

# The Virtual Climate Data Server (vCDS): An iRODS-Based Data Management Software Appliance Supporting Climate Data Services and Virtualization-as-a-Service in the NASA Center for Climate Simulation

*John L. Schnase<sup>1</sup>, Glenn S. Tamkin<sup>2,3</sup>, W. David Ripley, III<sup>2,3</sup>,  
Savannah Strong<sup>2,3</sup>, Roger Gill<sup>2,4</sup>, and Daniel Q. Duffy<sup>2</sup>*

<sup>1</sup> Office of Computational and Information Science and Technology

<sup>2</sup> NASA Center for Climate Simulation (NCCS)

<sup>3</sup> Computer Science Corporation (CSC), <sup>4</sup> Innovim, LLC  
NASA Goddard Space Flight Center  
Greenbelt, MD 20771

## Abstract

Scientific data services are becoming an important part of the NASA Center for Climate Simulation's mission. Our technological response to this expanding role is built around the concept of a Virtual Climate Data Server (vCDS), repetitive provisioning, image-based deployment and distribution, and virtualization-as-a-service. The vCDS is an iRODS-based data server specialized to the needs of a particular data-centric application. We use RPM scripts to build vCDS images in our local computing environment, our local Virtual Machine Environment, NASA's Nebula Cloud Services, and Amazon's Elastic Compute Cloud. Once provisioned into one or more of these virtualized resource classes, vCDSs can use iRODS's federation capabilities to create an integrated ecosystem of managed collections that is scalable and adaptable to changing resource requirements. This approach enables platform- or software-as-a-service deployment of vCDS and allows the NCCS to offer virtualization-as-a-service: a capacity to respond in an agile way to new customer requests for data services.

**Index Keyword Terms**—iRODS, climate data services, cloud computing, software appliance, virtualization

## 1. Introduction

The NASA Center for Climate Simulation (NCCS) provides large-scale compute engines, analytics, data sharing, long-term storage, networking, and other high-end computing services designed to meet the specialized needs of the Earth science communities. By doing so, NCCS brings NASA observational and model data products to climate research carried out by a wide range of national and international organizations [1].

Last year, we examined the potential of iRODS, the Integrated Rule-Oriented Data System, as a means of integrating, archiving, and delivering scientific data to the communities we serve. We built a testbed collection of independent iRODS data systems comprising observational and simulation data and used the testbed to learn about iRODS and understand how the technology might further our mission. We came away from that exercise believing that iRODS could provide a useful platform upon which to build a collection of scientific data services tailored to the needs of our customers [2].

This year, we have worked to build an operational iRODS capability for the NCCS. The result is a product, architecture, and approach we refer to as the Virtual Climate Data Server (vCDS), a software appliance specialized to the needs of climate data collections management. In the following sections, we describe our experiences with vCDS, including motivation and rationale for the approach, implementation details, and future plans regarding scientific data services in the NCCS.

## 2. Background

Data services and data publication are becoming increasingly important aspects of NCCS's mission. For example, our two major customers, NASA's Global Modeling and Assimilation Office (GMAO) and the Goddard Institute for Space Studies (GISS) are contributing products to the Intergovernmental Panel on Climate Change (IPCC) Fifth Assessment Report (AR5) [3]. Data products for the IPCC AR5 assessment need to be published to the broader community through the Earth System Grid (ESG) [4]. GMAO computes the Modern Era Retrospective-Analysis for Research and Applications (MERRA) data set in the NCCS, which in turn we convey (publish) to the Goddard Earth System Data Infor-

mation and Information Services Center (GES DISC). And the NCCS will be computing the Level 4 root-level soil moisture product for NASA’s Soil Moisture Active Passion (SMAP) mission when it launches in 2014.

As suggested by these examples, the diversity of our customer base as well as the diversity of the data itself are increasing. Our customers now include individual scientist, labs, research projects, flight missions, and even non-traditional, private-sector consumers of climate simulation products such as the insurance/re-insurance industry. The datasets involved may be products generated by a General Circulation Model (GCM), observational data, reanalysis data, or specialized derived products requiring the combination of simulation and observational data. Depending on the circumstances, management of the data may require short-term storage, long-term archival preservation, or some type of hierarchical staging to accommodate interactive visualization or use by an application.

### 3. vCDS Concept and Rationale

Our notion of a Virtual Climate Data Server has grown out of a simple use case that we developed to capture the essence of this new data challenge.

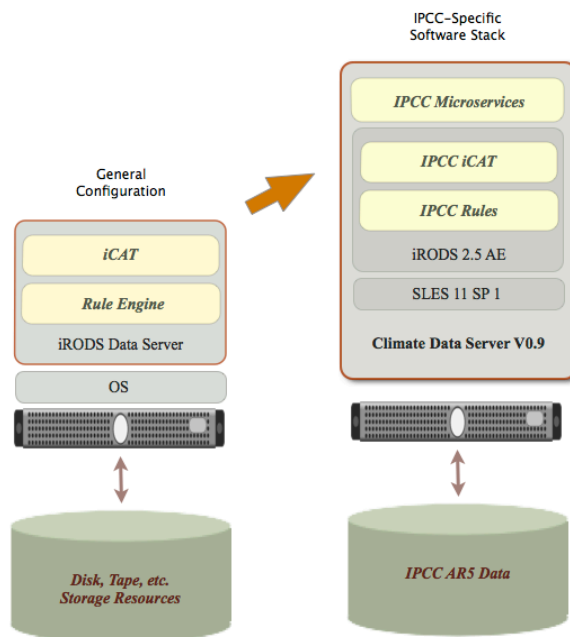
*A customer approaches the NCCS with a new dataset they want us to manage: What technology is needed to quickly meet that customer’s requirement under the following constraints:*

- *The solution should be: simple, fast, and affordable;*
- *provide core capabilities to get started, but extensible to accommodate future needs;*
- *be flexible, with the ability to use, optimize, and change deployment configurations in response to resource availability;*
- *allow the new dataset to be integrated into an existing data collection; and*
- *come with a help desk and user support?*

The answer we returned to repeatedly is that we need a data server software appliance specialized to the needs of a managed collection of climate-related scientific data. Furthermore, in order to be agile and responsive, this appliance needed to be able to use the flexible resource allocation capabilities afforded by cloud computing. To allow for ease of collections integration and support the full information lifecycle requirements of a scientific archive, it should be built around the iRODS technology. Hence, our notion of an iRODS-based Virtual Climate Data Server whose core functionality and suite of ancillary utilities would support our expanding climate data services mission.

