# IS-ENES – WP10

# D10.3 - Operational service package

## Abstract

| Grant Agreement number | **228203** | Proposal Number: | **FP7-INFRA-2008-1.1.2.21** |
|---|---|---|---|
| **Project Acronym:** | **IS-ENES** | | |
| **Project Coordinator:** | **Dr Sylvie JOUSSAUME** | | |

| Document Title: | Operational Service Package | | Deliverable: | D10.3 |
|---|---|---|---|---|
| **Document Id N°:** | D10.3 | **Version:** | 1.0 | **Date:** | 27.03.2013 |
| **Status:** | Final | | | | |

| Filename: | ISENES_D10.3_final |
|---|---|
| **Authors:** | |
| **Project Classification:** | Public, Confidential |

## Revision Table

| Version | Comments |
|---|---|
| Draft03 | Circulated to reviewers |
| Draft04 | Circulated to WP partners |
| Final | Final adjustments completed |

# 1      Executive Summary

The ENES data archive is now supported by stable software infrastructure, the ESGF Peer-to-Peer system (P2P). Development of P2P was led by PCMDI and put into place for the global CMIP5 archive, but IS-ENES partners contributed key components of the system and played a major role in requirements specification, early deployment and testing. The final system meets the projects gaols of a distributed archive infrastructure allowing data at multiple sites to be accessed through a central interface with a range of advanced data services. Access to data is controlled by security certificates, allowing usage to be monitored and reported. Checksums are provided along with tools which verify correct transfers for users, ensuring safe transport of large data volumes. Access to subsets of data is enabled through the widely used OPeNDAP standard, allowing user applications to extract small subsets of data where complete files are not needed. The system has been tested in extremis, distributing the peta-scale CMIP5 archive, and is running smoothly (see SA2 reports for details). Effort set aside for links to a "METAFOR discovery service" mentioned in the IS-ENES Description of Work (DoW) has been re-deployed, as the METAFOR project did not deliver an operations ready discovery service. More effort than expected was put into security systems and publication systems in order to ensure that the European archive integrated smoothly into the globally distributed ESGF system. The impact of the missing METAFOR discovery functionality which had been envisioned in the DoW is replaced by an extremely efficient and flexible search capability delivered by US collaborators.

# 2      System overview

The ENES data archive is designed as a federated system, linking resources at existing data centres with ongoing independent funding streams, into a unified archive with transparent access to data and a range of advanced data services. The overriding requirements of the scientific community were associated with efficient delivery of the CMIP5 archive. This could only be achieved by through close co-operation with the PCMDI team leading the CMIP5 archive support on behalf of the World Climate Research Programme (WCRP). Integrating the system with partners outside the consortium within a loose federation led by PCMDI brought many problems, but also many advantages through the shared development effort. The objectives of the many groups involved were generally well aligned because of the shared objective of supporting the CMIP5 archive, but differing national and institutional priorities still led to conflicts at times.

The system allows different institutions to perform a range of roles, listed below. The modularity of the system allows the federation to combine data providers with limited technical and user support resources available for this role and archive centres with more experience in data provision.

## 2.1   Serving institutional data

At the simplest level, an institution can use the system to publish their own data, so that these data appear in the search interface(s) supported by the archive centres.

## 2.2   Serving community data

Serving data from multiple institutions does not add a great deal of technical complexity, but does bring the additional responsibilities to ensure clarity of provenance information and the additional workload of dealing with data providers and managing streams of data from the providers.
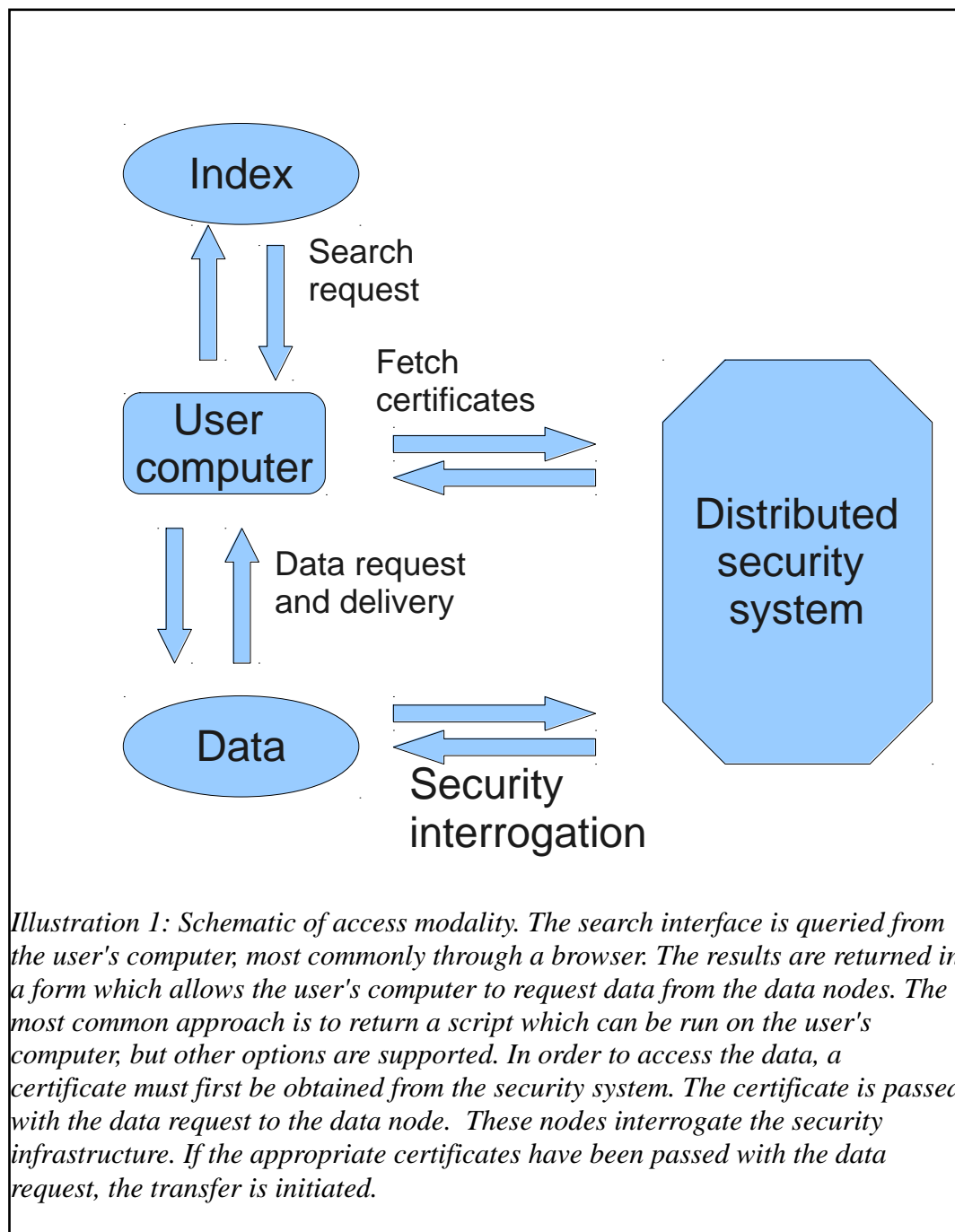
## 2.3    Providing a search interface

