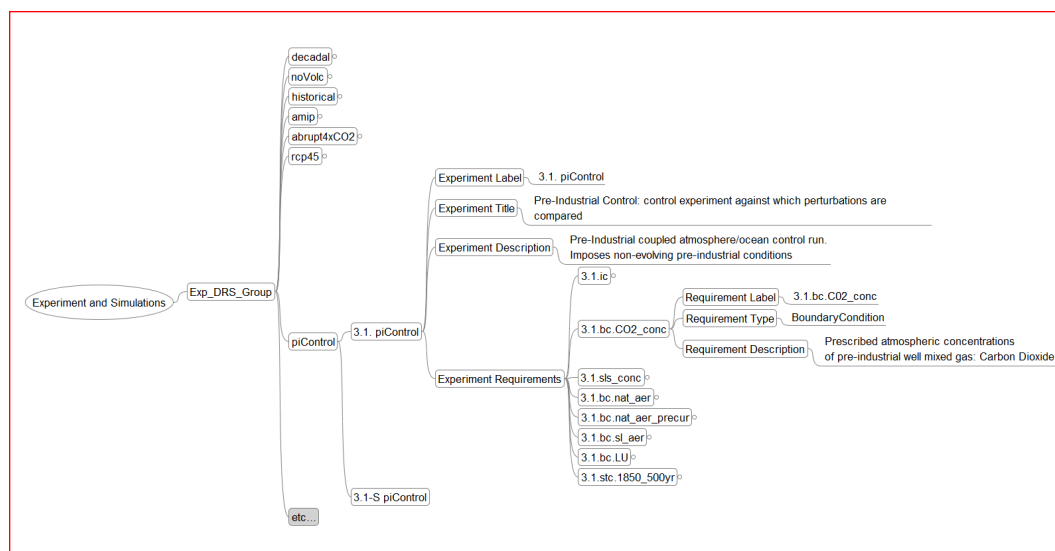
**Fig. 3.** A portion of the model grid CV. Same typographic rules as in Fig. 2 are applied. The parameter group *VerticalCoordinateSystem* gathers information about the vertical coordinate system used by the model.

distinct scientific purposes: study of a particular time period (e.g. mid-Holocene, Last Glacial Maximum or 20th century long-term experiments), analysis of the climate response to a given forcing scenario (e.g. volcanic eruptions, anthropogenic aerosols) or evaluation of model errors and statistical significance (e.g. atmosphere-only experiment to identify biases due to coupled mode). Each experiment type is characterized by a set of compulsory requirements and additional recommendations. But even among mandatory requirements, some flexibility remains in their concrete implementation. METAFOR work consisted firstly of encoding CMIP5-defined experiment- and simulation vocabulary as specific CV-XML documents so as to become machine readable. Secondly, it aimed at capturing the characteristics of a simulation that is left to the person configuring the simulation. Thirdly, it proposed a way to tell how the simulation described meets the experiment requirements it is intended to fit; this is ensured by introduction of the “Conformance” concept. In its current state, CV for conformance is quite restricted, asking how experiment requirements are met (if so) as per the mean. Possible choices are “via standard configuration”, “via model modifications”, “via inputs”, “via combination”, “not applicable” or “not conformant”. Its main function is to enforce a conformance check by metadata providers.

The Experiment CV-XML documents containing the specific CMIP5 experiment and simulation CV<sup>3</sup> were fixed once and for all and cannot be modified by the climate modellers; they are ready for ingestion into the CMIP5 Questionnaire (see next section) and conform to the CIM activity package class structure. In this CV, experiments are identified by a label, a title, a description and an associated list of requirements. Taking the pre-industrial control experiment as an example (see Fig. 4), *3.1\_pi-Control* stands for the experiment label; *Pre-Industrial Control: control experiment against which perturbations are compared* for the experiment title and *Pre-Industrial coupled atmosphere/ocean control run. Imposes non-evolving pre-industrial conditions* for the experiment description. In turn, each requirement has a label, a type and a description attached. To continue with the same example: the requirement with label *3.1.bc.CO2\_conc* has *BoundaryCondition* as requirement type and *Prescribed atmospheric concentrations of pre-industrial well mixed gas: Carbon Dioxide* as requirement description.

One or more simulations may support the realization of one particular experiment. Each simulation is identified by a short name, a long name, a description, its DRS member name (“rip” values standing for “realization – initialization

<sup>3</sup>See [http://METAFORclimate.eu/trac/browser/controlled\\_vocabularies/branches/cmip5/Activity](http://METAFORclimate.eu/trac/browser/controlled_vocabularies/branches/cmip5/Activity).



