# The Earth System Grid:
# Enabling Access to Multi-Model Climate Simulation Data

**The Earth System Grid Center for Enabling Technologies Team:**

**D N Williams[3,9], R Ananthakrishnan[1], D E Bernholdt[7,9], S Bharathi[8], D Brown[5], M Chen[7], A L Chervenak[8], L Cinquini[5], R Drach[3], I T Foster[1,9], P Fox[5], D Fraser[1], S Hankin[6], P Jones[4], C Kesselman[8], D E Middleton[5,9], J Schwidder[7], R Schweitzer[6], R Schuler[8], A Shoshani[2], F Siebenlist[1], A Sim[2], W G Strand[5], N. Wilhelmi[5]**

[1] Argonne National Laboratory, Argonne, IL, USA.

[2] Lawrence Berkeley National Laboratory, Berkeley, CA, USA.

[3] Lawrence Livermore National Laboratory, Livermore, CA, USA.

[4] Los Alamos National Laboratory, Los Alamos, NM, USA.

[5] National Center for Atmospheric Research, Boulder, CO, USA.

[6] National Oceanic and Atmospheric Administration (PMEL), Seattle, WA, USA.

[7] Oak Ridge National Laboratory, Oak Ridge, TN, USA.

[8] University of Southern California, Information Sciences Institute, Marina del Ray, CA, USA.

[9] E-mail: williams13@llnl.gov, don@ucar.edu, itf@mcs.anl.gov, bernholdtde@ornl.gov

## Abstract

The Earth System Grid (ESG) project is creating a data sharing environment that links international climate research centers and provides a range of users with model-generated simulations. ESG leverages disparate technologies to manage data at distributed sites in a manner that creates a unified virtual environment. By transforming distributed climate simulation data into a collaborative community resource, ESG is changing the way global climate research is conducted.

Since its production launch in 2004, ESG's most notable accomplishment was to supply climate simulation data generated by a number of distinct models to scores of scientists who contributed to the recent Intergovernmental Panel on Climate Change (IPCC) Fourth Assessment Report (A4). Other national and international stakeholders such as the Community Climate System Model developers, the Climate Science Computational End Station, and the North American Regional Climate Change Assessment Program also have endorsed ESG technologies for disseminating climate data to their respective user communities. In coming years, the recently created Earth System Grid Center for Enabling Technology (ESG-CET), will lead the effort to extend existing ESG methods to assist international scientists in their endeavors to understand the climate system and to predict its potential future change.

Climate scientists have a wide variety of practical needs, but the ability to efficiently access and manipulate data is an overarching problem. Researchers are increasingly required to access large complex datasets that are archived in different formats on disparate platforms scattered around the world—and to extract pieces of datasets to perform statistical or model diagnostic metrics "in place." A common virtual environment that will allow easy access to large climate model datasets and tools is keenly needed. The software infrastructure that supports this environment must not only provide access to climate model data, but also facilitate the use of visualization software, diagnostic algorithms, and related resources.

To this end, the U.S. Department of Energy (DOE) Scientific Discovery through Advanced Computing (SciDAC)-2 program established the Earth System Grid Center for Enabling Technologies (ESG-CET) (FIG. 1). ESG-CET is working to advance climate science by utilizing and developing computational resources and technologies for accessing and managing model data that are physically located in distributed multi-platform archives.

Fig 1 here

**ESG Impact on the Climate Community**

The Earth System Grid (ESG) seeks to address the challenges of climate science. Its infrastructure improves research efficiency by enabling rapid browsing of, access to, and analysis of, even the largest climate datasets. ESG's current data holdings include: the extensive Intergovernmental Panel on Climate Change (IPCC) Fourth Assessment Report (AR4) archive (now known as the Coupled Model Intercomparison Project 3 (CMIP3) multi-model database), the Community Climate System Model (CCSM) archive, and more recently the CCSM

Biogeochemistry (BGC) Carbon-Land Model Intercomparison Project (C-LAMP) archive. ESG also serves data from Parallel Climate Model (PCM) and Parallel Ocean Program (POP) simulations, and from the Cloud Feedback Model Intercomparison Project (CFMIP). Since ESG's production launch in 2004, users have downloaded some 1 million files containing a total of more than 350 terabytes of data. (1 terabyte (TB) = $10^{12}$ bytes.)

As part of "publishing" their data into an archive, model data providers are given access to ESG catalog metadata—information describing the data—so that they can register appropriate metadata for their data. This allows data providers to manage all information related to generating, defining, archiving, and retrieving model simulation runs. Providers may also restrict access to their data, according to a variety of criteria. ESG's long-term goal is to tie the metadata ingestion process to climate simulation workflow, so that model-simulation metadata can be added automatically to the ESG data holdings, thereby expediting the data publishing process and minimizing processing errors.

To illustrate the data publication process, we describe here the process followed for the CMIP3 archive, which contains the output of multiple coupled ocean-atmosphere model simulations. This process was coordinated by the Program for Climate Model Diagnosis and Intercomparison (PCMDI). PCMDI, which is located at the Lawrence Livermore National Laboratory, built on its considerable experience in managing data from international climate model intercomparison experiments.