The search interface is now relatively easy to implement from a technical perspective, but institutions that run it should ensure that appropriate support is available for users who have difficulty with the system.

## 2.4    Providing user registration and authentication services

Dealing with user registration requires extra care to protect confidential user information and to ensure that user access to the portal does not compromise security.

## 2.5    Providing authorisation services

Authorisation services are critical to the archive function. As all access data from a project must be authorised by the institute responsible for managing the IPR of that project, the authorisation service is a single point of failure.

*Illustration 1: Schematic of access modality. The search interface is queried from the user's computer, most commonly through a browser. The results are returned in a form which allows the user's computer to request data from the data nodes. The most common approach is to return a script which can be run on the user's computer, but other options are supported. In order to access the data, a certificate must first be obtained from the security system. The certificate is passed with the data request to the data node. These nodes interrogate the security infrastructure. If the appropriate certificates have been passed with the data request, the transfer is initiated.*

# 3   Sub-tasks

## 3.1   Security system

A federated security system underpins many of the advanced services delivered by federated archive. The security system needs to support systems required for delivery of the IS-ENES programme of work. The scale of the archive is such that data cannot be delivered through a web interface. The security system supports access both by browsers (via OpenID authentication), and desktop clients and libraries (via X.509 certificates[1]). The adoption of X.509 certificates caused

---

[1] **X.509** is a standard for a public key infrastructure maintained by the International Telecommunication Union

some initial problems for users who are unfamiliar with the certificate management, but is widely supported by third party software and hence supports a flexible and extensible range of archive services. Initial prototypes of the CMIP5 archive infrastructure, produced by the US Earth System Grid project, relied on tokenised URLs and would only have supported access through a browser interface. The initial release for the operational archive system in 2010 supported both the advanced certificate system and the primitive token based system, as some partners in the global federation had doubts about maturity of the certificate-based system. IS-ENES partners installed the certificate-based system from the start and demonstrated its reliability and the clear benefits of additional flexibility. Other nodes in the global federation gradually followed the IS-ENES example, leading to the eventual phasing out of the token-based approach.

While the use of X.509 certificates is widespread, there were difficulties particular to the implementation within the archive infrastructure. Establishing a reliable system to manage certificates from all participating institutions was difficult, particularly as many partners outside Europe had assigned very limited resources to the operation of archive nodes. It was also necessary to provide support for an extremely diverse user community. Security certificates are widely used in GRID computing in a context where users have specialist support from their institutions. The configuration issues which need to be dealt with on user machines are straight forward for a specialist but can be confusing for users who are not familiar with system configuration problems.

Three tools to support user certificate management are provided, to ensure that there is support for a broad range of user operating systems. Additionally, it is possible for users to obtain certificates using the globus toolkit (see below).

## Java webstart

The java webstart application allows users to obtain a secure certificate from a web interface. This approach will wok on any system with sufficiently recent implementation of the Java library, provided that the user has authority to run the web-start application. To obtain certificates users simply need to go to the following URL and follow instructions:

http://rainbow.llnl.gov/dist/esg-myproxy-logon/MyProxyLogon.jnlp

## MyProxyClient python script

For users who are familiar with the python language, the MyProxyClient package may be more comfortable. The software should be installed by someone with "sudo" access on the local computer, using the following command:

sudo easy_install MyProxyClient

The following command will then place a valid certificate in the "creds.pem" file:

myproxyclient logon -b -T -s <MyProxy server> -l <IDP User ID> -o creds.pem,

where <MyProxy server> is the URL of the MyProxy server of the user's identity provider and <IDP User ID> is their local user identity with that identity provider (not generally the same as the OpenID user name).

http://ndg-security.ceda.ac.uk/wiki/MyProxyClient

http://esgf.org/wiki/ClientAccessToESGFOPeNDAPServers

## Bash script collection (enesGetCert)

These scripts will prompt the user for their OpenID, IDP User ID and MyProxy server URL. The script is designed to run on Linux and Mac OS operating systems. One of these scripts is embedded in the download scripts issued by the search interface. Two variants are provided, one using a java utility and a second which uses curl. These two variants are intended to provide flexibility and ensure that support is offered for a broad range of platforms.

