

## Taking Part in ESG: Data Node Installation

Dean Williams and Karl E. Taylor (PCMDI)

Bryan Lawrence (BADC)

*Please note that the information in this document may become dated rather quickly as the ESG Federation is established. If the current date differs from the document version date, check the CMIP5 website to ensure that a later version has not been produced.*

### Hardware, storage, and network requirements:

To begin with, the hardware, storage, and network requirements for an ESG-CET Data Node are:

- 1) Hardware of your choice running the latest version of CentOS or Red Hat Enterprise Linux (RHEL). For example:
  - a. Dell PowerEdge 2850: 2 core x 3.2 GHz Intel, 4 GB RAM, RHEL 5.3, 26 TB disk
  - b. Dell R610 2-CPU E5540, 32 GB RAM, CentOS 5.3, 24 TB disk, direct attachment to 250 TB RAID archive
  - c. Other flavors of Linux/Unix will probably work, but will not be supported at this time. As time progresses, additional operating systems will likely be added to the support list.
- 2) Spinning disks (rather than tapes) are required for storage (to assure fast access), but the choice of specific hardware devices is up to the host.
- 3) The faster the network, the better. Some groups will have 2 x 10 Gbps network access while others will have less than 1 Gbps. Initially, PCMDI will have 2 x 10 Gbps network access and is currently working to connect to a 100 Gbps network.

### Software prerequisites:

Once your hardware (with operating system), storage and network are in place, make sure the following software is installed before initiating the Data Node installation script:

- 1) gmake (GNU make) – usually preinstalled on Redhat Linux
- 2) svn (subversion client) – usually preinstalled on Redhat Linux
- 3) JDK 6+ (Java 1.6+) – usually installed on Redhat Linux
- 4) An account on an ESG-CET MyProxy server, with Publisher role. This can only be created from the ESG-CET Gateway Portal at PCMDI. (You must request that PCMDI establish an account before you will be able to publish to the PCMDI Gateway)

- 5) X11 libraries and headers (xserver-xorg-dev, libx11-dev, libxt-dev, libxrender-dev). These are only used to build the publisher graphical user interface (GUI), which facilitates visual interactions and quick views of published data. In the absence of X11, the publisher scripts will still build, but the GUI will not be available.
- 6) Development library headers are also needed, namely zlib, termcap and readline. They are used when building the PostgreSQL database. It is recommended that the operating system be installed with all the development headers for the distribution.
- 7) A host with the following ports accessible:
  - HTTP: port 80 (for incoming traffic from the Gateway)
  - SSL: port 443
  - GridFTP: ports 2811, 50000 - 51000
  - There may be another port open between the Gateway and Data Nodes.

### **Installation script:**

The installation script contains SVN software download commands and bash scripts to obtain, build and configure the ESG-CET Data Node for you. The installer is queried to answer a few questions about the Data Node configuration for their particular environment (questions like, “Do you want to install an analysis and visualization product server?” or “Do you want a GridFTP server?”). (Default choices are provided to the installer and in most cases recommended.) The time it takes to install the ESG-CET Data Node is approximately 1 hour. Experienced installers have installed the Data Node software stack in 30 minutes. The ESG-CET Data Node script installation can be found at the following URL:

- [http://rainbow.llnl.gov/dist/datanode/ESG\\_DataNode\\_install.tgz](http://rainbow.llnl.gov/dist/datanode/ESG_DataNode_install.tgz)

### **Virtual Machine (VM) installation:**

The ESG-CET software stack may be installed on a traditional 'bare metal' linux distribution, as described above, or distributed via a virtual machine (VM). A VM is a software platform that provides the ability to run a complete and sovereign operating system (a guest OS) as an application inside another operating system (the host OS). The guest OS executes software applications identically to a physical machine. The ESG-CET VMs contain the fully installed CentOS operating system (with requisite libraries et. al. installed) and the ESG-CET Data Node software stack. (This eliminates the need to check for software prerequisites.) VMs take advantage of the benefits of hardware virtualization – specifically, better security, hardware insulation, portability, ease of backup, and protection against potential software conflicts. For ready operation with minimal configuration, follow the links to the CentOS ESG-CET Data Node VM installation files:

- [http://rainbow.llnl.gov/dist/datanode/ESG\\_CentOS\\_Linux\\_2.6.x\\_kernel.v4.dn.tgz](http://rainbow.llnl.gov/dist/datanode/ESG_CentOS_Linux_2.6.x_kernel.v4.dn.tgz)
- [http://rainbow.llnl.gov/dist/datanode/ESG\\_CentOS.v4.dn.md5](http://rainbow.llnl.gov/dist/datanode/ESG_CentOS.v4.dn.md5)

*[Please note that the first link, a VMWare CentOS virtual machine and is 1.6GB!]*

### **ESG-CET Software components to be installed:**

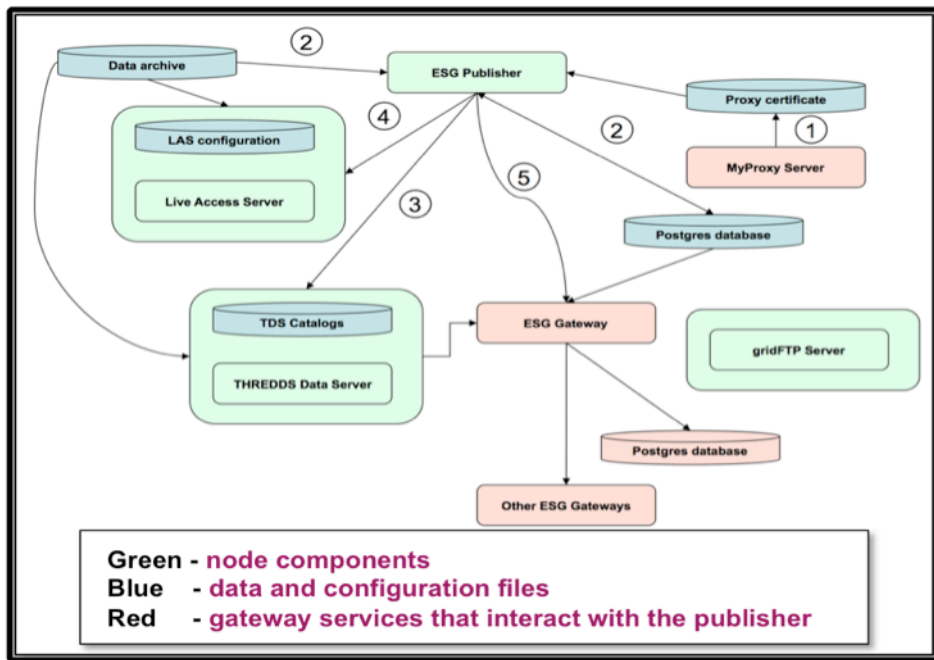
The components of the ESG-CET Data Node that will be installed are:

- 1) PostgreSQL (relational database) – contains the results of scanning and harvesting metadata from the input (netCDF CF) files. The database also contains other information required for system operation.
- 2) CDAT/CDMS – CDAT is based on the Python object-oriented language. CDMS is a sub-module of CDAT used by ESG-CET to scan data.
- 3) Python ESG-CET package – Publication modules, scripts, and GUI.
- 4) Apache Tomcat (Java servlet container) – web application server that runs the THREDDS application.
- 5) THREDDS Data Server (TDS) – a servlet for cataloging ESG-CET Data Node files and data resources used for reading catalogs and publishing to the PCMDI Gateway.
- 6) Live Access Server (LAS) – servlet for generating visualization products [*Note: LAS product service is optional and not yet ready for distribution.*]
- 7) Globus Toolkit / MyProxy – In our distributed system it is important to provide data security and integrity as well as user authentication. The MyProxy client is used for Single Sign-On to the PCMDI Gateway.
- 8) GridFTP server – GRID-enabled FTP server.
- 9) Test publication – Publish a test data file.

### Operational workflow:

As illustrated in the figure below, data is published using the Data Node software through a series of operations: The publisher

1. Obtains a proxy certificate.
2. Scans a set of files: a) associates files to datasets, and b) caches file metadata in the node database.
3. Generates THREDDS catalogs and reinitializes the THREDDS Data Server (TDS).
4. Optionally (and only after the LAS product service becomes available) reinitializes the Live Access Server (LAS) based on the TDS catalog.
5. Publishes to the PCMDI Gateway: a) contacts Gateway via web service, and b) allows the Gateway to harvest the newly created TDS catalogs.



### Support:

For ESG-CET software and installation help and support, please send your e-mail questions to:

- [esg-node-user@lists.llnl.gov](mailto:esg-node-user@lists.llnl.gov)