Some 20 participating modelling centers performed numerous control and climate-change scenario experiments prescribed by the IPCC (Meehl et al. 2007). Before transmitting their data to PCMDI, modelling groups transformed their model output into the requested network Common Data Form (netCDF) format and Climate and Forecast (CF) metadata convention*. To facilitate this conversion, PCMDI provided participating modelling centers with the Climate Model Output Rewriter (CMOR, pronounced "Seymour") software, which produces CF-compliant netCDF files.

Once model data conversion was complete, the modelling centers sent their simulation data to PCMDI on 1 TB disks. After the data passed preliminary quality-control checks, ESG tools were used to publish data into the CMIP3 multi-model archive. Since the completion of data publishing in 2005, approved users have been able to access these data in a variety of ways, including via File Transfer Protocol (FTP), the ESG data portal, the Live Access Server (LAS), and the Open-source Project for a Network Data Access Protocol (OPeNDAP).

The ESG data archives are freely available for access. While mainly utilized by research scientists, its user community also includes educators, students, and employees of both domestic and foreign governments. The analysis and interpretation of ESG data have lead to the production of hundreds of research publications and impact reports (http://www-pcmdi.llnl.gov/ipcc/subproject_publications.php). By facilitating such wide data access, ESG has provided a venue at which scientists, engineers, and other have offered concrete suggestions

leading to the enhancement of the scientific accuracy, portability, and performance of current-generation climate models.

Responding to user experience gained through the distribution of data thus far, ESG development continues. For example, ESG now offers multi-level searches and analyses—a capability that can reduce network traffic.

---

*NetCDF is a set of software libraries and machine-independent data formats that support the creation, access, and sharing of array-oriented scientific data. The CF conventions for climate and forecast metadata are designed to promote the processing and sharing of files created with the netCDF API.

**History of ESG**

Work on ESG began in the year 2000 with the "Prototyping an Earth System Grid" (ESG I) project, supported by the DOE. In this prototype phase, ESG developed Data Grid technologies for managing the movement and replication of large datasets, and applied these to an ESG-enabled data browser based on the PCMDI Climate Data Analysis Tools (CDAT). At a 2001 Supercomputing Conference, ESG demonstrated the potential for remotely accessing and analyzing climate data located at several sites across the U.S., and achieved cross-country transfer rates of more than 500 Mbit/s (Chervenak et al. 2003).

While the ESG I prototype provided a proof of concept, the 2002 SciDAC funded Earth System Grid II project, "Turning Climate Datasets into Community Resources" (Bernholdt et al. 2005), made this concept a reality. These efforts targeted the development of technologies pertinent to metadata (e.g., metadata extraction from netCDF files and catalogue services), security (web-based user registration and authentication), data transport (GridFTP-enabled "OPeNDAP-g" for dataset aggregation, high-performance access with subsetting, file transport and integration with mass storage systems), and web portal technologies to provide interactive access to climate data holdings. (OPeNDAP provides software that makes data accessible to remote applications without file downloads—independent of the storage format.)

In its effort to become an integral part of the work of climate scientists (Bernholdt et al. 2005), ESG began to distribute CCSM and PCM model data in mid-2004. This first production system

brought major advances in model archiving, data management, and sharing of distributed climate

data. Originally distributed among three sites (i.e., NCAR, LBNL, and ORNL), this system now

supports some 6,000 registered international users and manages over 160 TB of data. The first

ESG architecture design (FIG. 2) integrated a wide-range of Grid and standard information

technology (IT) tools.

The ESG portal is the main access point to the system. The portal provides a central location for

authentication, authorization and accounting services and brokers user data requests among the

distributed data nodes.  The portal also provides the interface through which authorized

providers publish data. A Replica Location Service (RLS) server at each data node indexes the

physical files (or "replicas") available at that site. Online files are served via an LAHFS server

(Lightweight Authorized HTTP File Server), while files on deep storage are requested and served

via the SRM (Storage Resource Manager), which retrieves them from the archive and transfers

them to a central disk cache, where they are made available by another LAHFS server. All ESG

system components are continuously monitored, and system administrators and users are

notified whenever a service becomes unavailable.

In some cases, individual files or groups of files are too large to be transferred via the ESG portal.

For such cases, ESG developed a DataMover tool to effect robust large-scale data movement by

interacting with SRMs and replicating thousands of files between specified mass storage

systems.  A client-side version of this tool, DataMover-Lite (DML), automates multi-file

requests for data transfers from SRMs into the clients' file systems. The DML's user-friendly

interface allows easy monitoring of file transfers to a client machine (FIG. 2**)**.

By late 2004 at PCMDI, ESG began distribution of climate model data relevant for the IPCC's

Fourth Assessment Report (AR4). Subsequently, this database was designated as the Climate

Model Intercomparison Project, Phase 3 (CMIP3) multi-model archive to emphasize its

continuity with some 15 years of international efforts (e.g. the precursor CMIP1 and CMIP2

coupled model intercomparisons) in coupled climate model simulations. The CMIP3 designation

also emphasizes PCMDI's commitment to extend these holdings to include additional data from

current-generation coupled climate models (Meehl et al. 2007). As a result of this extensive data-

management experience, ESG developed useful technologies to provide climate scientists with

virtual access to distributed data resources.

With its combined data delivery strategy (via FTP, ESG web portal, OPeNDAP, and LAS), the

CMIP3 archive has distributed more than 300 TB of data to users. The data portal now supports

some 1,500 registered analysis projects and manages over 35 TB of data (~80,000 files).  More

than 400 peer-reviewed publications have been authored on analyses of the CMIP3 data.