The getcert.jar file was developed in IS-ENES, with a string reliance on the jglobus library (

http://dev.globus.org/wiki/CoG_jglobus)

## Globus toolkit

User who have the Globus toolkit (http://www.globus.org/toolkit/) installed, can obtain a certificate with the following command:
myproxy-logon -s <myproxy-server> -l <username> -o credentials.pem

| ENES identity providers (OpenID) and certificate servers (MyProxy) | | |
|---|---|---|
| STFC | Index node | esgf-index1.ceda.ac.uk/ |
| | OpenID | ceda.ac.uk/openid/ |
| | MyProxy | myproxy.ceda.ac.uk |
| | User registration | badc.nerc.ac.uk/reg/user_register_info.html |
| DKRZ | Index node | esgf-data.dkrz.de/esgf-web-fe/ |
| | OpenID | esgf-data.dkrz.de/openid |
| | MyProxy | esgf-data.dkrz.de |
| | User registration | esgf-data.dkrz.de/esgf-web-fe/createAccount |
| IPSL | Index node | http://esgf-node.ipsl.fr |
| | OpenID | http://esgf-node.ipsl.fr |
| | MyProxy | esgf-node.ipsl.fr |
| | User registration | esgf-node.ipsl.fr/esgf-web-fe/createAccount |

### 3.2   Check-summing

A script developed at DKRZ is now a key part of the global ESGF archive system. Once users accessing data through the portal have determined the selection of data they require, a script is generated which they can run to transfer the data to their local machine. The script, developed by DKRZ, verifies the checksum of every file transferred, using the MD5 protocol. Ongoing discussions with PCMDI have ensured that checksums for all files are recorded in the archive metadata. Experience with long distance transfer of large files suggests that file corruption is a significant problem when moving terabytes of data, and corruption can happen while preserving file size.  IS-ENES partners led the way in deploying systems with comprehensive check-sums.

On publication, checksums are included in the archive catalogue.  The catalogue checksums may be used by client software (such as synchro-data and replication tools described in D10.4) and are also

used directly by the portal when users request a download script. The script passed to users contains addresses of files matching their request together with the checksums. When the script is run checksums are computed after transfer of files to the users' local machines and verified against the archive checksum.

## 3.3    Quality control

Considerable effort has been invested in standardisation of quality control procedures and reporting. The scale of the CMIP5 archive provides many challenges, with over 50 climate models providing data for more than 800 variables. A key requirement has been the creation of a system for recording the results of quality control tests and in a form which is convenient and useful for both users and data providers. IS-ENES has developed a system in which a repository of quality control results can be viewed through a dedicated user interface. Once all quality control tests have been dealt with and critical problems resolved it is possible to assign a Digital Object Identifier (DOI) to the data to act as a permanent reference. When this stage is reached, the DOI can be linked to the model documentation. Also a report of the quality status is provided in the form of a Metafor CIM data quality document. These documents are harvested by the CIM portal and are visible to users based on e.g. the CIM viewer integrated in ENES index nodes. Display of model documentation is discussed further in D10.4.

## 3.4    Initial publication under the "gateway system"

In 2011 much of the European data in the CMIP5 archive was published through the federated system which had been jointly developed by NCAR and PCMDI up until 2010. This system has since been withdrawn, but some components contributed by IS-ENES partners have been moved into the new Peer-to-Peer system. Both systems use the same implementation of the UNIDATA THREDDS data server, and IS-ENES put considerable effort into establishing and maintaining well defined interfaces between components to ensure that the transition from the "gateway system" to the "P2P" system did not interfere with the operation of European services.

## 3.5    Transition to P2P

During 2012 the archive infrastructure has gone through a major upgrade, led by PCMDI. The new system features a more responsive search interface and streamlined authorisation of users, removing two performance bottlenecks. IS-ENES contributed with early deployment of evaluation versions of the system and coordination of the requirements specification. Table 1 indicates the dates of significant system deployments.

| Institute | Date | Component/version | Comments |
|---|---|---|---|
| IPSL | 10/2011 | v1.1.1-bay_ridge-release | Test federation. Data node only. |
| IPSL | 12/2011 | v1.2.0-bensonhurst-release | Test federation. All components working. |
| STFC | 03/2012 | Data node / v1.3.1 | Initial deployment of P2P compatible data node |
| DKRZ | 03/2012 | Index Node | Test federation component |
| IPSL | 04/2012 | v1.3.2-borough_park-release | Transition preparation phase completed. All components operational. Publication towards both systems. |
| IPSL | 06/2012 | v1.4.0-Brighton Beach | Transition done. All components operational. |

| | | release | Publication toward the new system only. All previously published datasets has been republished. |
|---|---|---|---|
| DKRZ | 06/2012 | Index Node | Production tests |
| STFC | 06/2012 | Index node / v1.3.4 | Initial deployment of index node and P2P web interface in test federation |
| DKRZ | 07/2012 | Data Nodes | Upgrade to use new P2P node installer |
| DKRZ | 07/2012 | Data Nodes | Configuration for CORDEX test node |
| STFC | 08/2012 | Index node / v1.4.0 | Production P2P system |
| STFC | 08/2012 | Data node / v1.4.0 | Upgrade production data node |
| DKRZ | 08/2012 | Index Node | To run parallel services with old system |
| DKRZ | 10/2012 | Index Node | Migration of users to new system |
| STFC | 10/2012 | Data node (replicas) / v1.4.1 | Separate data node deployed for serving CMIP5 replicated data |
| DKRZ | 12/2012 | Gateway | Decommissioning of old gateway. |
| LIU | | | CORDEX test node |
| STFC | 12/2012 | Index node / v1.4.2 | Upgraded P2P web interface |

*Table 1: Deployment summary for P2P infrastructure*

For the final deployment, all previously published datasets had to be republished.