### 4. vCDS Architecture

The basic configuration of an iRODS data server consists of a specific version of iRODS installed on a particular operating system running on particular hardware. Moving toward a virtual appliance model has been a two-step process in which we (1) encapsulate the operating system and iRODS as a virtual machine image, then (2) specialize that image with functionality required for managing climate data. Our approach to specialization has been to build general-purpose scientific “kits” – such as NetCDF, HDF, and GeoTIF – that sit in the vertical stack above iRODS and below application-specific climate kits such as IPCC, MERRA, and SMAP.



Our initial focus has been building a vCDS to manage IPCC AR5 NetCDF data. Details about these components are provided below, but in summary, the core elements include the following:

- *Application-specific microservices* — Canonical archive operations, particularly the mechanisms required to ingest Open Archive Information System (OAIS)-compliant Submission Information Package (SIP) metadata for IPCC NetCDF objects.
- *Application-specific metadata* — OAIS-compliant constitutive (application-independent) Representation Information (RI) and Preservation Description Information (PDI) metadata for IPCC NetCDF objects.
- *Application-specific rules* — IPCC NetCDF triggers and workflows.

- *A specific release of iRODS* — In the current version we have used iRODS 2.5 that has been augmented with what we refer to as Administrative Extensions (AE).
- *A specific operating system* — In our case, SLES 11 SP 1.

At the time of this writing we have built vCDS Version 0.9. Collectively, we refer to the functionality associated with vCDS as a the vCDS V0.9 “product suite.” It includes (1) NetCDF/IPCC Kits, (2) Administrative Extensions, utilities that enable (3) Repetitive Provisioning, and mechanisms for (4) Deployment and Distribution.

The technology readiness level (TRL) of V0.9 is approximately 7, meaning we have completed system prototyping and demonstration in an operational environment and the system is at or near scale of an operational system, with most functions available for demonstration and test [5]. The software used in the vCDS V0.9 stack is shown in the table below.

Name	Version	Notes
iRODS	2.5	Core iRODS installation. Includes i-commands.
Extrods	1.1.0.1-beta	Officially provided iRODS web UI.
PHP	5.2.14	Required for iRODS web UI.
Apache web server	2.2.10	Required to serve iRODS web UI.
FUSE library	2.7.2	Base FUSE library required for iRODS FUSE interface.
Postgresql	8.3.14	Required RDBMS for iCAT.
UnixODBC	2.2.12	Required for iRODS communication to iCAT.
PyRods	2.5	Community provided Python wrapper for iRODS libraries.
EmbedPython	2.5	Community provided iRODS extension that allows for Python based microservice development.
Python	2.6	Core Python environment, needed for PyRods and EmbedPython.
iRODS-NCCS	0.7	Custom Python based microservices for data handling.
iRODS-web	0.7	Java application for viewing iRODS audit history and usage statistics.
Apache Tomcat	7.0.14	Java Servlet container that serves iRODS-web-stats application.
PyGreSQL	4.0	Python postgres driver
JDBC	2.5.5	Java database connectivity required for NetCDF
Java Runtime	1.6.0_24	Java runtime, required for iRODS-web-stats.
Ncdump-hdf	4.2.5	Library for interacting with HDF files, required for iRODS-nccs
SLES	11,sp1	Base OS.

#### 4.1. NetCDF/IPCC Kits

The core functionality of vCDS V0.9 resides in the NetCDS and IPCC kits that contain the iRODS microservices, rules, configuration settings, and software utilities required to implement the system’s canonical operations. These functions include the following:

- Basic system-level operations of an archive: Create, Read, Update, and Delete (CRUD);
- Rules to identify IPCC NetCDF files;
- Microservices to manage OAIS Information Objects:
  - Submission Information Packages (SIPs)
  - Archive Information Packages (AIPs)
  - Distribution Information Packages (DIPs);
- The iRODS iCAT along with optimizations to manage and view metadata.

These capabilities are central to managing the producer/consumer relationships of an archive. Managing metadata separate from the storage objects themselves is key to this since doing so enables discovery, long term curation, reuse of the data, and use of the data for unintended purposes: it is what distinguishes an archive or managed collection of scientific data from the typical bit storage functions of a filesystem.

To demonstrate how we have approached managing metadata in vCDS V0.9, we first describe the internal metadata structure of the NetCDF file format and how it has been specialized by the IPCC for use in their data products. Then, we show how we externalize the embedded NetCDF metadata in a way that makes the vCDS system OAIS compliant.

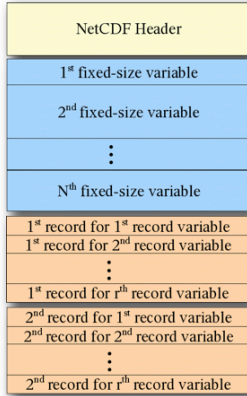
##### 4.1.1 NetCDF Metadata

IPCC climate model outputs are stored as NetCDF files. NetCDF (Network Common Data Form) is a set of software libraries and self-describing, machine-independent data formats that support the creation, access, and sharing of array-oriented scientific data. The project homepage is hosted by the Unidata program at the University Corporation for Atmospheric Research (UCAR). They are also the chief source of NetCDF software, standards development, updates, etc. The format is an open, international standard of the Open Geospatial Consortium [6].

The figure below provides an example of the metadata typically embedded in an IPCC file. The NetCDF header contains general information about the accompanying data as well as specific information required by a NetCDF-aware application to index into the accompanying records to use the files’s data. The organization of this embedded header information is stipulated by the research community through the Climate Model Intercomparison Project (CMIP5) Data Reference Syntax (DRS) and Controlled Vocabularies specification [7].