In late 2006, ESG entered a new phase as the ESG Center for Enabling Technologies (ESG-

CET), with funding from SciDAC-2.  The primary goal of ESG-CET is to generalize the existing

system to support a more international, broadly distributed, and diverse collection of archive

sites and data types. A secondary goal is to extend ESG beyond access to raw data, by developing server-side capabilities that allow a user to conduct common analysis procedures on data where it physically resides before downloading the derived products to the client's site. ESG-CET views such capabilities as essential if more people are to use of the petabytes (1 petabyte (PB) = $1x10^{15}$ bytes) of climate data that are potentially available. Thus, ESG-CET intends to develop the "petascale" data-access capabilities that are rapidly becoming necessary for climate scientists.

New to ESG in 2007 is C-LAMP intercomparison data from different land biogeochemistry (BGC) schemes embedded in the CCSM coupled ocean-atmosphere climate model (Hoffman et al. 2007). The C-LAMP experimental output data now are being archived on an ORNL site that is modelled after the ESG CMIP3 database.

ESG continues to help the climate community manage a rapidly growing data environment. In 2007, the North American Regional Climate Change Assessment Program (NARCCAP) endorsed ESG methods for disseminating high-resolution regional climate model data. Archived at both PCMDI and NCAR, these data will be distributed through NCAR's ESG portal. Registration administrators will be required for access.

**Future Directions**

In coming years, the ESG-CET will scale up existing ESG capabilities to meet the needs of several ambitious scientific projects:

- The Coupled Model Intercomparison Project, Phase 5 (CMIP5) for scientists contributing to the IPCC Fifth Assessment Report (AR5)

- The Computational Climate End Station (CCES) at the DOE Leadership Computing Facility at Oak Ridge National Laboratory (ORNL), and

- Other wide-ranging climate model evaluation activities.

These projects will provide a focus for ESG technologies and their future development. ESG-CET will seek to connect a large number of users with geographically distributed climate model archives via client-server infrastructure, and to provide them with advanced data analysis tools.

In addition, ESG will broaden to include new types of model data (e.g., biogeochemistry, dynamic vegetation, etc.), to provide more powerful (server-side) access and analysis services, to enhance interoperability among commonly used climate analysis tools, and to enable end-to-end simulation and analysis workflow (FIG. 3).

Fig 3 here

*Collaborators*

Future ESG-CET activities also will be framed by relationships with other institutions that share common data-management interests.  For example, members of the Global Organization for Earth System Science Portal (GO-ESSP) are developing a software infrastructure for discovery,

acquisition and analysis of climate model data. The GO-ESSP members who have agreed to take leading roles as gateways and/or nodes in the CMIP5 (IPCC AR5) testbed are the Program for Climate Model Diagnostic Intercomparison (PCMDI), the National Center for Atmospheric Research (NCAR), the Oak Ridge National Laboratory (ORNL),and the Los Alamos National Laboratory (LANL).

Other GO-ESSP members that will play a vital role in this include: the Geophysical Fluid Dynamics Laboratory (GFDL), the British Atmospheric Data Centre (BADC), the World Data Center for Climate (WDCC), and the University of Tokyo Center for Climate System Research.

Because this consortium will extend beyond U.S.-bound partnerships, it may require additional software deployment to accommodate components from the U.K. NERC DataGrid (NDG), the European Union (EU) MetaFor project, and the German C3-Grid initiative.

### *Future usage*

Commonly, under the current ESG system, a user first uses hyperlinks through a web browser, then queries a database from a specific archive to retrieve desired records. The user then can retrieve desired data via the ESG data portal, a DataMover-Lite tool, or a "Get" operation in a web browser (wget). Once data are downloaded to the user's site, data regridding, reduction, and diagnosis often follow.

Under the current ESG system, this process often requires many data movements and places a substantial burden on network, storage, and computing resources. With the next-generation ESG architecture, the user instead will browse, search, and discover (determine the existence, presence, or properties of) distributed data on remote sites. The user then either can request these raw data, or can regrid and analyze them, *where the physical data actually reside*, before downloading. This approach will reduce network traffic and allow the researcher to focus primarily on science, rather than on the organization and movement of data (FIG. 4). (The approach also places new demands on data hosting sites, which we must manage.)

Fig 4 here

Initial ESG efforts focused primarily on web portal-based access to climate data. In future ESG-CET work, the existing web-portal capabilities will be augmented with a strong emphasis on applications that can tap new services to streamline data download, as well as provide powerful high-level analysis and visualization capabilities on the user's platform. Future ESG services will allow popular climate analysis and visualization tools (e.g., CDAT, NCL, GrADS, Ferrret, IDL, Matlab, etc.) to be used directly within the system. The user interface for DataMover-Lite (DML) will also be developed further.

*Functional specification and architecture design*

In order to meet the petascale needs of the climate community, the ESG-CET architecture must allow a large number of distributed sites with varying capabilities to federate and/or work as stand-alone entities. Federation implies a virtual trust relationship so that users authenticate once

to gain access to data across multiple systems and organizations. To accomplish this goal, the future ESG-CET architecture will be based on three tiers of data services (FIG. 5).

Tier 1 services provide shared functionality across the overall ESG-CET federation include user registration and management, common access control metadata services, a common set of notification services about data changes, and global monitoring services for detecting problems.