## 3.6 Data ingestion and version control

Robust data ingestion procedures are required to ensure that large data volumes are transferred into the archive without corruption of data or metadata. The "drslib" library, developed at STFC through IS-ENES, is now part of the standard distribution. Drslib supports adoption of the standardised directory structure and implements the CMIP5 version control protocol. Drslib was developed for CMIP5, but is being re-configured in preparation for CORDEX and other future projects.

## 3.7 Replication software

The distributed archive infrastructure has many benefits, but the speed of access to large datasets held at remote locations is not adequate for all purposes. Consequently, the ENES archive needs to bring copies of large datasets to Europe. Replication software has been developed and deployed with the following functions:

•      A discovery component to locate and find files (and new versions of files) to be replicated (this is controlled by a priority list of file collections to be replicated next)

•      A download component, responsible for parallel download of files based on multiple transfer protocols (e.g. http, and gridftp) – this also includes data integrity checks

•      A storage component, responsible for storage of the replicated data sets in a consistent file hierarchy (e.g. using soft or hard links)

•      A publication component – publishing the replicated data as replicas to a index node

An adapted version of the DKRZ replication software is also in production at PCMDI (which e.g. modified step 3 to their local storage layout needs) and is available at: git://esgf.org/esgf-

replication.git

IPSL uses the synchro-data tool described in D10.4 (IPSL are only completing steps 1 to 5 above and then providing a copy for users who have direct access to their data store). STFC uses a collection of scripts[2] which perform a similar function to ESGF-replication. This diversity of approaches reflects the difficulties in finding a general solution to fit with all local operation procedures. Within IS-ENES there has, however, been considerable cross-fertilization between the centres leading to a consistent approach. The adoption of the ESGF-Replication software by PCMDI, in preference to solutions being developed by ESGF partners in the US, is an indicator of the success of this aspect of IS-ENES development work. The software is not only being used to bring data to Europe, it is also being used by PCMDI to transfer large volumes of European data to PCMDI to provide faster access for US scientists.

## 3.8  Metadata display

Since the 4[th] IPCC Assessment Report, both the IPCC and WCRP have called for improved documentation of climate models used in international climate assessments. Components of an advanced documentation system were developed in the FP7 METAFOR project, in close collaboration with the US CURATOR project[3]. IS-ENES intended to take over operational maintenance of the climate model metadata and associated services after the end of the METAFOR project. The partnership between METAFOR and CURATOR lead to a globally agreed standard for model documentation which has been adopted by the WCRP for the CMIP5 archive, but left some gaps in the software infrastructure. This infrastructure is also required to interact with the CMIP5 archive infrastructure which, under the leadership of PCMDI has undergone radical changes since the end of the METAFOR project. Work to support access to model metadata in the new architecture has continued in the G8 Research Initiative project ExArch and has been undertaken by IPSL. IS-ENES has contributed through the deployment of early versions of the P2P system, and ensuring that visibility of metadata remains a priority.

The infrastructure for metadata collection and presentation is highly modular. Information has been entered into the system through an online questionnaire developed in FP7 METAFOR and run operationally at STFC[4]. Completed questionnaire entries are harvested by the repository, which is currently being developed at IPSL with support from the G8 ExArch project and will be taken into full operational deployment in the FP7 IS-ENES2 project. The repository supports a JAVA Application Programming Interface (API), and provides a library for incorporation into portals. This library powers the meta-data user interface in the ESGF ENES portals.

Illustration 7 shows the link to the model documentation which appears in the search results of the ESGF ENES portal and illustration 8 shows an example of the model documentation which is available.

---

[2]http://home.badc.rl.ac.uk/mjuckes/stfcEsgfRep/
[3]Supported by NOAA through the NOAA Environmental Software Infrastructure and Interoperability
  (NESSI: http://www.esrl.noaa.gov/nesii/ ) group.
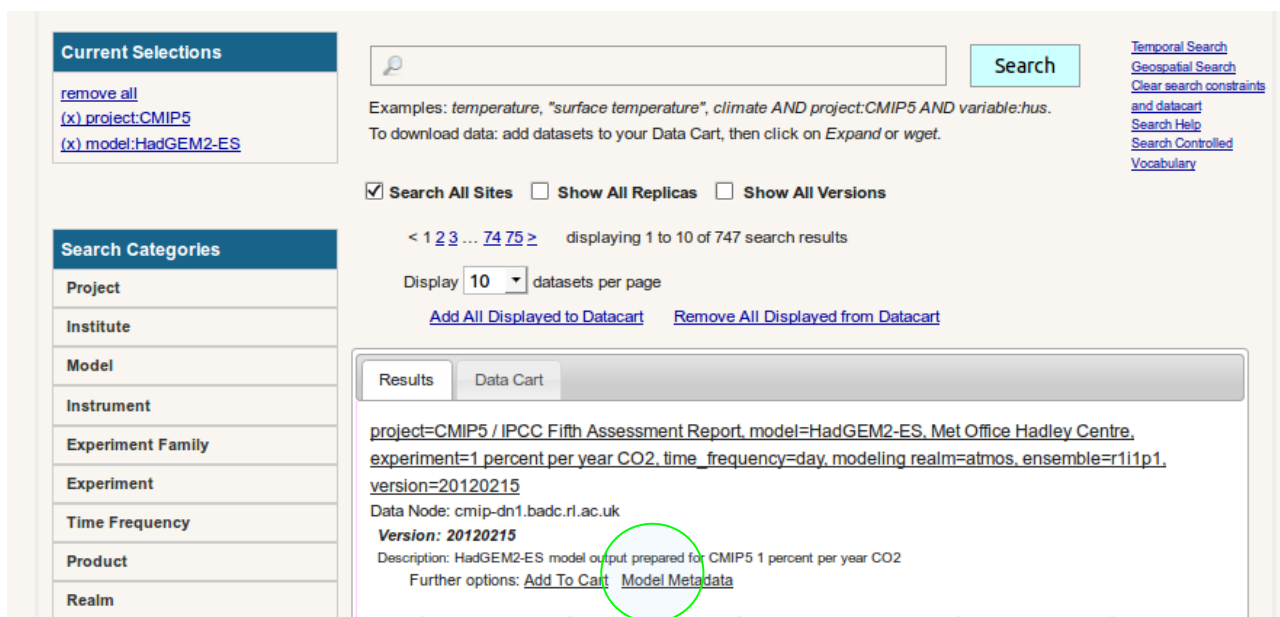[4]http://q.cmip5.ceda.ac.uk/