**Fig. 4.** Tree diagram showing information necessary to identify and document an experiment. Example shown is the CMIP5 pre-industrial experiment. The experiment is identified by a label, a title, an associated description and the list of requirements to be fulfilled by the simulations that instantiate this experiment. Each requirement is in its turn identified by a label, a type and a description. Value (text) for these attributes is fixed once and for all by the CMIP5 experiment protocol. Notice that this tree diagram is just illustrative (it is not a CV mindmap).

method – physics” identifier; see Taylor et al., 2011b), the name of the model used, the hardware platform on which it has been executed, the start date, time extent, or end date. Among these attributes only model name and the DRS member name is controlled vocabulary (defined within the CMIP5 experiment protocol, as mentioned above). When an experiment requires ensemble runs, one simulation is in its turn described as composed of one or several simulation members, each one being unambiguously identified by its DRS member name (“rip” value). Ensemble type (with the following possible values: *Experiment Driven*, *Initial Condition*, *Perturbed boundary Conditions*, *Perturbed Physics* or *Mixed*) is an additional attribute important for capturing in a standard way the perturbation applied to the ensemble members. Figure 5 illustrates how these attributes are filled in for an ensemble simulation labelled *decadal1959* that is an instance of the *1.1 decadal experiment*.

## 5 From controlled vocabulary to metadata

### 5.1 Creating instances of CMIP5 metadata

To collect metadata for CMIP5 numerical models, simulations and experiments, METAFOR has constructed what was initially intended to be a “simple questionnaire”. However, it rapidly became clear that a traditional questionnaire based on a linear collection of information would be completely inappropriate for the task, given the amount of information to be collected and given that much of this information would have to be shared and compared, for instance across

two simulation descriptions. Moreover, a simple, linear text based questionnaire would have required a huge effort of “by hand” treatment in order to translate information harvested into CIM-instances that ultimately feed the CMIP5 metadata database (see Sect. 5.2 for details on information workflow). Thus a more complex tool was needed, and clearly that tool had to be based on the controlled vocabularies defined for CMIP5 and described in Sect. 4. The name has remained, but the “CMIP5 Questionnaire” should be thought of as a complex metadata entry tool, reproducing the CIM syntax structure and syntax and able to make links between metadata objects referring each other.

The resulting questionnaire provides support for harvesting all aspects a modeller controls when he or she performs a CMIP5 experiment (see Fig. 6): the model(s) used (including the coupling system), its associated grids, the computational platform it has been run on, the different simulations performed and the experiment they are related to, the input data files and, optionally, the CF standard names of the variables in the file used as a model component input. It allows users to interactively produce CIM metadata documents (see Lawrence et al., 2012, for an explanation of the term “document” in this context) without any knowledge of CIM structures. The CMIP5 Questionnaire has been built using the python Django web framework (<http://www.djangoproject.com/>), deployed at the British Atmospheric Data Centre (BADC) and is available online at <http://q.cmp5.ceda.ac.uk/>.

An illustration of how the Model CV is exploited to build the CMIP5 Questionnaire pages is shown in Fig. 7a. The end result is that the structure of the model component pages in





the questionnaire – in terms of, for example, the hierarchy presented and the order of the parameters asked about – is completely controlled by the originating CV mindmap. This flexibility has, of course, been crucial in the development of the questionnaire.

Figure 7a shows the page corresponding to the *SeaIceThermodynamics* component taken as the example when discussing the Model CV definition process (Fig. 2). The navigation tree on the left provides a hierarchical view of the possible component structure of an earth system model. It strictly reflects the CV structure of the eight realm components as fixed in the mindmaps. The first three frames (from the top of the page) are for generic questions, common to all components (either realm or child): user-defined component names (the component type, here *SeaIceThermodynamics* being fixed) and which grid is used by the current component. The next three frames, zoomed in Fig. 7b, contain questions entirely driven by the CV for that component. For example, the fifth frame that asks a question about the SchemeType (*snow-aging*, *snow-ice* or *Other*) mirrors the *Snow* parameter group.

As explained above, the CMIP5 Questionnaire helps the modellers to describe their model using the CV. The questionnaire is also extensible, however, offering the possibility for the user to define parameter-value attributes for each component, and indeed arbitrary additional component structures. Obviously, such flexibility is not in line with the current main scope of standardization. Nevertheless, we considered it important to allow the user to add information that has not been anticipated by the METAFOR CV. Moreover, additional user inputs can help identifying parts of the CV that will need to be completed or changed in an after-CMIP5 perspective.

The questionnaire also uses the specific CV defined for the simulation descriptions. The way a given simulation meets the CMIP5 requirements of an experiment is described by a so-called “Conformance” (see Sect. 4.2). Conformance can be reached via modifications of model inputs, changes in the model parameters, slight modifications of the code itself, or via a combination of those. A simulation may not even fully conform to its experiment (for instance when the data producer realizes afterwards, when checking the long list of requirements, that his simulation missed one of them). In this later case “not conformant” is the minimal amount of information to provide. Figure 8 illustrates how the conformance of a simulation named *PICTL* to requirements of the *Pre industrial Control* experiment is captured by the questionnaire.