### IPCC NetCDF Header Metadata

```
netcdf cSoil_Lmon_GISS-E2-
R_piControl_r11p1_398101-400512 {
dimensions:
time = UNLIMITED ; // (300 currently)
lat = 90 ;
lon = 144 ;
bnds = 2 ;
variables:
double time(time) ;
time.bounds = "time_bnds" ;
time.units = "days since 3981-1-1" ;
time.calendar = "365_day" ;
time.axis = "T" ;
time.long_name = "time" ;
time.standard_name = "time" ;
double time_bnds(time, bnds) ;
double lat(lat) ;
lat.bounds = "lat_bnds" ;
lat.units = "degrees_north" ;
lat.axis = "Y" ;
lat.long_name = "latitude" ;
lat.standard_name = "latitude" ;
double lat_bnds(lat, bnds) ;
double lon(lon) ;
lon.bounds = "lon_bnds" ;
lon.units = "degrees_east" ;
lon.axis = "X" ;
lon.long_name = "longitude" ;
lon.standard_name = "longitude" ;
double lon_bnds(lon, bnds) ;
float cSoil(time, lat, lon) ;
cSoil.standard_name = "soil_carbon_content" ;
cSoil.long_name = "Carbon Mass in Soil Pool" ;
cSoil.units = "kg m-2" ;
cSoil.original_name = "dummy" ;
```



```
cSoil.cell_methods = "time: mean area:
mean where land" ;
cSoil.cell_measures = "area: areacella" ;
cSoil.history = "2011-03-21T20:30:30Z
altered by CMOR: replaced missing value
flag (-1e+30) with standard missing value
(1e+20).";
cSoil.missing_value = 1.e+20f ;
cSoil.FillValue = 1.e+20f ;
cSoil.associated_files = "baseURL:
http://cmip-pcmdi.llnl.gov/
CMIP5/dataLocation gridspecFile:
```

```
gridspec_land_fx_GISS-
E2R_piControl_r0i0p0.nc areacella:
areacella_fx_GISS-E2R_
piControl_r0i0p0.nc";

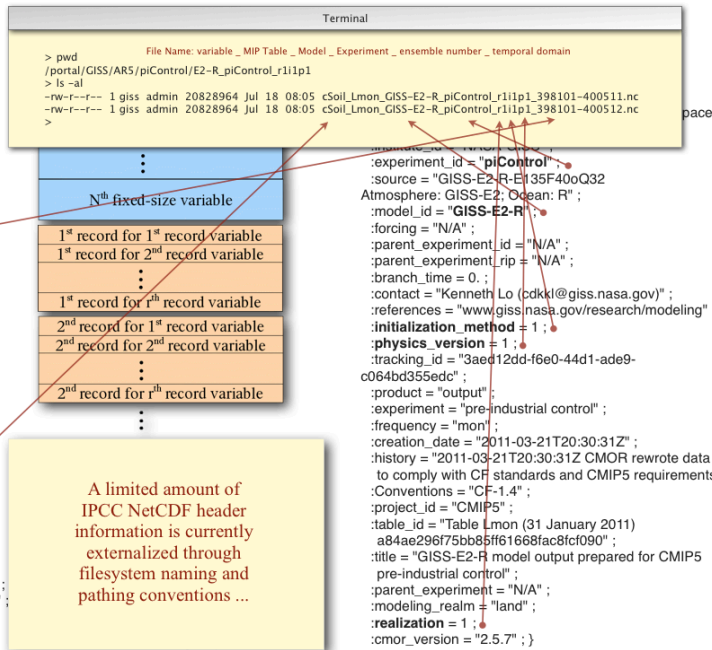
// global attributes:
institution = "NASA/GISS (Goddard Institute for Space
Studies) New York, NY" ;
institute_id = "NASA-GISS" ;
experiment_id = "piControl" ;
source = "GISS-E2-R-E135F40oQ32
Atmosphere: GISS-E2; Ocean: R" ;
model_id = "GISS-E2-R" ;
forcing = "N/A" ;
parent_experiment_id = "N/A" ;
parent_experiment_rip = "N/A" ;
branch_time = 0. ;
contact = "Kenneth Lo (cdkl@giiss.nasa.gov)" ;
references = "www.giss.nasa.gov/research/modeling" ;
initialization_method = 1 ;
physics_version = 1 ;
tracking_id = "3aed12dd-f6e0-44d1-ade9-
c064bd355edc" ;
product = "output" ;
experiment = "pre-industrial control" ;
frequency = "mon" ;
creation_date = "2011-03-21T20:30:31Z" ;
history = "2011-03-21T20:30:31Z CMOR rewrite data
to comply with CF standards and CMIP5 requirements." ;
Conventions = "CF-1.4" ;
project_id = "CMIP5" ;
table_id = "Table Lmon (31 January 2011)
a84ae296f75bb85ff61668fac8fcf090" ;
title = "GISS-E2-R model output prepared for CMIP5
pre-industrial control" ;
parent_experiment = "N/A" ;
modeling_realm = "land" ;
realization = 1 ;
cmor_version = "2.5.7" ; }
```

It is generally the case that facilities such as the NCCS externalize a small amount of this embedded metadata through file and path naming conventions as a way of organizing their NetCDF collections. The next two figures demonstrate how this is currently done in the NCCS and how with an iRODS-based vCDS we are able to externalize all of the embedded

metadata, storing it in the iCAT. Doing so makes it possible to search over these encapsulated attributes without having to open individual files in the collection. The primary work of the NetCDF/IPCC kits in VCDS V0.9 is managing the extraction of this core, application- and use-independent embedded metadata.