*Illustration 7: Screenshot of search results in the ESGF IS-ENES STFC portal. The green ellipse shows the link extended model metadata which is provided with the search results.*
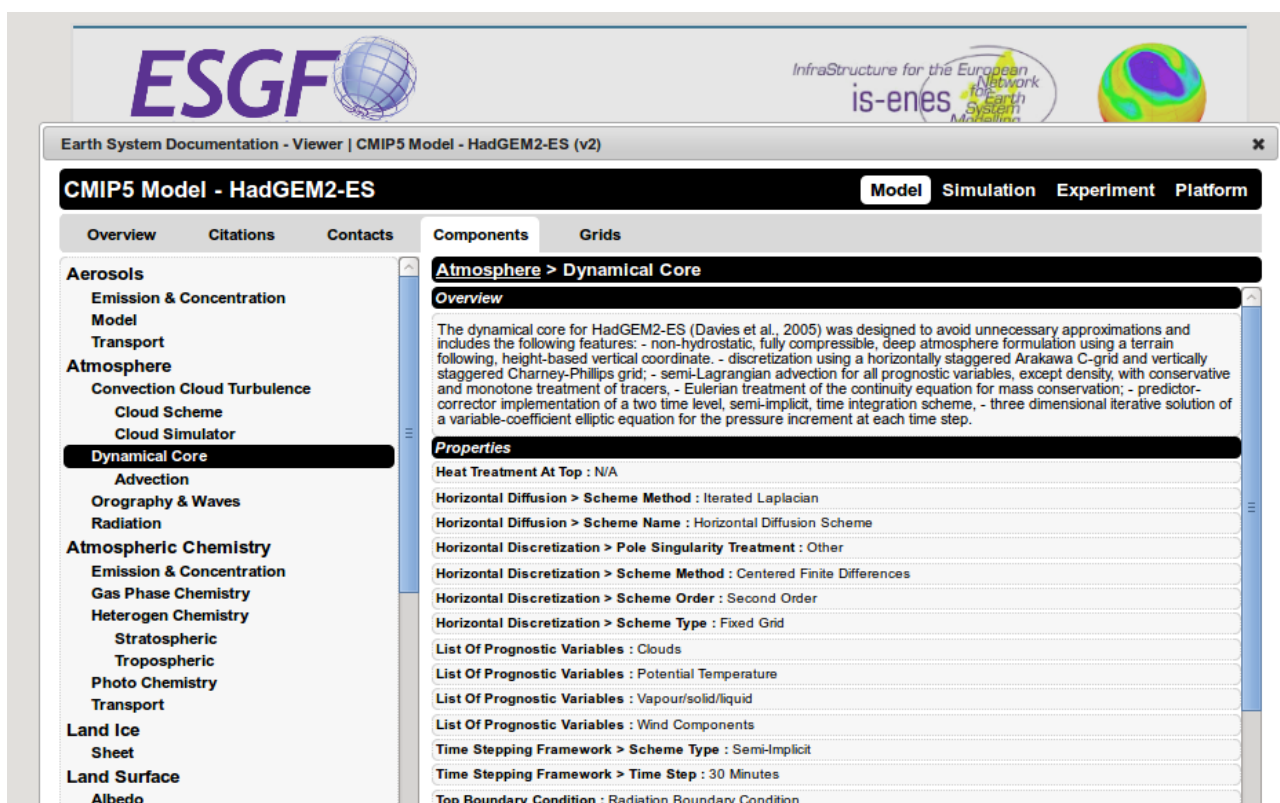


*Illustration 8: Screenshot of model documentation, showing view with "Model" selected in the main menu bar, "components" in the second and "Dynamical core" in the tree of components.*

## 3.9

## 3.10  CORDEX Support

### CORDEXwriter: Formatting data for the CORDEX archive

A set of scripts is designed to produce the CORDEX CORE and Tier 1 output in netcdf format.

All scripts are bash-shell scripts and use the CDO and NCO packages.

Input can be in NetCDF or GRIB format.  There are 3 main scripts, dealing with daily monthly and seasonal data respectively. The scripts will form the data into the required time slices and ensure that global attributes and dimensions are correctly placed in the files. For some variables time averaging is also performed.