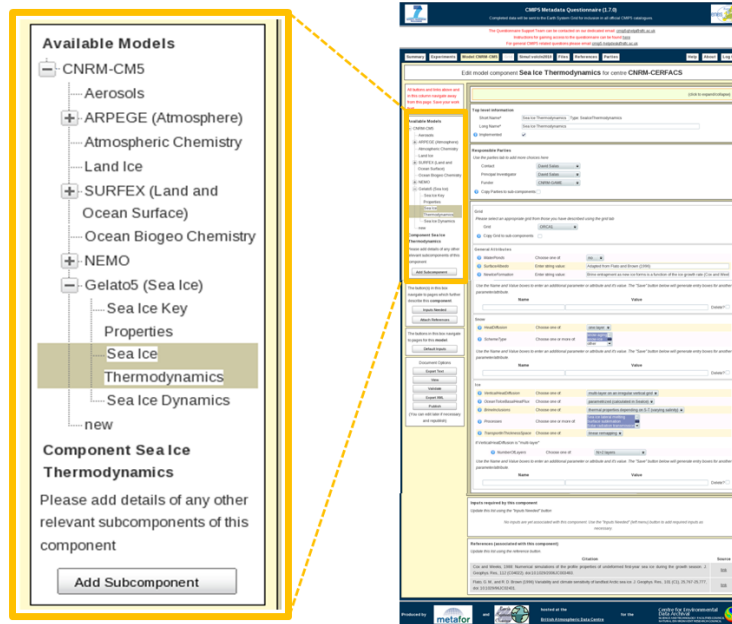
## 5.2 The information pipeline

It is clear that the METAFOR CV has been built with the intent to go beyond simple vocabulary collection usage. Indeed, it is targeted at automatic ingestion by downstream tools (the CMIP5 Questionnaire – discussed in Sect. 5.1) and for inclusion into OWL (Web Ontology Language)

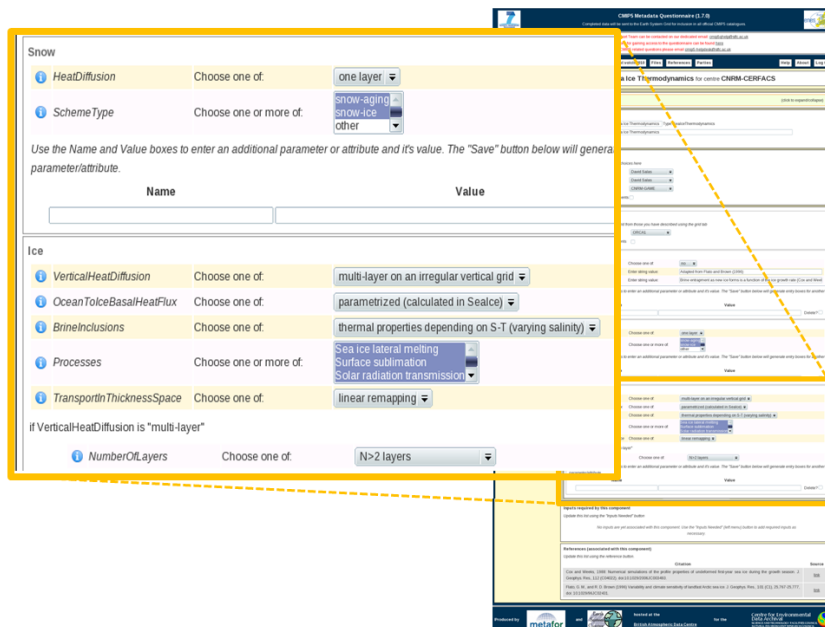
ontologies, e.g. as used in the ESG/CURATOR portal then in use. The Experiment and Simulation CV were fixed by the CMIP5 protocol and are not likely to evolve in the CMIP5 time frame (hence it was created and stored directly in XML without extra tooling). The CV built for model description is, on the other hand, intentionally managed in a different way (i.e. in mindmaps, see Sect. 4), independently from the software tools using them. The objective is to ensure separation of concerns between building and usage so that the semantic database (the Model CV) and the tools using them (the questionnaire) or hosting them (the CIM) can evolve on their own timeline. However, the mindmap format cannot directly feed these downstream tools: format conversion into a machine-readable format was required. To that end, we developed the software to support the information pipeline illustrated in Fig. 9. This tool chain can be found on the METAFOR SVN repository at <http://metaforclimate.eu/trac>.