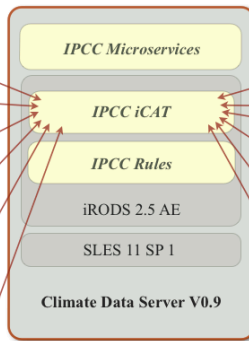
### IPCC NetCDF Header Metadata

```
netcdf cSoil_Lmon_GISS-E2-
R_piControl_r11p1_398101-400512 {
dimensions:
time = UNLIMITED ; // (300 currently)
lat = 90 ;
lon = 144 ;
bnds = 2 ;
variables:
double time(time) ;
time.bounds = "time_bnds" ;
time.units = "days since 3981-1-1" ;
time.calendar = "365_day" ;
time.axis = "T" ;
time.long_name = "time" ;
time.standard_name = "time" ;
double time_bnds(time, bnds) ;
double lat(lat) ;
lat.bounds = "lat_bnds" ;
lat.units = "degrees_north" ;
lat.axis = "Y" ;
lat.long_name = "latitude" ;
lat.standard_name = "latitude" ;
double lat_bnds(lat, bnds) ;
double lon(lon) ;
lon.bounds = "lon_bnds" ;
lon.units = "degrees_east" ;
lon.axis = "X" ;
lon.long_name = "longitude" ;
lon.standard_name = "longitude" ;
double lon_bnds(lon, bnds) ;
float cSoil(time, lat, lon) ;
cSoil.standard_name = "soil_carbon_content" ;
cSoil.long_name = "Carbon Mass in Soil Pool" ;
cSoil.units = "kg m-2" ;
cSoil.original_name = "dummy" ;
```



## IPCC NetCDF Header Metadata

```
netcdf cSoil_Lmon_GISS-E2-
R_piControl_r11p1_398101-400512 {
dimensions:
time = UNLIMITED ; // (300 currently)
lat = 90 ;
lon = 144 ;
bnds = 2 ;
variables:
double time(time) ;
time.bounds = "time_bnds" ;
time.units = "days since 3981-1-1" ;
time.calendar = "365_day" ;
time.axis = "T" ;
time.long_name = "time" ;
time.standard_name = "time" ;
double time_bnds(time, bnds) ;
double lat(lat) ;
lat.bounds = "lat_bnds" ;
lat.units = "degrees_north" ;
lat.axis = "Y" ;
lat.long_name = "latitude" ;
lat.standard_name = "latitude" ;
double lat_bnds(lat, bnds) ;
double lon(lon) ;
lon.bounds = "lon_bnds" ;
lon.units = "degrees_east" ;
lon.axis = "X" ;
lon.long_name = "longitude" ;
lon.standard_name = "longitude" ;
double lon_bnds(lon, bnds) ;
float cSoil(time, lat, lon) ;
cSoil.standard_name = "soil_carbon_content" ;
cSoil.long_name = "Carbon Mass in Soil Pool" ;
cSoil.units = "kg m-2" ;
cSoil.original_name = "dummy" ;
```



With an iRODS-based Climate Data Server, all of the IPCC NetCDF metadata is externalized through the iCAT database in accordance with the OAIS Reference Model ...

```
gridspec_land_fx_GISS-
E2R_piControl_r0i0p0.nc areacella:
areacella_fx_GISS-E2R_
piControl_r0i0p0.nc" ;
// global attributes:
institution = "NASA/GISS (Goddard Institute for Space
Studies New York, NY" ;
institute_id = "NASA-GISS" ;
experiment_id = "piControl" ;
source = "GISS-E2-R-E135F40oQ32
Atmosphere: GISS-E2; Ocean: R" ;
model_id = "GISS-E2-R" ;
forcing = "N/A" ;
parent_experiment_id = "N/A" ;
parent_experiment_rip = "N/A" ;
branch_time = 0 ;
contact = "Kenneth Lo (cdkkl@giss.nasa.gov)" ;
references = "www.giss.nasa.gov/research/modeling" ;
initialization_method = 1 ;
physics_version = 1 ;
tracking_id = "3aed12dd-f6e0-44d1-ade9-
c064bd355edc" ;
product = "output" ;
experiment = "pre-industrial control" ;
frequency = "mon" ;
creation_date = "2011-03-21T20:30:31Z" ;
history = "2011-03-21T20:30:31Z CMOR rewrite data
to comply with CF standards and CMIP5 requirements." ;
Conventions = "CF-1.4" ;
project_id = "CMIP5" ;
table_id = "Table Lmon (31 January 2011)
a84ae296f75bb85f61668fac8fc090" ;
title = "GISS-E2-R model output prepared for CMIP5
pre-industrial control" ;
parent_experiment = "N/A" ;
modeling_realiz = "land" ;
realization = 1 ;
cmor_version = "2.5.7" ; }
```

### 4.1.2 OAIS Compliance

An Open Archival Information System (OAIS) is an archive consisting of an organization of people and systems that has the responsibility to preserve information and make it available for a designated community [8]. The OAIS reference model addresses a full range of archival information preservation functions including:

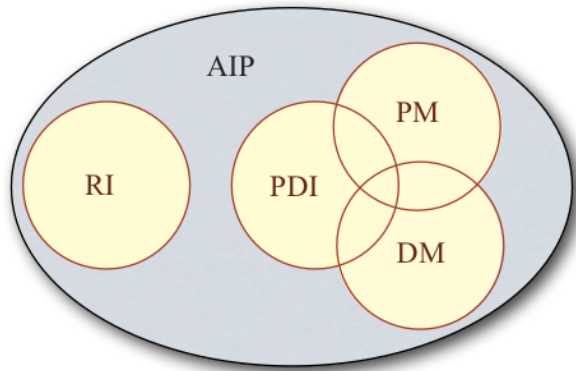
- ingest, data management, access, and dissemination;
- the migration of digital information to new media and forms;
- the data models used to represent the information, the role of software in information preservation, and the exchange of digital information among archives.
- the identification of both internal and external interfaces to the archive functions;
- it identifies a number of high-level services at these interfaces;
- it provides various illustrative examples and some 'best practice' recommendations;
- and it defines a minimal set of responsibilities for an archive to be called an OAIS.

One of the goals for the IPCC/NetCDF kits was to enable them to create collections that are compliant with the OAIS standard. As a starting point, that means the metadata about objects in the collection needs to be categorized into the types of metadata recognized by the OAIS standard. Metadata in an OAIS-compliant collection is organized around the concept of an Archive Information Package (AIP), which contains the following

classes of metadata:

- *Representation Information (RI)*: Metadata that explains how to interpret the raw data;
- *Preservation Description Information (PDI)*: Preservation related information, such as:
  - Provenance – Describes source, custody trail, and history,
  - Context – Describes relationships with internal/external data,
  - Reference – Describes unique identifiers,
  - Fixity – Describes protection from unauthorized alteration;
- *Policy Metadata (PM)*: Parent organization archive administration information.
  - Fixed (from organization business model)
  - Negotiable (from data producers and consumers)
- *Discovered Metadata (DM)*: Producer/consumer information that will foster maximal use of the archive.