All ESG-CET sites share a common user database. (Access to specific data collections and related resources, such as IPCC data, will need to be approved by the data "owners.") A user should be able to find data of interest throughout the whole federation, independently of the site where the user initiates a data search.

Tier 2 comprises a limited number of ESG gateways to broker requests for some or all of the data registered with ESG. Having multiple gateways allows for gateways that provide access to only a subset of data—for example, a CMIP5 gateway—and also supplies fault-tolerance to the overall system. Gateway-deployed services include the user interface for searching and browsing metadata, for requesting data products (including analysis and visualization tools), and for orchestrating complex workflows. Because the software deployed at gateways is likely to require considerable expertise to maintain, these gateways probably will be operated directly by ESG-CET engineering staff.

Tier 3 includes the actual data holdings and the back-end services used to access data, which reside on a potentially large number of federated ESG Nodes. Tier 3 resources typically host the services needed to publish data to ESG, and to execute data-product requests formulated through an ESG gateway. Because researchers and engineers at local institutions with varying levels of expertise operate ESG nodes, the deployed software stack is kept to a minimum and is supplied with detailed and exhaustive documentation. A single ESG gateway may serve data requests to many associated ESG nodes: for example, as part of the CMIP5 (IPCC AR5) database, more than 20 institutions are expected to operate ESG data nodes.

### *ESG-CET future component design*

The component layers of ESG-CET help solve the challenges faced by petascale archives. In some cases, these approaches are innovative and have been adopted by other scientific application areas. An in-depth description of the multiple technologies comprising the ESG-CET component design follows.

### *Metadata*

The design of the metadata database is at the heart of ESG, since the metadata model underlies other major ESG components, and especially the search and browse facilities and the publishing system. Thus, this is one of the first tasks that ESG-CET has addressed.

The current software focuses on gridded datasets generated from climate model experiments. The next version will retain that focus, while anticipating the need to expand the scope of data served

to related subject areas, such as the assimilation of observations and the prediction of climate-change impacts via models. Similarly, ESG will continue to add value by supporting derived and virtual datasets. (Virtual datasets have all the properties of regular datasets except, by definition, they have no location information. Derived datasets are data products resulting from some type of transformation of one or more datasets.)

In addition, ESG is beginning to design a new search capability based on the concept of "faceted classification" (Adkisson 2003) that will provide several important features. At a given point, the user will see search terms and categories that apply within the current context, and will be able to avoid queries that return empty result sets. Similarly, organizing metadata around facets will provide important flexibility, since the categories can be changed and updated without impacting the rest of the system.

ESG is working with related metadata projects to ensure consistency with emerging community standards. For example, members of the Earth System Curator (ESC) and MetaFor projects are participating in the design process, and ESG is exploring how the respective schemas intersect. Because both ESC and MetaFor emphasize the viewpoint of the data producer, they have developed schemas that allow a rich description of the structure of models and model components; in contrast, ESG takes the viewpoint of the end user, who is typically more concerned with the scientific aspects of the simulation. The union of these data models thus will provide a richer and more comprehensive database, and will ensure that ESG can interface with software systems derived from ESC.

*Federated metadata*

For ESG-CET, a major challenge in realizing a global, petascale, and distributed architecture is the design of a federated system that allows participating sites to publish data sets and their associated metadata.

ESG-CET's vision for the common architecture includes a single master metadata catalogue that is hosted at the Global Services (Tier 1) layer of ESG (FIG. 5). In practice, this service may be deployed at a particular ESG gateway (e.g., at NCAR, ORNL, or PCMDI), but should be gateway-independent. All updates to metadata must be performed on this master catalogue; its contents are updated periodically to replicas maintained at each ESG gateway node. This feature allows users to issue metadata queries at any gateway. Having multiple updated catalogues to answer queries provides load balancing: the master catalogue does not become a bottleneck. Also, even if the master catalogue is unavailable, queries can be satisfied at any gateway using the updated catalogues.

ESG also envisions metadata "harvesting," by a federated catalogue. In this architecture, each ESG data node (Tier 3) will generate metadata for the data sets that it publishes, and then will store this metadata in a local catalogue. The master metadata catalogue may harvest this local metadata.

*Security*

Secure access to data and resources is crucial, but must not create an undue burden for data users and administrators. An emerging security requirement is for browser and client Single Sign On (SSO) to gain access to data. The SSO will allow the ESG portal functionality to be split among multiple servers while enforcing authenticated access by the browser or client. Moreover, in future years it will allow ESG to leverage the identity management systems of selected partner organizations so that users will only have to authenticate within their home domain.

Existing open source technology for enabling Single Sign On is relatively easy to use and to integrate within the ESG infrastructure. Extending the architecture to include credential repositories as well as centralized security services will provide an interface so that clients can operate easily and securely within ESG. However, alternative technologies may still be necessary to maintain compatibility with other resources.

*Product services*

ESG-CET serves customers with a broad range of sophistication. These users range from numerical modellers who want "raw" model output to users who only want to quickly visualize the gross features of model output.

Petascale ESG data holdings require substantial server-side reduction in order to satisfy certain users – both through straightforward subsetting and dissemination and through averaging and other analysis operations.