### Configuring ESGF for CORDEX

The CORDEX project will produce regional climate projections down-scaled from the CMIP5 global projections. Managing the data from CORDEX requires some additional features in the archive configuration. Several meetings were held with CORDEX scientists, resulting in a clear statement of data requirements for the CORDEX regional climate model data. An ESGF data node was then configured to publish data complying with the requirements and test data published to the ESGF federation.

### CORDEX "MIP" tables

"Model inter-comparison project" tables (MIP tables) are a de-facto standard for describing the data to be archived in a model inter-comparison project. The system has been developed by PCMDI, and the tables are used by the CMOR software developed and maintained at PCMDI. All CMIP5 data in the ENES archive has been produced in a standard file format using the CMOR tool. For CORDEX a more flexible approach has been adopted (by the CORDEX community) with a less rigorous file format specification. Nevertheless, key information is encoded in the CORDEX MIP tables:

http://www2-pcmdi.llnl.gov/cmor/tables/copy4_of_cmip5-tables/ – produced by IPSL.

## 3.11  COMPSs

The Barcelona Supercomputing Centre (BSC) provided the COMPSs (COMP Superscalar) programming framework in order to provide a tool to help the development of scientific workflows suitable for the modelling community and to optimize the orchestration of the execution exploiting the inherent parallelism of applications when running them on distributed infrastructures as Grids. COMP Superscalar has been introduced during the project as an evolution of GRIDSs; COMPSs. However, with respect to its predecessor, COMP Superscalar runtime is formed by a set of components, each one in charge of a given functionality. This componentised runtime follows the Grid Component Model (GCM), a component model especially designed for the Grid whose reference implementation is provided by the ProActive open source project (proactive.inria.fr).
**A first example application for multi-model ensemble mean of surface temperature has been ported to COMPSs to COMPSs (**

Figure 1). CDO utilities have been used to compute the multi-model ensemble mean. After exploring the computing task, five main steps have been identified: (1) reference grid definition, (2) temporal slice selection, (3) regridding to reference grid, (4) temporal mean computation, and (5) model mean computation. Only task (1) and (5) are computed by the master node. The three main operation of the application have been identified as candidate workers to be run on the computing nodes.
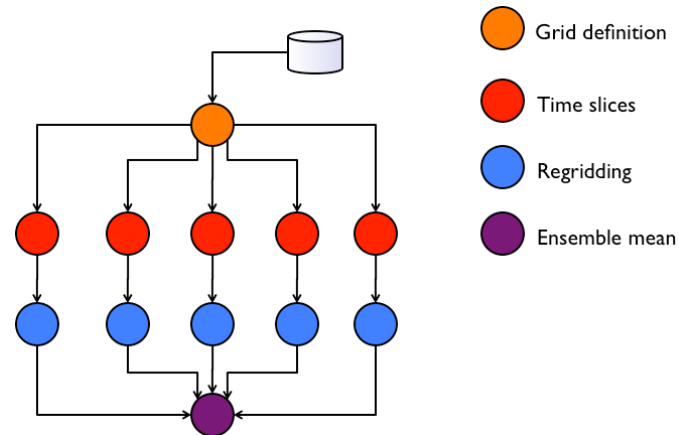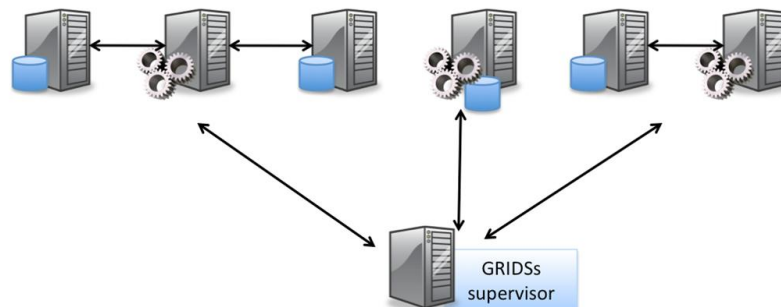


**Figure 1 – Multimodel ensemble mean implemented in COMPSs**

As depicted in Figure 2 several hosts of the grid have been configured as computational, archive or



archive+computational nodes; in this way the behaviour of the COMPSs supervisor will be simulated scheduling the tasks based on the files locality information and/or on the status of the links connecting the computational and archive node; to this aim the evaluation of the integration in COMPSs of the dynamic status of the network

The results of the tests with one master node coordinating the execution of the tasks on the computing nodes, demonstrated that the computing tasks were too small to make the application scalable considering that the runtime overhead (coordination and transfer time of the input data) is more than the benefit of parallelizing the original sequential application.

**Figure 2 – COMPSs testbed**

To overcome these issues a second version has been developed that allows the COMPSs scheduler to send the tasks where the input dataset are, avoiding data transfers and also the possibility of using bigger archives as inputs. This version improved as expected the speed up of the application but still stressing the fact that the computing part of the execution is small to exploit a bigger number of nodes.