The vCDS V0.9 iCAT database was extended to accommodate this OAIS metadata classification. The database schema was designed to reflect the uniqueness and repetitiveness of various NetCDF metadata items. For example, while the global attributes for each file are considered unique, the dimension, variable and function



definitions are often common across many files. The metadata idiosyncrasies influencing the database schema include the following:

- coordinate variables (input to model) have an associated dimension and function as well as a bound that defines the grid size for the associated dimension. These metadata are not unique to a given file;
- output variables (output from model) have an associated function that is not unique to a given file;
- all variables are considered to be OAIS Representation Information;
- all global attributes are considered to be OAIS Preservation Description Information. As such they additionally have an OAIS-subcategory.

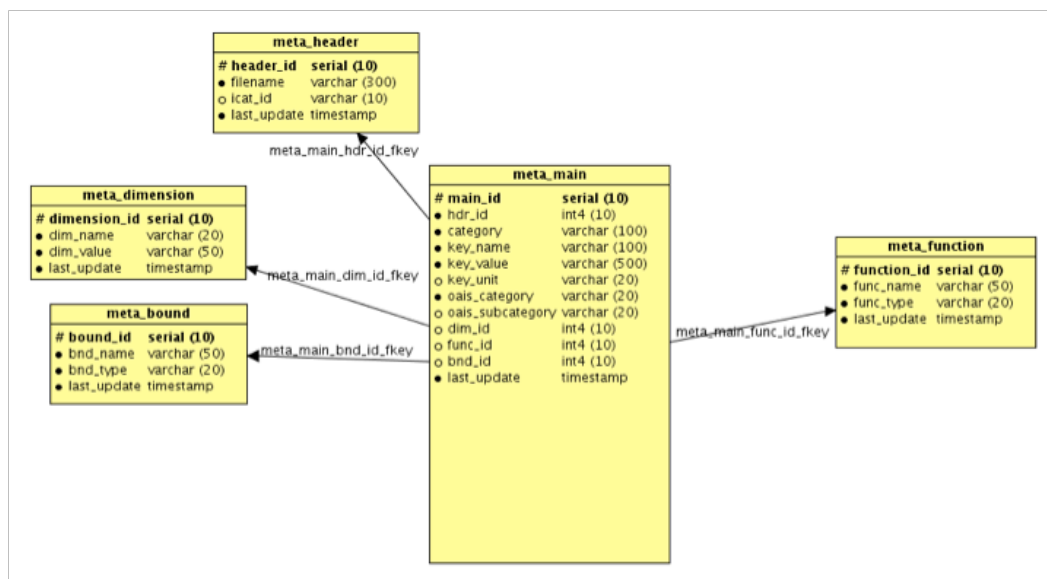
To accommodate the repetitive nature of some metadata, we implemented the following 5-table design:

- **Meta\_main** — This table contains records for each metadata item considered unique to a given file. It is composed of variable and global attribute metadata

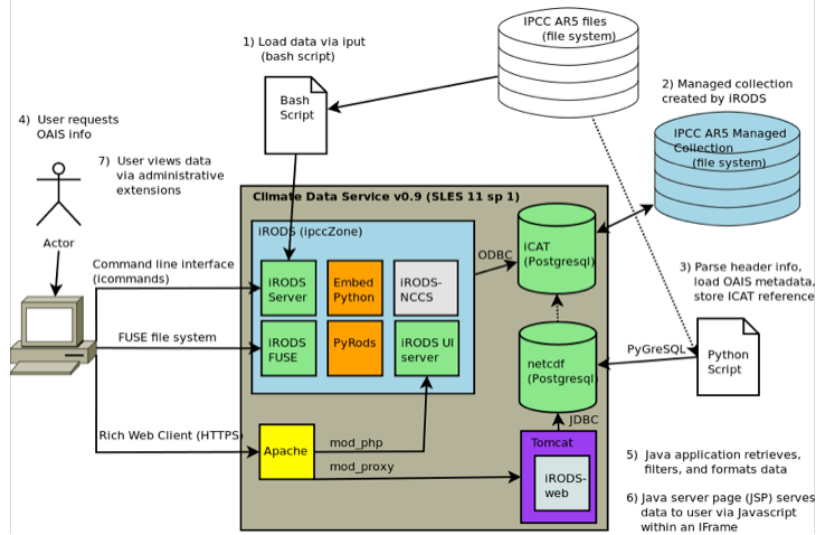
defined as name, value, and unit 3-tuples. OAIS metadata categories and subcategories as well as foreign keys to the other 4 tables are also included as required per record type.

- **Meta\_header** — This table contains the filename and an index used as a foreign key by **meta\_main** to link all the records associated with the file in **meta\_main** back to the filename in **meta\_header**. This table also includes a foreign key reference to the filename in the iCAT.
- **Meta\_dimension** — This table contains only unique entries for dimension name and dimension size. Many files may have the same set or subset of dimensions. It also contains an index used as a foreign key by each coordinate variable record in **meta\_main** to link the coordinate variable with its associated dimension.
- **Meta\_function** — This table contains only unique entries for function name and function type. Many files may have the same functions. It also contains an index used as a foreign key by each coordinate and output variable record in **meta\_main** to link the variable with its associated function.
- **Meta\_bound** — This table contains only unique entries for bound function name and bound function type. Many files may have the same bound functions. It also contains an index used as a foreign key by each coordinate variable record in **meta\_main** to link the coordinate variable with its associated bound function.

Note that, in the database schema diagram shown below, closed dots identify record entries defined as 'NOT NULL' while open dots identify entries not defined as such.



In vCDS V0.9, the “Create” operation extracts the embedded, application-independent metadata from IPCC’s NetCDF files and saves those values in tables that recognize OAIS metadata categories. This, in essence, begins construction of an OAIS Submission Information Package (SIP). At the moment of object creation, the NetCDF header information comprises primarily Representation Information (metadata on how to interpret the raw data in the accompanying file) and Preservation Description Information (PDI). The intent in future work is to fill out other classes of metadata, such as Policy Metadata (PM) and Discovered Metadata (DM) as more detailed policies are developed.