ESG has developed usage scenarios to better understand optimal experiences for different users. The Live Access Server (LAS), for example, is being adapted as a generalized workflow engine for converting raw data into analysis products and visualizations. The architecture will be adaptable to the range of users and products. ESG can generate a wide range of visualization types (1D, 2D, and 3D visualizations with customized contour levels, palettes, etc.) and basic statistical operations, such as determining extrema. The analysis capabilities have been incorporated into beginning-level server-side (i.e., remote) functionality that is accessible through external protocols. ESG-CET is working with another SciDAC-2 Center, the Center for Enabling Distributed Petascale Science (CEDPS) (Baranovski et al. 2007) to develop methods for processing such requests in a parallel computing environment.

*Querying and browsing the catalogue*

An interface is under development in which ESG-CET's web portal and client applications can query, browse, and discover data from ESG's distributed data archives. Some ideas for querying the data are coming from the Open Geospatial Consortium (OGC) and search results will be returned as THREDDS. The OGC is an international industry consortium developing open-source interface specifications. THREDDS specifies a tree-like data structure, which can represent a hierarchy of datasets by use of the Extensible Mark-up Language (XML).

The current ESG Portal uses THREDDS XML with Extensible Stylesheet Language Transformations (XSLT) to render the portal interface in Hypertext Markup Language (HTML). Unfortunately, when data files are published, static THREDDS XML is produced, making the

catalogue structure rigid. A more flexible approach would be to dynamically generate THREDDS

XML with data organized in different hierarchies. to allow users to browse in different ways (by

experiment, by variable, by model, etc.).

*User interface*

Interactive user access to ESG-CET's data archive will probably come by way of a web browser

and/or an application-based interface. In addition to basic improvements in the ESG portal, we

are also exploring how to integrate the Live Access Server (LAS) user interface into the ESG

portal.


Some of the new technology libraries might offer interfaces that would allow for the injection of

more dynamic elements without the need to change the ESG portal framework. In view of these

possibilities, ESG plans to reuse as much of the existing code as possible. However as the design

for the new interface evolves, the current approach will need to be reevaluated.

*Publishing*

To publish ESG data, a provider needs to stage a set of model runs, extract the necessary

metadata, and run checks to ensure that the metadata meet both ESG and project-specific

standards. Metadata that are not extractable from the raw model data must be added manually.

The data provider also must specify access privileges for the datasets and must authenticate with

a gateway to obtain publishing privileges. Every published object (file or other aggregation) gets a

unique identifier that remains with the object if it is replicated or moved to a different site. This

entire publishing process can be scripted, and the provenance of the data (origin and history of

subsequent owners) can be fully documented.

**Summary**

The ESG-CET team now is designing new services that will expose users to products and services that can derive data from spinning disk or tertiary storage, migrate data from one location to another, and enable a number of workflow capabilities. ESG-CET intends for its approaches and technologies to also impact other SciDAC-2 scientific application areas that entail petascale data management.

**Acknowledgements**

**References**

Adkisson, Heidi P. 2003: Use of faceted classification,
http://www.webdesignpractices.com/navigation/facets.htm (accessed 2 November 2003).

Allcock, B., Bresnahan, J., Kettimuthu, R., Link, M., Dumitrescu, C., Raicu, I. and Foster, I.
The Globus Striped GridFTP Framework and Server. *Supercomputing 2005*, 2005.

Ananthakrishnan, R., Bernholdt, D.E., Bharathi, S., Brown, D., Chen, M., Chervenak, A.L.,
Cinquini, L., Drach, R., Foster, I.T., Fox, P., Fraser, D., Halliday, K., Hankin, S., Jones,
P., Kesselman, C., Middleton, D.E., Schwidder, J., Schweitzer, R., Schuler, R., Shoshani,
A., Siebenlist, F., Sim, A., Strand, W.G., Wilhelmi, N., Su, M. and Williams, D.N.,
Building a global federation system for climate change research: the earth system grid
center for enabling technologies (ESG-CET). *Journal of Physics: Conference Series*, 78.
2007.

Baranovski, A., Bharathi, S., Bresnahan, J., Chervenak, A., Foster, I., Fraser, D., Freeman, T.,
Gunter, D., Jackson, K., Keahey, K., Kesselman, C., Konerding, D.E., Leroy, N., Link,
M., Livny, M., Miller, N., Miller, R., Oleynik, G., Pearlman, L., Schopf, J.M., Schuler,
R. and Tierney, B. Enabling distributed petascale science. Journal of Physics: Conference
Series, 78. 2007.

Bernholdt, D., Bharathi, S., Brown, D., Chanchio, K., Chen, M., Chervenak, A., Cinquini, L.,
Drach, B., Foster, I., Fox, P., Garcia, J., Kesselman, C., Markel, R., Middleton, D.,
Nefedova, V., Pouchard, L., Shoshani, A., Sim, A., Strand, G. and Williams, D., 2005: The
Earth System Grid: Supporting the next generation of climate modeling research.
*Proceedings of the IEEE*, **93** (3). 485-495.

Chervenak, A., Deelman, E., Foster, I., Guy, L., Hoschek, W., Iamnitchi, A., Kesselman, C., Kunst, P., Ripenu, M., Schwartzkopf, B., Stockinger, H., Stockinger, K. and Tierney, B. Giggle: A Framework for Constructing Scalable Replica Location Services SC2002: High Performance Networking and Computing. http://www.globus.org/research/papers.html#giggle, 2002.