Another example application implements a workflow that orchestrates the execution of a high resolution global atmosphere ocean models in grid environments. Figure 3 depicts the execution graph orchestrated by COMPSs. Also in this case, the availability of data locally to a data node has been simulated.
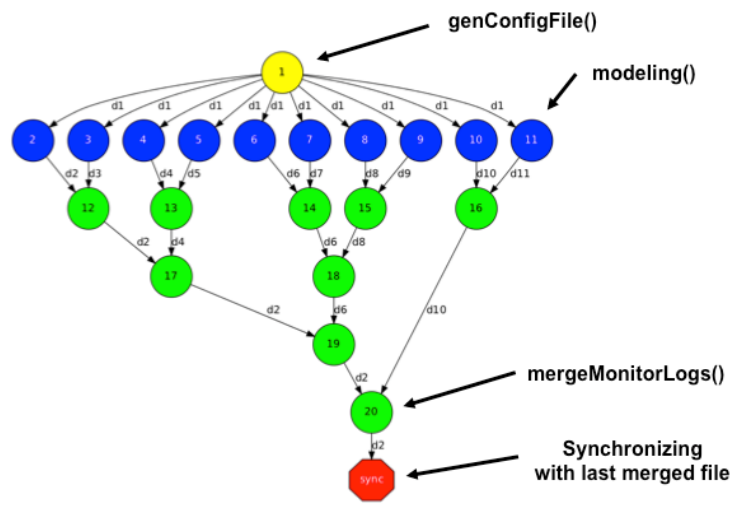


**Figure 3 – COMPSs tasks workflow of a high resolution coupled model orchestration**

An expected improvement was to evaluate the use of the ESG data nodes as source but no real requirements have been received asking for a COMPSs support in this case.

## 3.12 Dashboard

During 2012, the Dashboard system has been both improved (several bugs have been fixed) and extended (new modules and functionalities have been added). Two new releases have been integrated into the ESGF software stack during the year.

In terms of new features, the Dashboard is now able to collect information about:

1.      CPU and memory to provide local metrics statistics;

2.        the node-types to provide deployment maps at peer-group and federation level;
3.        registered users to show the *number of users/Identityprovider* and the total number of users;

It is important to remark that points 1 and 2 have been completed and are now in production. Point 3 is in production too and it mainly allows to manage the ESGF Identity providers (future extensions will allow to include additional types of identity providers). During 2012 a new Dashboard module related to the data usage statistics has been also designed and implemented. Presently, it is in a pre-production stage and it will be finalised early in 2013.

An important addition to the ESGF-Dashboard software has been the ESGF-Desktop part. It represents a pure Javascript module providing desktop-like functionalities through a web interface. It has been completely designed and implemented during 2012 and it integrates several components. The more relevant ones are a Twitter gadget (to share comments and news), a RSSFeedViewer (to display the RSSFeeds related to each data node), a Terminal (to run some commands) and a Dashboard (for monitoring purposes). From an architectural point of view the ESGF-Desktop can be considered the GUI of the Dashboard.

Most part of this work has been carried out during the 1-month visit at LLNL/PCMDI (from April 13 to May 13) by Dr. S. Fiore from CMCC. The visit has been supported by the IS-ENES project.

# 4      Standards and interface specifications

IS-ENES has strongly supported modular development of components in the ESGF infrastructure and adoption of well defined interfaces between components in order to enable parallel development of systems. The pressures of CMIP5 deadlines and diverse priorities of partners in the global federation have prevented some of these documents from reaching maturity.

## CMIP5 Data Reference Syntax

A well structured and detailed specification of the data format for contributions to the CMIP5 archive underpins the whole system. The document is authored by Karl Taylor (PCMDI), but IS-ENES partners contributed significantly to its formulation.

http://cmip-pcmdi.llnl.gov/cmip5/docs/cmip5_data_reference_syntax.pdf

## CF standard names

The CF Standard name table defines the names used for variables in the CMIP5 archive, and includes a precise definitions to avoid ambiguities arising from different usages in different modelling centres. STFC moderates the global discussions leading to the adoption of names. Over 800 variable names were required for CMIP5.

http://cf-pcmdi.llnl.gov/documents/cf-standard-names/standard-name-table/22/cf-standard-name-table.html

## ESGF Security Interface Control Document

A clearly defined security architecture is critical to the smooth functioning of the distributed archive. IS-ENES contributed to the interface specifications:

http://esg-pcmdi.llnl.gov/esgf/esgf-security-interface-control-documents/esgf-security-interface-control-document-1.0/esgfInterfaceControlDocument-0.8.pdf

**ESGF THREDDS XML profile**

The ESGF THREDDS profile is critical to the stability and interoperability of the distributed system.

https://is-enes-wiki.dkrz.de/farm/SA2/JRA4%20Collaboration?action=AttachFile&do=get&target=ESGF_THREDDS_profile-v02.pdf

# 5 Summary

The IS-ENES data services package will link European climate data providers into a global network managed by the Earth System Grid Federation (ESGF). The infrastructure is designed to facilitate integration with nationally funded services so as to gain the added value of international integration without impairing existing service provision. IS-ENES has contributed to the development of software, both directly through design and delivery of packages and through participation in testing and evaluation of the infrastructure.

# 6 Outlook

The infrastructure deployed in IS-ENES has set some important benchmarks and established a distributed system which has successfully delivered large data volumes to a broad scientific community. The core system of standardised terminology to describe data, standard catalogue formats supporting access to files and subsets of files, and a common security infrastructure is robust. A number of key problems are discussed briefly here.