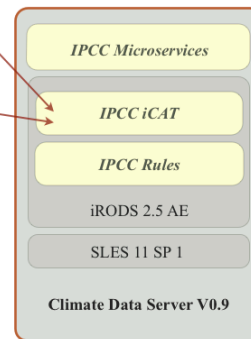
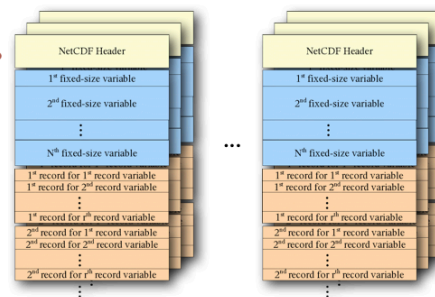
### Submission Information Package (SIP)

**Representation Information (RI):** How to interpret the raw data.

axis	standard_name	coordinates	history
bounds	long_name	_FillValue	long_name
units	associate_files	missing_value	original_name
calendar	cell_measures	flag_values	original_units
formula_terms	cell_methods	flag_meanings	standard_name
positive	comments	grid_mapping	

**Preservation Description Information (PDI):** Preservation related information.

Provenance	Context	Reference	Fixity
institution	experiment	tracking_id	<-byte_size>
institute_id	experiment_id	Conventions	<-check_sum>
contact	source	table_id	
history	model_id	creation_date	
	forcing	status_type	
	parent_experiment_id	<-file_name>	
	parent_experiment_rtp	<-activity>	
	branch_time		
	comment		
	references		
	initialization_method		
	physics_version		
	product		
	experiment		
	frequency		
	project_id		
	title		
	parent_experiment		
	modeling_realms		
	realization		
	<-temporal_subset>		



### 4.2. Administrative Extensions

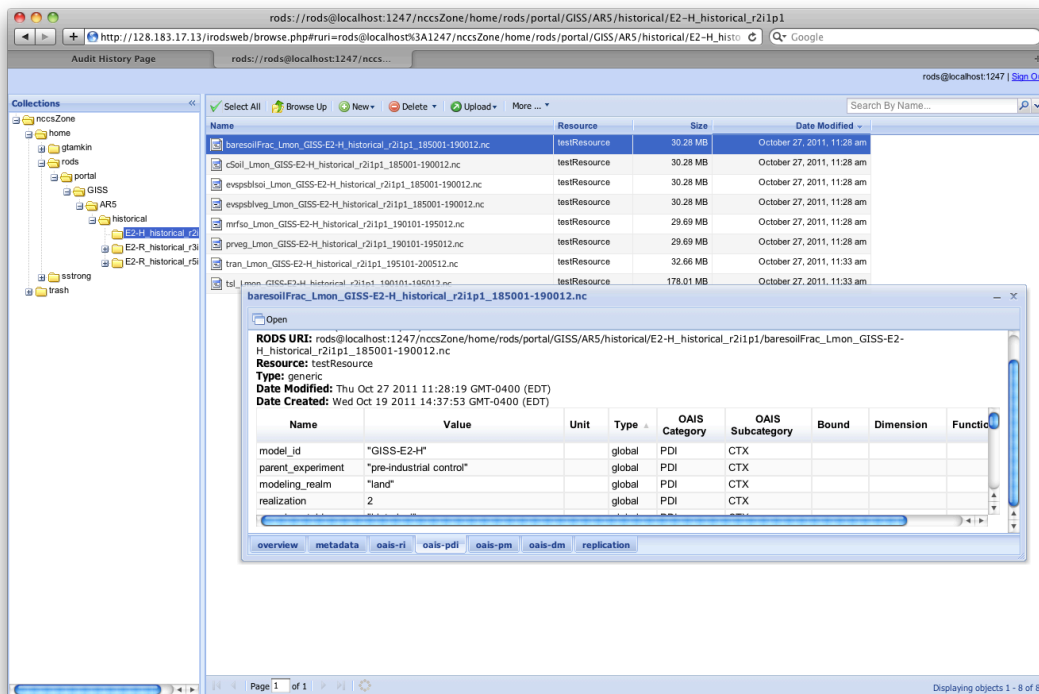
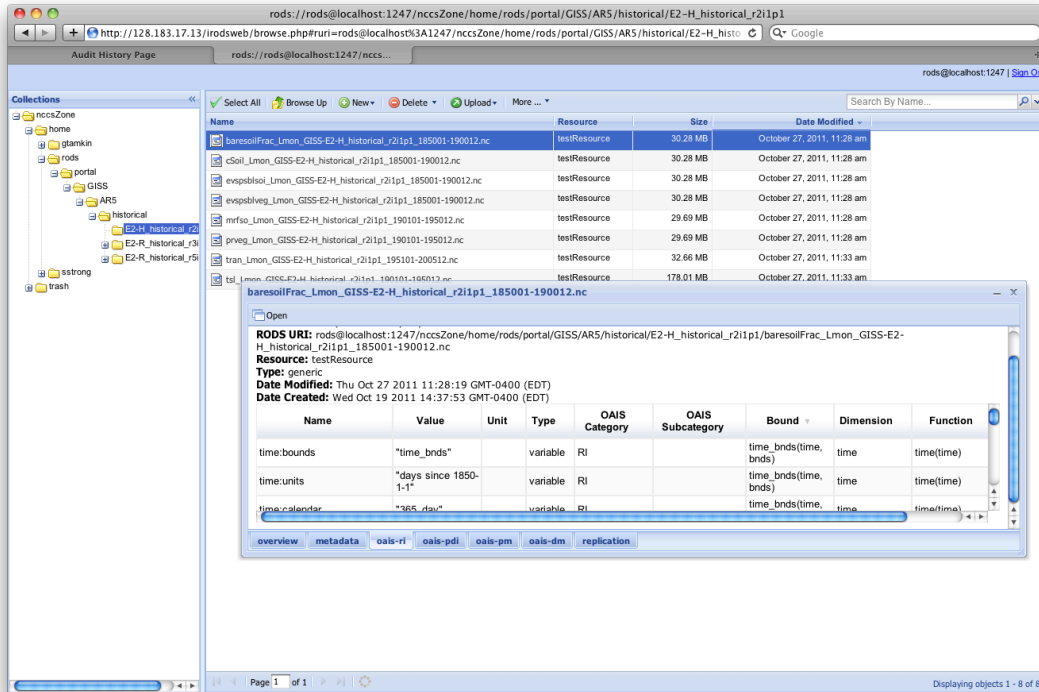
A second element in the vCDS product suite is a collection of capabilities we refer to as Administrative Extensions (AE). These include iRODS Postgres extensions and utilities to log system-level object provenance and provide QA for OAIS metadata compliance plus associated Rich Web Browser GUI extensions.