Chervenak, A., Schopf, J.M., Pearlman, L., Su, M.-H., Bharathi, S., Cinquini, L., D'Arcy, M., Miller, N. and Bernholdt, D. Monitoring the Earth System Grid with MDS4. 2nd IEEE Intl. *Conference on e-Science and Grid Computing (e-Science 2006),* Amsterdam, Netherlands, 2006.

Chervenak, A., et al., 2003: High-performance remote access to climate simulation data: A challenge problem for data grid technologies. *Parallel Computing*, **29** (10). 1335-1356.

Committee on Archiving and Accessing Environmental and Geospatial Data at NOOA, Environmental Data Management at NOAA: Archiving, Stewardship, and Access. *The National Academies Press*, Washington, D.C., http://www.nap.edu, 2007.

Foster, I., Globus Toolkit Version 4: Software for Service-Oriented Systems. *Journal of Computational Science and Technology*, 21 (4). 523-530. 2006.

Gerald A. Meehl, Curt Covey, Thomas Delworth, Mojib Latif, Bryant McAvaney, John F. B. Mitchell, Ronald J. Stouffer, and Karl E. Taylor, 2007: The WCRP CMIP3 multi-model dataset: A new era in climate change research. *Bulletin of the American Meteorological Society*, **88** (9), 1383-1394.

Hoffman, Forrest M., Curtis C. Covey, Inez Y. Fung, James T. Randerson, Peter E. Thornton, Yen-Huei Lee, Nan A. Rosenbloom, Reto C. Stöckli, Steven W. Running, David E.

Bernholdt, and Dean N. Williams, 2007: Results from the Carbon-Land Model

Intercomparison Project (C-LAMP) and availability of the data on the Earth System Grid

(ESG). *Journal of Physics: Conference Series*, **78**. doi:10.1088/1742-6596/78/1/012026.

Shoshani, A., Sim, A. and Gu, J., Storage Resource Managers: Essential Components for the

Grid. *Kluwer Academic Publishers*, 2003.

**Figure Captions**

FIG. 1. The ESG-CET collaboration includes participation from Argonne National Laboratory (ANL), Los Alamos National Laboratory (LANL), Lawrence Berkeley National Laboratory (LBNL), Lawrence Livermore National Laboratory (LLNL), National Center for Atmospheric Research (NCAR), Oak Ridge National Laboratory (ORNL), Pacific Marine Environmental Laboratory (PMEL), and the University of Southern California (USC) Information Sciences Institute.

FIG. 2. Schematic of the first-generation ESG architecture showing the U.S. repositories. Climate model data is located on deep archives at the Lawrence Berkeley National Laboratory (LBNL) National Energy Research Scientific Computing Center (NERSC), the National Center for Atmospheric Research's (NCAR) Mass Storage System (MSS), and the Oak Ridge National Laboratory (ORNL) High Performance Storage System (HPSS), and on the Los Alamos National Laboratory (LANL) fast-access rotating disks. Also depicted is a provider publishing data by means of a web browser and a data user accessing published data via either a web browser or DataMover-Lite (DML).

FIG. 3. A high-level ESG-CET "roadmap" showing the planned evolution of the ESG system from terascale to petascale data management. Also shown are the scientific data management and analysis requirements in relationship to the ESG development timeframe. Note that a distributed testbed for CMIP5 (IPCC AR5) needs to be in place by early 2009.
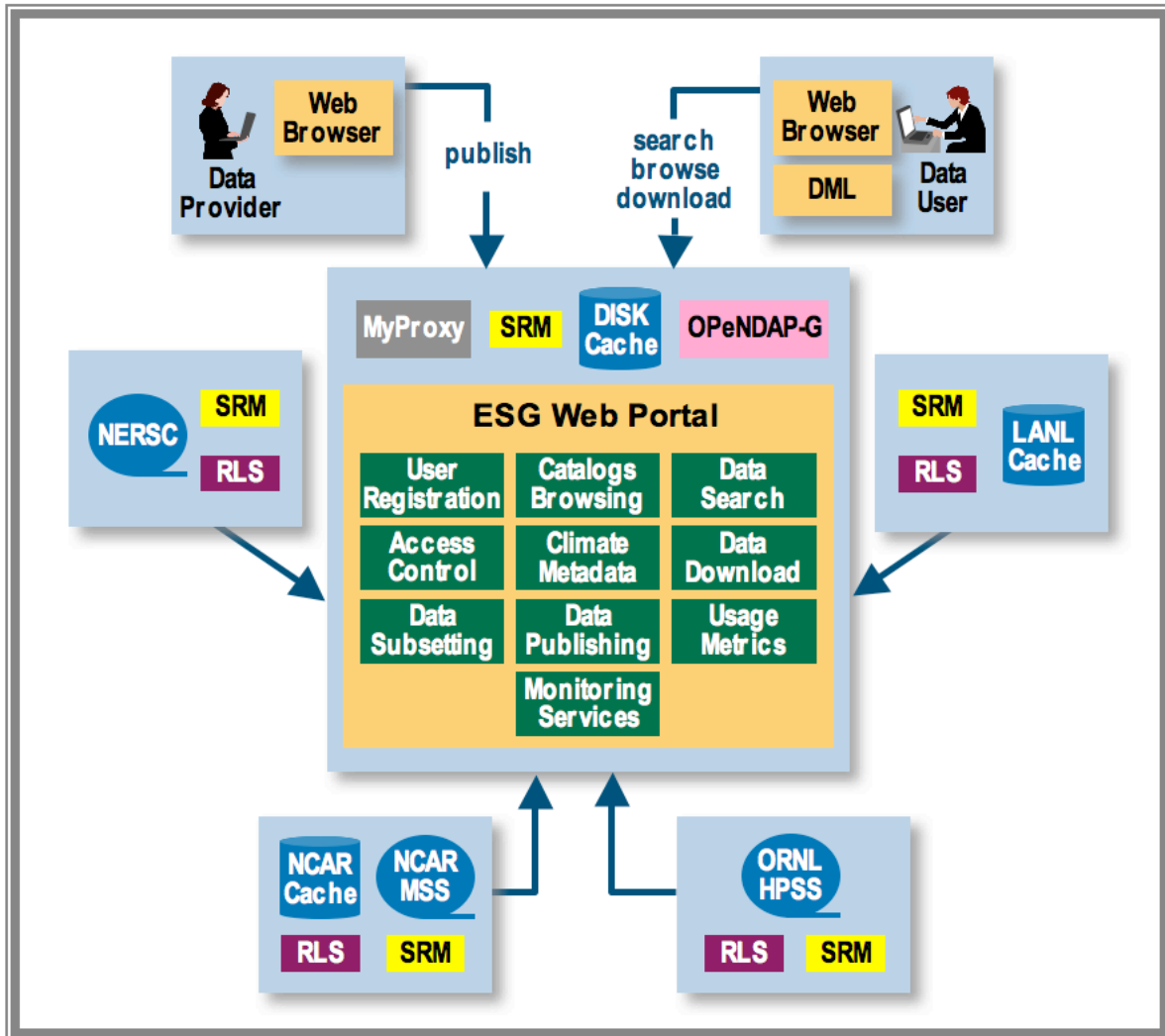
FIG. 4. In this example a user searches the ESG-CET portal from a remote site. Required data are found at the distributed data node sites (e.g., the National Center for Atmospheric Research (NCAR) deep storage archive, and the British Atmospheric Data Centre