## 6.1 Governance

The ENES distributed archive is part of a globally federated system. The IS-ENES project has provided a clear management structure for the European components of the global system, but managing priorities with partners outside Europe has been problematic, despite the clear shared priorities in delivering the CMIP5 archive. Efforts to establish a clear global governance structure for the emerging archive system are being pushed forward by IS-ENES partners and will be further supported by IS-ENES2. A white paper on governance of the distributed archive has been produced (http://home.badc.rl.ac.uk/mjuckes/documents/ENES_data_federation_wp_22Feb2013.pdf )

## 6.2 Documentation

The CMIP5 archive has extensive model documentation, but fragility in the initial systems used to collect information from the modelling groups (the CMIP5 questionnaire), which were run in a quasi-operational mode to meet the deadlines imposed by the CMIP5 process at a time when development was still fluid, led to some frustration in the data provider community. At the time that CORDEX data requirements were being discussed, there were no resources available for a substantial revision of the CMIP5 questionnaire and the CORDEX community made no commitment to provide a similar level of documentation. Lack of detailed documentation may lead to frustration among users and limit the impact of the CORDEX data archive, but the decision reflects the feeling that encouraging groups to contribute data is, in this first coordinated regional downscaling experiment, the priority. The UK PIMMS project has made some progress in adapting the CMIP5 questionnaire for a wider data provider community, and adding flexibility to avoid the frustrations that were common during the CMIP5 meta-data entry phase. Steps should be taken to encourage CORDEX data providers to use this tool.

## 6.3 Quality Control

There is broad agreement that quality control of data is needed, but little common ground on how to manage quality control of millions of files in a globally distributed archive. IS-ENES set some new standards and created a framework for carrying out tests and linking the results into the archive so that they are clearly visible to users.

# 7 ANNEX 1: Foreground

## 7.1 *Software packages and scripts*

Version control of integrated archive system is incomplete, but all components and tools developed by IS-ENES have well defined versions.

| EnesGetCert | |
|---|---|
| Category | User tool; bash script |
| Location | |
| Documentation | |
| Licence | |
| Dependencies | Java version 1.5 or greater; |

| MyProxyClient | |
|---|---|
| Category | User tool; python library and script |
| Location | |
| Documentation | |
| Licence | |

| CORDEXwritter | |
|---|---|
| Category | User tool; bash script collection |
| Location | http://exporter.nsc.liu.se/864ebccfc365436d872467a34bf1b707 |
| Documentation | In tar file at above location |
| Licence | Open source |

| Dashboard | |
|---|---|
| Category | ESGF component |
| Location | |
| Documentation | |
| Licence | |

## 7.2 Documentation

| CMIP5 Data Reference Syntax | |
|---|---|
| Category | CMIP5 data protocol |
| Location | http://cmip-pcmdi.llnl.gov/cmip5/docs/cmip5_data_reference_syntax.pdf |
| Publication date | 9 March, 2011 |
| Authors | Karl E. Taylor, V. Balaji, Steve Hankin, Martin Juckes, Bryan Lawrence, and Stephen Pascoe |

| Security interface | |
|---|---|
| Category | ESGF protocol |
| Location | http://esg-pcmdi.llnl.gov/esgf/esgf-security-interface-control-documents/esgf-security-interface-control-document-1.0/esgfInterfaceControlDocument-0.8.pdf |
| Publication date | 14 April, 2010 |
| Authors | Rachana Ananthakrishnan, Luca Cinquini, Phil Kershaw, Neill Miller. |

| ESGF THREDDS profile | |
|---|---|
| Category | ESGF protocol |
| Location | https://is-enes-wiki.dkrz.de/farm/SA2/JRA4%20Collaboration?action=AttachFile&do=get&target=ESGF_THREDDS_profile-v02.pdf |
| Publication date | 20 February, 2013 |
| Authors | Stephen Pascoe |

| CORDEX ESGF facet mappings | |
|---|---|
| Category | ESGF protocol |
| Location | http://home.badc.rl.ac.uk/mjuckes/isenes/docs/cordexGatewayFacets_v1.pdf |
| Publication date | December 2011 |
| Authors | Martin Juckes |

# 8 Glossary

| Term | Description |
|---|---|
| CF standard names | Naming convention for climate data – http://cf-pcmdi.llnl.gov/documents/cf-standard-names/standard-name-table/22/cf-standard-name-table.html |
| CMIP5 | Coupled Model Intercomparison Project, Phase 5 – http://cmip-pcmdi.llnl.gov/cmip5/ |
| CMOR | Climate Model Output Re-writer – software library – http://www2-pcmdi.llnl.gov/cmor |
| CORDEX | Coordinated Regional Climate Downscaling Experiment – http://wcrp-cordex.ipsl.jussieu.fr/ |
| ESGF | Earth System Grid Federation |
| GLOBUS | Opensource grid software – http://www.globus.org/ |
| GRIB | Binary file format – http://www.wmo.int/pages/prog/www/DPS/FM92-GRIB2-11-2003.pdf |
| IPR | Intellectual Property Rights |
| MD5 | Cryptographic hash function – http://tools.ietf.org/html/rfc1321 |
| METAFOR | Common metadata for climate modelling digital repositories, FP7-INFRASTRUCTURES project 211753 |
| MIP tables | Model Intercomparison Project tables maintained by PCMDI |
| NCAR | US National Centre for Atmospheric Research |
| NetCDF | Network Common Data Format - http://www.unidata.ucar.edu/software/netcdf/ |
| OPeNDAP | Software library – http://www.unidata.ucar.edu/software/netcdf/ |
| OpenID | Single sign on protocol – http://openid.net/ |
| P2P | Peer-to-peer architecture of the ESGF software stack |
| PCMDI | Program for climate model diagnostics and intercomparison – http://www-pcmdi.llnl.gov/ |
| PIMMS | Portable Infrastructure for the Metafor Metadata System – http://proj.badc.rl.ac.uk/pimms |
| Sudo | Linux command |
| THREDDS | Data distribution package – http://www.unidata.ucar.edu/projects/THREDDS/ |
| URL | Uniform Resource Locator |
| WCRP | World Climate Research Program – http://www.wcrp-climate.org/ |
| X.509 | Public key infrastructure standard – http://www.ietf.org/rfc/rfc2459.txt |