To view the NetCDF metadata by OAIS metadata categories, pre-defined SQL queries were added to the iRODS Rich Web Client.

Displays of RI and PDI are depicted in the following images. First is a sample coordinate variable query result showing selected metadata entries from meta\_main (Name, Value, Unit, Type, OAIS Category, OAIS Subcategory), from meta\_bound (Bound), from meta\_dimension (Dimension), from meta\_function (Function) and from meta\_header (ICAT id). Note that these coordinate variables are categorized as Representation Information (RI).

Next is a sample global attribute query result showing selected metadata entries from meta\_main (Name, Value, Unit, Type, OAI Category, OAI Subcategory) and from meta\_header (ICAT id). Note that

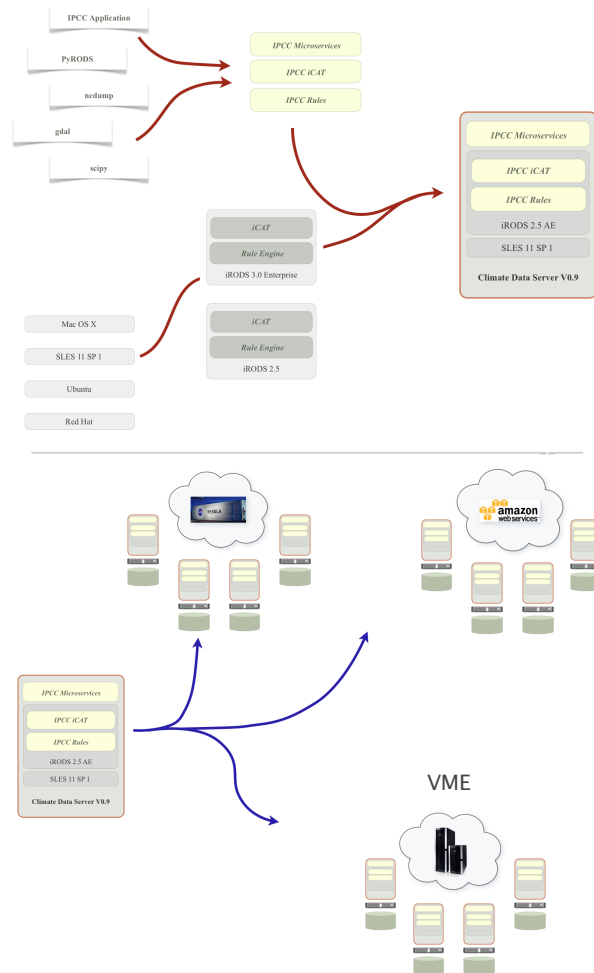
global attribute have no dimension, no function and no bound. Also note that these global attributes are categorized as Preservation Description Information (PDI) each with subcategory Context (CTX).





### 4.3 Repetitive Provisioning

We work in a virtualized environment that includes MacBooks running VMware Fusion, a vSphere dev/test server farm, a NASA cloud computing environment called Nebula, as well as Amazon's Elastic Compute Cloud (EC2). As one of the major accomplishments of this software release, we developed RPM Package Manager (RPM) scripts to build software stacks for SLES 11 SP1, iRODS 2.5 with Administrative Extensions, and vCDS V0.9 virtual images.

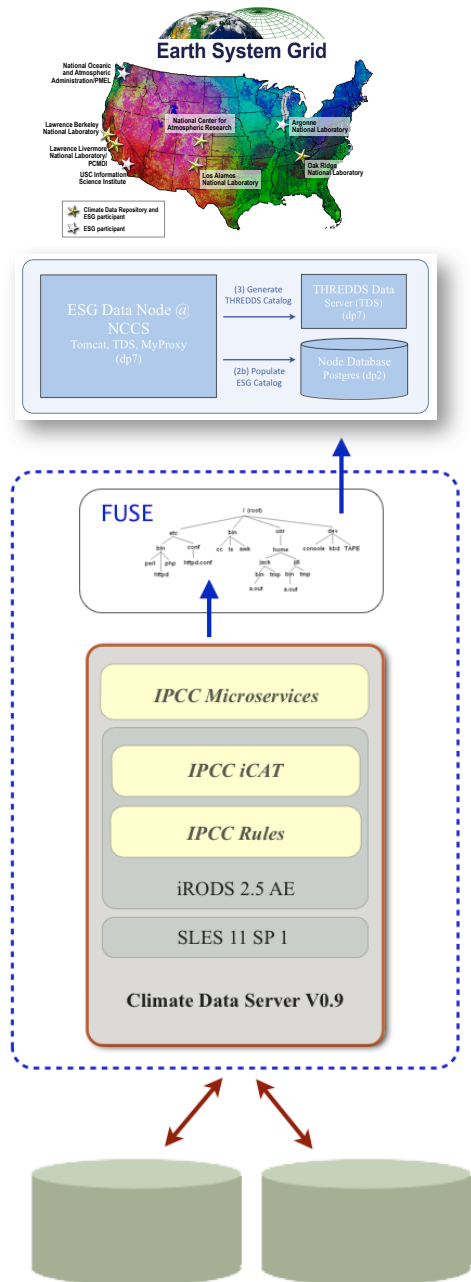


### 4.4 Deployment and Distribution

RPM makes it easy to set up automated build and install procedures consisting of many packages for an entire operating system. When these images are provisioned into a virtual cloud resource, capabilities can be delivered as Infrastructure-as-a-Service (IaaS) (e.g. SLES 11 SP1), Platform-as-a-Service (PaaS) (iRODS 2.5 AE), and Software-as-a-Service (SaaS) (vCDS V0.9). Collectively, our ability to provision into these various resource classes enables virtualization-as-a-service (VaaS), which is huge unmet need in our domain.