(BADC) fast access disk. Using popular climate analysis tools (e.g., CDAT, Ferret, NCL), the user regrids the data *where they physically reside* before transferring the reduced data subset to the PCMDI gateway, where further intercomparison diagnostics are performed on the disparate data sets. Once the diagnostics are complete, the desired products are returned to the user's platform.

FIG. 5. Tiered ESG-CET architecture showing the tri-level data services and one of the initial ESG Gateways specific to the CMIP5 (IPCC AR5) application. Initially, three ESG Gateways are planned: one at PCMDI focused on the CMIP5 (IPCC AR5) needs, one at ORNL to support the Computational Climate End Station project, and one at NCAR to serve the CCSM and PCM model communities.

**FIG. 1. The ESG-CET collaboration includes participation from Argonne National Laboratory (ANL), Los Alamos National Laboratory (LANL), Lawrence Berkeley National Laboratory (LBNL), Lawrence Livermore National Laboratory (LLNL), National Center for Atmospheric Research (NCAR), Oak Ridge National Laboratory (ORNL), Pacific Marine Environmental Laboratory (PMEL), and the University of Southern California (USC) Information Sciences Institute.**

**FIG. 2. Schematic of the first-generation ESG architecture showing the U.S. repositories. Climate model data is located on deep archives at the Lawrence Berkeley National Laboratory (LBNL) National Energy Research Scientific Computing Center (NERSC), the National Center for Atmospheric Research's (NCAR) Mass Storage System (MSS), and the Oak Ridge National Laboratory (ORNL) High Performance Storage System (HPSS), and on the Los Alamos National Laboratory (LANL) fast-access rotating disks. Also depicted is a provider publishing data by means of a web browser and a data user accessing published data via either a web browser or DataMover-Lite (DML).**
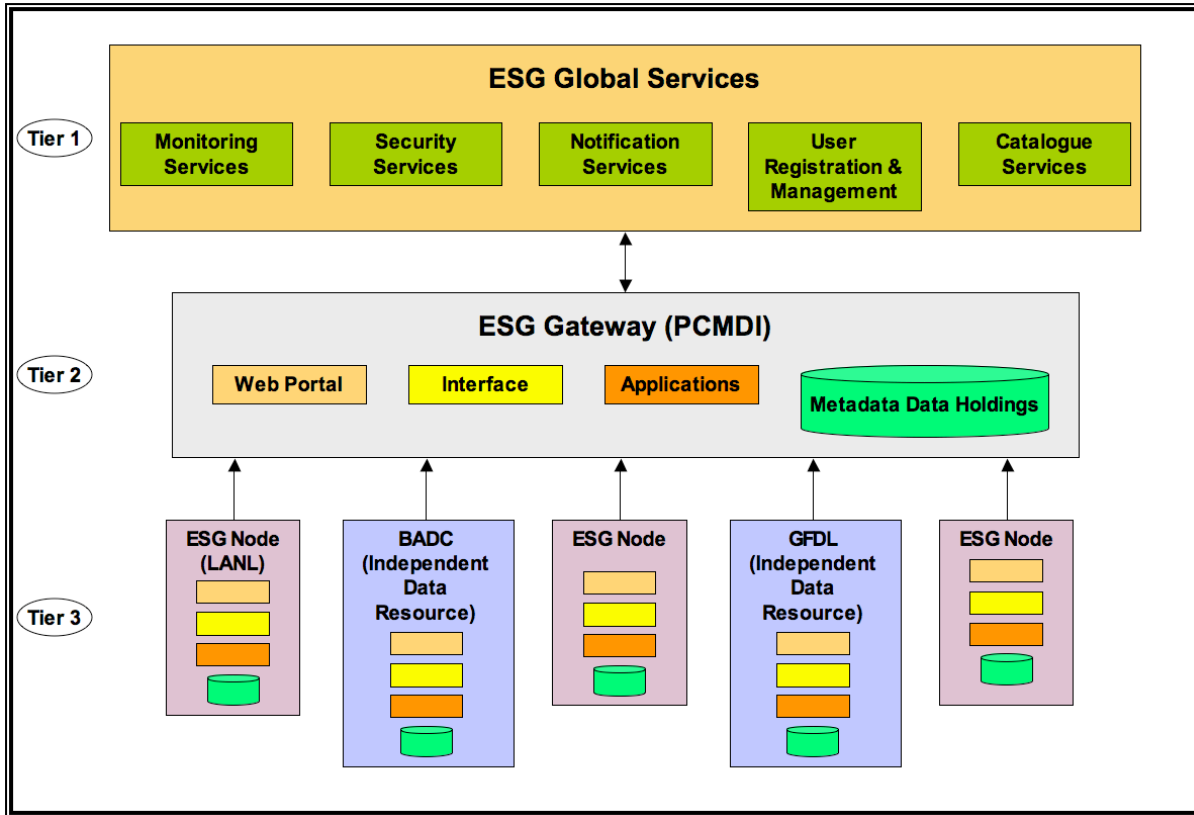
**FIG. 3. A high-level ESG-CET "roadmap" showing the planned evolution of the ESG system from terascale to petascale data management. Also shown are the scientific data management and analysis requirements in relationship to the ESG development timeframe. Note that a distributed testbed for CMIP5 (IPCC AR5) needs to be in place by early 2009.