To satisfy CMIP5 Questionnaire needs, a simple XML-CV structure was defined to encode the Model CV based on the mindmap rules and constraints described earlier. A mindmap validator (top-left grey box in Fig. 9), written in XSLT and invoked by Python, was implemented to check that a specific mindmap (top-right red box in Fig. 9) conforms to the defined encoding rules (see Sect. 4.1). If a feature in the mindmap missed a rule (e.g. an element coded as leaf parameter having a child element) the person responsible for the CV mindmap is asked to make appropriate corrections. Once the validation step is passed, a mindmap translator (top-right grey box in Fig. 9) rewrites the mindmap information into an XML file (middle-right red box in Fig. 9), suitable for ingestion in the questionnaire (middle orange box in Fig. 9). These CV-XML documents are then imported into Django tables and are used to automatically build the questionnaire graphical interface part related to the model description. Once filled in, the questionnaire supports three levels of validation (validate in the middle-left, Fig. 9): (i) the CV constraints are directly enforced while filling in the component description (e.g. a page cannot be saved if text is provided where a numeric value is expected); (ii) when documents are exported as XML files, a validation against the CIM XSD (<http://www.w3.org/TR/xmlschema11-1/>) is automatically enforced; (iii) a Schematron (<http://www.schematron.com>)-based validation is performed to check deeper levels of coherency between the different parameters. The schematron validation ensures that parameters relevant only for a given condition are only filled when this condition is met. For example, in the description of the vertical grid, a *SurfaceReference* is asked only if the *VerticalCoordinateType* is mass-based. The pages of the questionnaire being non-dynamic, the schematron function is to check coherency between responses given by the person filling in the questionnaire.

To ensure usage by the ESGF gateway interfaces and faceted browsing (Williams et al., 2011), a tool was developed to convert the METAFOR Model CV into an OWL ontology (bottom-right red box in Fig. 9). This ontology was



**Fig. 7a.** How the model pages of the CMIP5 Questionnaire automatically inherit from the CV mindmap organization. Components' hierarchy (realm and child components) determines the model navigation tree (left column, enhanced in the zoom).



**Fig. 7b.** Continuing Fig. 7a. Each model component mindmap provides the content of the corresponding questionnaire page, and parameter groups in the component mindmap determine the frames in the page; mindmap leaf parameters define the requested information lines in the frames; list of possible CV values for a given parameter forms the content of drop-down menus (enhanced in the zoom).

also used to guide the mapping tool which allowed the conversion of CIM documents into gateway RDF (Resource Description Framework, <http://www.w3.org/TR/rdf-mt/>) triple stores (Lawrence et al., 2012). The conversion of Model CV into OWL was then the decisive step for the final adoption of METAFOR CV as CMIP5 metadata for models and

simulations (bottom-right yellow box in Fig. 9). Finally, CIM-compliant documents, conforming to the CMIP5 DRS, were broadcast as “atom feeds”, and the corresponding meta-data were ready to be included in the CMIP5 metadata catalogue deployed on the ESG portal.

**Fig. 8.** Illustration of the “Conformance” concept in the case of a *PICTL* simulation performed with CNRM-CM5 model in the framework of 3.1 *piControl* CMIP5 experiment. For the three requirements shown, the conformance is ensured *via inputs*, which means that the input files used contain the forcings requested. Note that a free-text area to enter additional details is always provided.

Since the original tool chain was developed, a new tool chain has been deployed. The CIM-compliant XML documents are now stored in a database, and extracted and displayed in client portals via JavaScript code which loads the documents across the net, and then displays them.

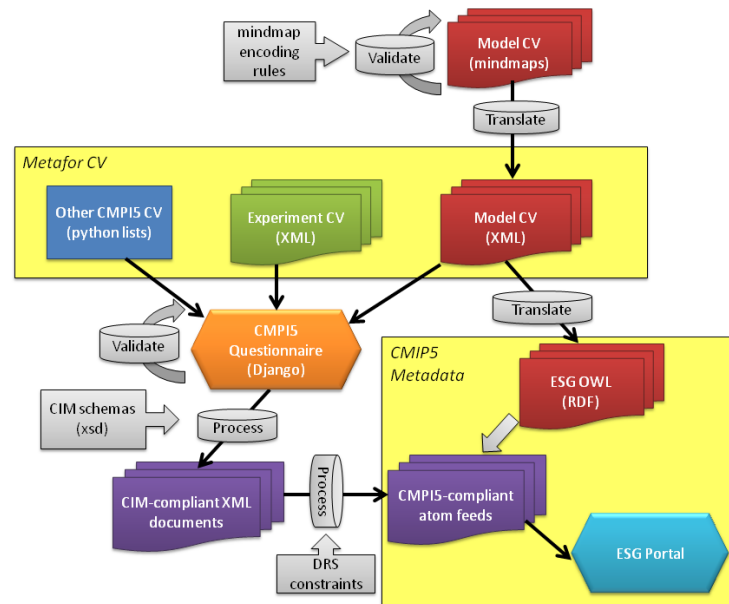
## 6 Summary and further work

CMIP5 was conducted by 20 modelling groups that produced about 90 000 years of simulation for a total volume of several petabytes. A CMIP5 climate data user is faced with a large amount and large diversity of data sets archived in CMIP5 data-node centres. In this context, the METAFOR mission was to provide a metadata system to support data preservation, data reuse (both in time and by different research communities), data readability and discovery, and to guarantee the data quality (or conformity). Until now, such an integrated metadata system for climate modelling was missing. The controlled vocabulary for model and simulation should be considered as necessary raw material for such a system.