## 5. Operational Deployment

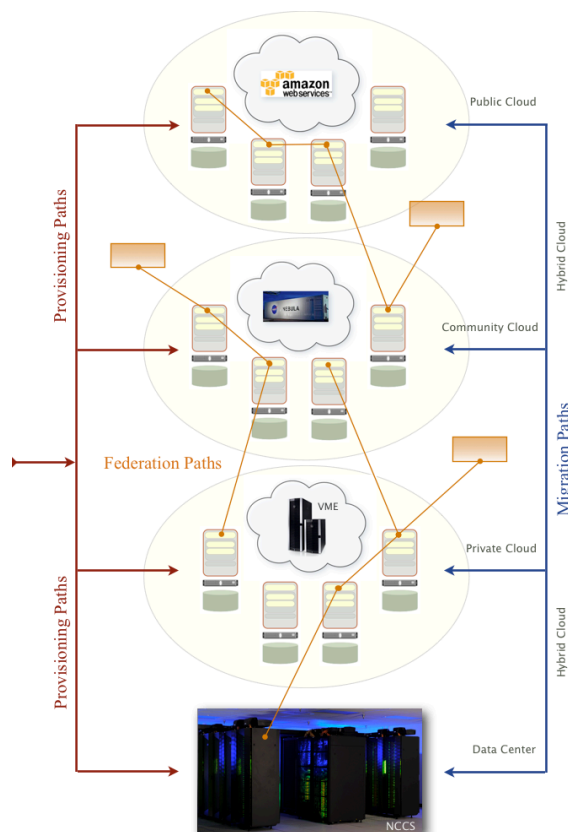
vCDS V0.9 is being deployed in the Amazon Elastic Compute Cloud where it will be hardened for operational use. These end-of-system-development activities will essentially elevate vCDS V0.9 to V1.0 at TRL 8/9. Its first application will be to manage a collection of IPCC AR5 data in EC2 for publication through the Earth System Grid. Since the ESG gateway requires a filesystem view of the data it serves, we are using FUSE (File System in User Space) to expose the vCDS collection to ESG.



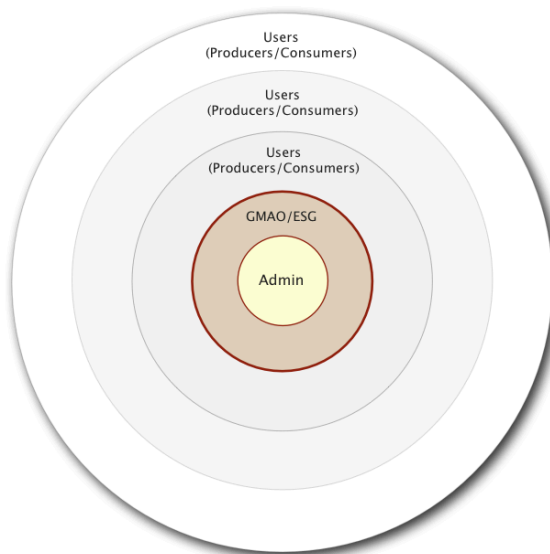
## 6. Discussion

Taken together, the elements of this work that we refer to as the vCDS product suite — the NetCDF/IPCC kits, administrative extensions, and utilities for automatic provisioning, and deployment and distribution — enable an approach to scientific collections management in which virtualization is a driving concept. It supports access to a tiered array of cloud services that are flexible, adaptable, scalable, and stageable to NCCS “bricks and mortar” facilities as needed. We can provision capabilities into any resource class, migrate images from one resource class to another, and use the iRODS federation mechanism to assemble virtual collections that cross resource classes. This approach provides an agile entry point into the NCCS for new customers with data-centric requirements and enables virtualization-as-a-service.

With the vCDS approach, we are trying to enable the full information lifecycle management of OAIS-compliant scientific data collections. A vCDS manages data as a distinguished collection for a person, project, lab, or other logical unit. A vCDS can manage a collection across multiple storage resources using rules and microservices to enforce collection policies. And a vCDS can federate with other vCDSs to manage multiple collections over multiple resources thereby creating what can reasonably be thought of as an ecosystem of managed collections.



## 7. Conclusions



Up to now, we have been focusing on publishing IPCC AR5 data to the Earth System Grid. In OAIS parlance, GMAO and GISS are our first data producers, ESG (an application) is the first consumer. The first customer will be the vCDS Collection Administrator, and the IPCC AR5 dataset will be our first vCDS managed collection. Follow-on work will focus on expanding the array of managed collections and broadening our community of users, which means expanding vCDS policies, rules, and microservices.

## References

- [1] NASA Center for Climate Simulation (NCCS), <http://www.nccs.nasa>.
- [2] Schnase, J.L., Tamkin, G., Fladung, D., Sinno, S., and Gill, R. 2011. Federated observational and simulation data in the NASA Center for Climate Simulation Data Management System Project. 2011 iRODS Users Group Meeting, <http://iren-web.renci.org/irods-meeting/nasa.pdf>.
- [3] Intergovernmental Panel on Climate Change (IPCC), <http://www.ipcc.ch>.
- [4] Earth System Grid (ESG), <http://www.earthsystemgrid.org/about/overview.htm>.
- [5] Technology Readiness Level (TRL), [http://esto.nasa.gov/files/TRL\\_definitions.pdf](http://esto.nasa.gov/files/TRL_definitions.pdf).
- [6] Network Common Data Form (NetCDF), <http://www.unidata.ucar.edu/software/netcdf>.
- [7] CMIP5 Data Reference Syntax and Controlled Vocabularies, [http://cmip-pcmdi.llnl.gov/cmip5/docs/cmip5\\_data\\_reference\\_syntax.pdf](http://cmip-pcmdi.llnl.gov/cmip5/docs/cmip5_data_reference_syntax.pdf).
- [8] Open Archive Information System (OAIS) Reference Model, <http://public.ccsds.org/sites/cwe/rids/Lists/CCSDS%206500P11/Attachments/650x0p11.pdf>.