**

**FIG. 4. In this example a user searches the ESG-CET portal from a remote site. Required data are found at the distributed data node sites (e.g., the National Center for Atmospheric Research (NCAR) deep storage archive, and the British Atmospheric Data Centre (BADC) fast access disk. Using popular climate analysis tools (e.g., CDAT, Ferret, NCL), the user regrids the data *where they physically reside* before transferring the reduced data subset to the PCMDI gateway, where further intercomparison diagnostics are performed on the disparate data sets. Once the diagnostics are complete, the desired products are returned to the user's platform.**

**FIG. 5.** Tiered ESG-CET architecture showing the tri-level data services and one of the initial ESG Gateways specific to the CMIP5 (IPCC AR5) application. Initially, three ESG Gateways are planned: one at PCMDI focused on the CMIP5 (IPCC AR5) needs, one at ORNL to support the Computational Climate End Station project, and one at NCAR to serve the CCSM and PCM model communities.

**Websites**

| Name | URL |
|------|-----|
| Scientific Discovery through Advanced Computing (SciDAC)-2 | http://www.scidac.gov/ |
| Earth System Grid Center for Enabling Technologies (ESG-CET) | http://www.earthsystemgrid.org/ |
| Coupled Model Intercomparison Project 3 (CMIP3) multi-model database | http://www-pcmdi.llnl.gov/ipcc/about_ipcc.php |
| Community Climate System Model (CCSM) archive | http://www.earthsystemgrid.org/ |
| CCSM Biogeochemistry (BGC) Carbon-Land Model Intercomparison Project (C-LAMP) archive | http://esg2.ornl.gov/ |
| Program for Climate Model Diagnosis and Intercomparison (PCMDI) | http://www-pcmdi.llnl.gov |
| Climate Model Output Rewriter (CMOR) | http://www2-pcmdi.llnl.gov/software-portal/cmor/ |
| Open-source Project for a | http://www.opendap.org/ |

| | |
|---|---|
| Network Data Access Protocol (OPeNDAP) | |
| NetCDF | http:/www.unidata.ucar.edu/software/netcdf/ |
| CF | http://cf-pcmdi.llnl.gov/ |
| Climate Data Analysis Tools (CDAT) | http://www-pcmdi.gov/software-portal/cdat/ |
| North American Regional Climate Change Assessment Program (NARCCAP) | http://www.narccap.ucar.edu/ |
| Earth System Curator (ESC) | http://www.earthsystemcurator.org/ |
| MetaFor | http://www.cgam.nerc.ac.uk/pmwiki/PRISM/index.php/Main/METAFORPage |
| Open Geospatial Consortium (OGC) | http://www.opengeospatial.org/ |
| THREDDS | http://www.unidata.ucar.edu/projects/THREDDS/ |
| Extensible Mark-up Language (XML) | http://www.xml.org |
| THREDDS XML with Extensible Stylesheet Language Transformations (XSLT) | http://www.w3.org/TR/xslt/ |
| Hypertext Markup Language | http://www.w3.org/MarkUp/ |

| | |
|---|---|
| (HTML) | |
| Live Access Server (LAS) | http://ferret.pmel.noaa.gov/Ferret/LAS/home/ |
| Center for Enabling Distributed Petascale Science (CEDPS) | www.cedps.net |