This paper introduced the controlled vocabulary developed both for generic description of earth system models and as input for the tool developed to collect this description for CMIP5 models and simulations (the CMIP5 Questionnaire). The mindmap technology used facilitated the CV development, ensuring a wide engagement of the scientific community in this process, hiding away the complexity of the underlying ontological concepts (the CIM). The metadata

pipeline, which starts from the mindmaps, serves both the metadata entry tool (the CMIP5 Questionnaire) and metadata catalogues such as the ESGF gateways. The cornerstone of the METAFOR CV has indisputably been the engagement of a large number of modellers from the climate community since the early stage of the CV elaboration process. The CV collection produced at the end of the METAFOR project gathers thousands of terms for which hierarchical arrangement is equally important as the terms themselves. Even if it can be improved further, METAFOR CV is the first one to address the whole climate modelling chain. Available in CMIP5 metadata catalogues and supporting data discovery tools, the hope is to provide essential services to climate data users.

There are two significant pieces of work yet to be done before the CV can be easily governed and maintained. Firstly, a conversion tool taking the CV XML back to the mindmap format would support the ability to convert between all CV formats. This tool would allow use of the CV XML as the primary preservation and governance artefact, generating mindmaps from those XML instances for websites and human-mediated discussions for example. Secondly, we need to formalize an http interface for the CV following appropriate standards (see Leadbetter et al., 2011). Secondly, the maintenance and governance of the controlled vocabulary and of the associated metadata pipeline needs addressing. Gathering feedback from the questionnaire users and finding ways to benefit from this feedback to make the CV evolve



**Fig. 9.** Key components of the CV and information pipeline from METAFOR CV (top yellow box) to CMIP5 metadata (bottom-right yellow box).

should also be strongly considered. The intended focus is the set-up of a real standard, which requires a governance committee to emerge (as planned in the framework of EU-FP7 ENES2 project). For now, the METAFOR work is extended within the UK JISC-funded PIMMS project (Portable Infrastructure for the METAFOR Metadata System).

Although populating the CMIP5 Questionnaire was not mandatory, but highly recommended by the CMIP5 panel (i.e. not blocking for CMIP5 model outputs publication), about 70 % of the modelling groups contributing to CMIP5 provided metadata through the questionnaire and more than a thousand CIM documents are stored in the CMIP5 documentation repository by this way. An overall 78 % of the published documents are attached to the description of experiments and simulations and 12 % describe the models and their grids. Only 2 groups provided a description for their model but not for the simulations they performed, and 8 groups did not provide any metadata at all. Further diagnostics to measure quality and completion rate of the metadata documents would be advisable. A key point for a wider acceptance of the metadata harvesting procedure (the questionnaire and the underlying CIM) is certainly to limit the effort asked to metadata providers. With the CMIP5 Questionnaire, the effort required from the provider was indisputably too strong. The logic of the information flow and connections between formal concepts was viewed as somewhat complex. Lessons are to be learned from the METAFOR experience in the context of CMIP5 that should be reinvested in future projects.

While the application of CMIP5 has dominated most of the development thus far, next generations of the questionnaire are currently being developed by the ES-DOC community (Earth System Documentation, <http://earthsystemcog.org/projects/es-doc-models/>). Initiated during the METAFOR project, specific CV is being developed to describe the models and simulations used in the ENSEMBLES EU project (<http://www.ensembles-eu.org>). The US NCPP project (National Climate Predictions and Projections, <http://earthsystemcog.org/projects/ncpp/>) and EU EURO-CORDEX (Coordinated Downscaling Experiment – European Domain, <http://www.euro-cordex.net>) are also agreeing on statistical and dynamical downscaling CV for regional climate studies. Finally, one can expect that the METAFOR CV for global climate models will be reused in upcoming or recent EU FP7 initiatives dedicated to climate services as the SPECS project (Seasonal-to-decadal climate Prediction for the improvement of European Climate).

*Acknowledgements.* METAFOR was funded by the EU 7th Framework Programme as an e-infrastructure (project #211753). The support of the EU FP7 IS-ENES (project #228203) is also acknowledged. This work benefited significantly from the engagement of other METAFOR members and colleagues from the US Earth System Curator project. We also appreciated guidance from the METAFOR advisory committee, in particular Wilco Hazeleger and Karl Taylor.

Edited by: M. Kawamiya



The publication of this article is financed by CNRS-INSU.

## References

- Balaji Institute: Gridspec – A standard for the description of grids used in Earth System models, available at: <http://www.gfdl.noaa.gov/~vb/gridstd/gridstd.html> (last access: 18 March 2014), GFDL 2007.
- Boucher, O. and Pham, M.: History of sulfate aerosol radiative forcings, *Geophys. Res. Lett.*, 29, 22.1–22.4, doi:10.1029/2001GL014048, 2002.
- Callaghan, S. A., Treshansky, A., Moine, M.-P., Guilyardi, E., Alias, A., Balaji, V., Bojariu, R., Cofiño, A. S., Denvil, S., Elkington, M., Ford, R., Kolaninski, M., Lautenschlager, M., Lawrence, B. N., Steenman-Clark, L., and Valcke, S.: The METAFOR project: preserving data through metadata standards for climate models and simulations, in: INTL-DPIF '10 Proceedings of the 1st International Digital Preservation Interoperability Framework Symposium, article No. 6, doi:10.1145/2039263.2039269, 2010.
- Cariolle, D. and Déqué, M.: Southern hemisphere medium-scale waves and total ozone disturbances in a spectral general circulation model, *J. Geophys. Res.*, 91, 10825–10846, 1986.
- Cariolle, D., Lasserre-Bigory, A., Royer, J.-F., and Geleyn, J.-F.: A general circulation model simulation of the springtime Antarctic ozone decrease and its impact on mid-latitudes, *J. Geophys. Res. Atmos.*, 95, 1883–1898, 1990.
- Dunlap, R., Mark, L., Rugaber, S., Balaji, V., Chastang, J., Cinquini, L., DeLuca, C., Middleton, D., and Murphy, S.: Earth system curator: metadata infrastructure for climate modeling, *Earth Sci. Inf.*, 1, 131–149, doi:10.1007/s12145-008-0016-1, 2008.
- Ford, R. W. and Riley, G. D.: The Bespoke Framework Generator, in: *Earth System Modelling*, Vol. 3, Coupling Software and Strategies, Series: SpringerBriefs in Earth System Sciences, ISBN 978-3-642-23359-3, 2011.
- Guilyardi, E., Balaji, V., Callaghan, S., DeLuca, C., Devine, G., Denvil, S., Ford, R., Pascoe, C., Lautenschlager, M., Lawrence, B. N., Steenman-Clark, L., and Valcke, S.: The CMIP5 model and simulation documentation: a new standard for climate modeling metadata, *CLIVAR Exchanges*, 16, 42–46, 2011.
- Lawrence, B. N., Balaji, V., Bentley, P., Callaghan, S., DeLuca, C., Denvil, S., Devine, G., Elkington, M., Ford, R. W., Guilyardi, E., Lautenschlager, M., Morgan, M., Moine, M.-P., Murphy, S., Pascoe, C., Ramthun, H., Slavin, P., Steenman-Clark, L., Tossaint, F., Treshansky, A., and Valcke, S.: Describing Earth system simulations with the Metafor CIM, *Geosci. Model Dev.*, 5, 1493–1500, doi:10.5194/gmd-5-1493-2012, 2012.
- Leadbetter, A., Clements, O., and Lowry, R.: Emerging standards in vocabulary server access methods, *Geophys. Res. Abstr.*, EGU2011-A-2143, EGU General Assembly 2011, Vienna, Austria, 2011.
- Levitus, S.: Climatological atlas of the world's oceans, NOAA Professional Paper 13, 173 pp., available at: [ftp://ftp.nodc.noaa.gov/pub/data.nodc/woa/PUBLICATIONS/levitus\\_atlas\\_1982.pdf](ftp://ftp.nodc.noaa.gov/pub/data.nodc/woa/PUBLICATIONS/levitus_atlas_1982.pdf) (last access: 18 March 2014), 1982.
- Redler, R., Valcke, S., and Ritzdorf, H.: OASIS4 – a coupling software for next generation earth system modelling, *Geosci. Model Dev.*, 3, 87–104, doi:10.5194/gmd-3-87-2010, 2010.
- Taylor, K. E. and Doutriaux, C.: CMIP5 Model Output Requirements: File Contents and Format, Data Structure and Metadata, available at: [http://cmip-pcmdi.llnl.gov/cmip5/docs/CMIP5\\_output\\_metadata\\_requirements.pdf](http://cmip-pcmdi.llnl.gov/cmip5/docs/CMIP5_output_metadata_requirements.pdf) and [http://pcmdi-cmip.llnl.gov/cmip5/docs/standard\\_output.pdf](http://pcmdi-cmip.llnl.gov/cmip5/docs/standard_output.pdf) (last access: 18 March 2014), 2010.
- Taylor, K. E., Stouffer, R. J., and Meehl, G. A.: An Overview of CMIP5 and the Experiment Design, *B. Am. Meteorol. Soc.*, 93, 485–498, doi:10.1175/BAMS-D-11-00094.1, 2011a.
- Taylor, K. E., Balaji, V., Hankin, S., Juckes, M., Lawrence, B. N., and Pascoe, S.: CMIP5 Data Reference Syntax (DRS) and Controlled Vocabularies, available at: [http://pcmdi-cmip.llnl.gov/cmip5/docs/cmip5\\_data\\_reference\\_syntax.pdf](http://pcmdi-cmip.llnl.gov/cmip5/docs/cmip5_data_reference_syntax.pdf) (last access: 18 March 2014), 2011b.
- Valcke, S., Balaji, V., Craig, A., DeLuca, C., Dunlap, R., Ford, R. W., Jacob, R., Larson, J., O'Kuinghtons, R., Riley, G. D., and Vertenstein, M.: Coupling technologies for Earth System Modelling, *Geosci. Model Dev.*, 5, 1589–1596, doi:10.5194/gmd-5-1589-2012, 2012.
- Voltaire, A., Sanchez-Gomez, E., Salas y Mélia, D., Decharme, B., Cassou, C., Sénési, S., Valcke, S., Beau, I., Alias, A., Chevalier, M., Déqué, M., Deshayes, J., Douville, H., Fernandez, E., Madec, G., Maisonnave, E., Moine, M.-P., Planton, S., Saint-Martin, D., Szopa, S., Tyteca, S., Alkama, R., Belamari, S., Braun, A., Coquart, L., and Chauvin, F.: The CNRM- CM5.1 global climate model: Description and basic evaluation, *Clim. Dynam.*, 40, 2091–2121, doi:10.1007/s00382-011-1259-y, 2011.
- Williams, D. N., Lawrence, B. N., Lautenschlager, M., Middleton, D., and Balaji, V.: The Earth System Grid Federation: Delivering globally accessible petascale data for CMIP5, in: *Proceedings of the 32nd Asia-Pacific Advanced Network Meeting*, 121–130, New Delhi, doi:10.7125/APAN.32.15, 2011.

## Appendix A

### List of climate scientists involved in the METAFOR consultation process

The METAFOR project members would like to express their sincere thanks to all the climate scientists who contributed in a significant way to the METAFOR controlled vocabulary elaboration process, sharing their knowledge without restriction and providing excellent guidance and recommendations (in alphabetical order, Table A1 below).

**Table A1.** List of climate scientists who contributed to the METAFOR controlled vocabulary, whether during face to face interviews, phone calls or e-mail exchanges.

---

Abrahams, Luke, UKCA, UK
Balaji, V., GFDL, USA
Boone, Aaron, CNRM, France
Bopp, Laurent, LSCE-IPSL, France
Braesicke, Peter, UKCA, UK
Bruehl, Christoph, UKCA, UK
Buja, Lawrence, NCAR, USA
Decharme, Bertrand CNRM, France
Déqué, Michel CNRM, France
Elkington, Mark, MetOffice, UK
Fichefet, Thierry, UCL-LLN, Belgium
Gibelin, Anne-Laure, CNRM, France
Goosse, Hugues, UCL-LLN, Belgium
Griffies, Stephen, GFDL, USA
Guilyardi, Eric, LOCEAN-IPSL, France
Hagemann, Stefan, MPI, Germany
Horowitz, Larry, GFDL, USA
Hourdin, Frédéric, LMD-IPSL, France
Kageyama, Masa, LSCE-IPSL, France
Khodry, Myriam, IPSL, France
Krinner, Gerhard, LGGE, France
Lawrence, Bryan, NCAS-BADC, UK
Madec, Gurvan, LOCEAN-IPSL, France
Malyshev, Sergey, GFDL, USA
Mann, Graham, Univ. of Leeds, UK
Marti, Olivier, LSCE-IPSL, France
Peuch, Vincent-Henri, CNRM, France
Polcher, Jan, LMD-IPSL, France
Ritz, Catherine, LGGE, France
Salas Y. Melia, David, CNRM, France
Slawitch, Ross, Univ. Maryland, USA
Strand, Gary, NCAR, USA
Van Velthoven, Peter, KNMI, the Netherlands
Vancoppenolle Martin, UCL-LLN, Belgium
Wyman Bruce, GFDL, USA

---

## Appendix B

### CMIP3 text-based questionnaire

Model Information of Potential Use to the IPCC Lead Authors and the AR4.

CNRM-CM3 (version used for IPCC AR4)

2 August 2005

Model identity:

A. Institution, sponsoring agency, country: Centre National de Recherches Météorologiques, Météo France, France

B. Model name (and names of component atmospheric, ocean, sea ice, etc. models): CNRM-CM3

Atmosphere: ARPEGE-Climat version 3

Ocean: OPA 8.1

Sea ice: GELATO 2

C. Vintage (i.e. year that model version was first used in a published application): 2004

D. General published references and web pages:

[http://www.cnrm.meteo.fr/scenario2004/references\\_eng.html](http://www.cnrm.meteo.fr/scenario2004/references_eng.html)

E. References that document changes over the last 5 years (i.e. since the IPCC TAR) in the coupled model or its components. We are specifically looking for references that document changes in some aspect(s) of model performance.

– descriptions of previous versions of the ARPEGE-Climat model can be found in the following publications:

– Déqué et al. (1994),

– Déqué and Piedelièvre (1995),

– Royer et al. (2002).

F. IPCC model version's global climate sensitivity (KW-1 m<sup>2</sup>) to increase in CO<sub>2</sub> and how it was determined (slab ocean expt., transient expt–Gregory method, ±2 K Cess expt., etc.): not yet available

G. Contacts (name and email addresses), as appropriate, for:

1. coupled model: David Salas y Melia, david.salas@meteo.fr

2. atmosphere : Michel Déqué, michel.deque@meteo.fr

3. ocean : David Salas y Melia, david.salas@meteo.fr

4. sea ice: David Salas y Melia, david.salas@meteo.fr
5. land surface: Hervé Douville, herve.douville@meteo.fr
6. vegetation: Hervé Douville, herve.douville@meteo.fr
7. other?

Besides atmosphere, ocean, sea ice, and prescription of land/vegetated surface, what can be included (interactively) and was it active in the model version that produced output stored in the PCMDI database?

#### A. Atmospheric chemistry?

- Ozone transport with simplified chemistry as described in Cariolle and Déqué (1986) and Cariolle et al. (1990).

#### B. Interactive biogeochemistry?

- no

#### C. What aerosols and are indirect effects modelled?

- The distributions of marine, desertic, urban aerosols, sulfate aerosols are specified. Marine and desertic aerosols are constant in all experiments. Urban aerosols vary according to estimates between 1860 and 2000. Sulfate aerosols are specified in all experiments according to Boucher and Pham (2002) data, see <http://www-loa.univ-lille1.fr/boucher/sres/> for more details. Note that only the direct effect of anthropogenic sulfate aerosols was taken into account.

#### D. Dynamic vegetation?

- no

#### E. Ice sheets?

- fixed

[...]

Component model characteristics (of current IPCC model version):

#### A. Atmosphere

1. Resolution: triangular truncation T63 with “linear” reduced Gaussian grid equivalent to T42 quadratic grid
2. Numerical scheme/grid (advective and time-stepping schemes; model top; vertical coordinate and number of layers above 200 hPa and below 850 hPa):

- semi-Lagrangian semi-implicit time integration with 30 min time step, 3-hour time step for radiative transfer;
  - top layer 0.05 hPa, progressive hybrid sigma-pressure vertical coordinate with 45 layers, 23 layers above 200 hPa, usually 7 layers below 850 hPa (less in regions of high orography)
3. List of prognostic variables (be sure to include, as appropriate, liquid water, chemical species, ice, etc.). Model output variable names are not needed, just a generic descriptive name (e.g. temperature, northward and eastward wind components, etc.)
    - temperature, northward and eastward wind components, specific humidity, ozone concentration, surface pressure
  4. Name, terse descriptions, and references (journal articles, web pages) for all major parameterizations. Include, as appropriate, descriptions of:
    - a. Clouds:
      - statistical cloud scheme for stratiform clouds based on Ricard and Royer (1993). Convective cloud cover based on the mass-flux transport
    - b. Convection
      - mass-flux convective scheme with Kuo-type closure based on Bougeault (1985) boundary layer based on Louis et al. (1982) with modifications by Mascart et al. (1995). SW, LW radiation based on Fouquart and Morcrette parameterizations implemented in a former version of the ECMWF model (Morcrette JJ, 1990; Morcrette JJ, 1991)
    - c. any special handling of wind and temperature at top of model:
      - relaxation of temperature, linear (Rayleigh) friction for wind
- Simulation details (report separately for each IPCC simulation contributed to database at PCMDI)  
Picntrl/Run\_1

This pre-industrial control simulation was initialized from a coupled simulation of a previous version of CNRM coupled model that initialized an ocean at rest with temperature and salinity profiles specified from Levitus (1982) climatology, integrated for 30 years with a relaxation of surface temperature to the monthly mean Reynolds climatology for 1950. The CNRM-CM3 version was then integrated for 70 years with pre-industrial 1860 greenhouse gases concentrations as a spin-up. After this spin-up period, results were stored from nominal years 1930 to 2429.