

Climate data storage in e-INIS

*G. Quigley, B. Coghlan, J. Ryan (TCD).
A McKinstry (ICHEC). K. Rochford (DIAS).
5th EGEE User Forum, Uppsala, Sweden.*

- **CMIP5**
 - Overall aims
 - Met Éireann and ICHEC involvement
 - Projected storage requirements
- **Grid Ireland and e-INIS**
 - National e-Infrastructure
- **National Datastore**
 - Architecture
 - Application to CMIP5
- **Progress to date**

- **Coupled Model Intercomparison Project Phase 5**
- **Standard experimental protocol for studying the output of coupled ocean-atmosphere general circulation models.**
- **Community-based infrastructure in support of climate model diagnosis, validation, intercomparison, documentation and data access.**
 - address outstanding scientific questions that arose as part of the IPCC AR4 process
 - improve understanding of climate
 - provide estimates of future climate change that will be useful to those considering its possible consequences

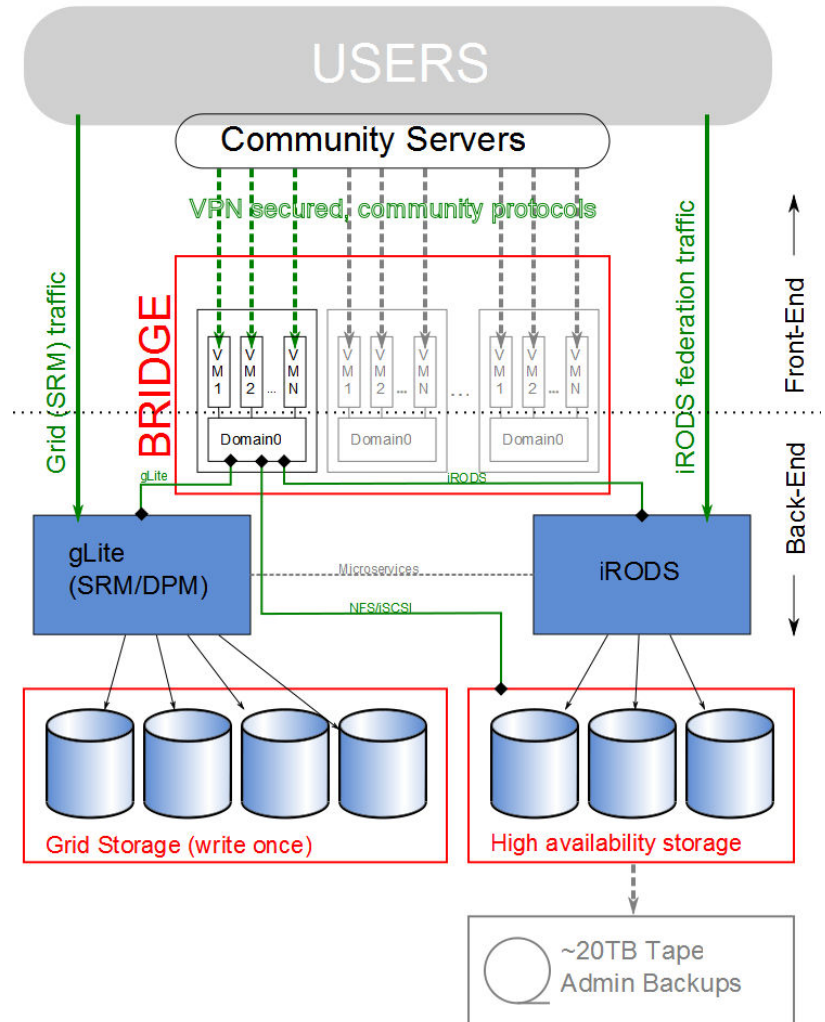
- **Framework for coordinated climate change experiments for the next 5 years**
 - simulations for assessment in the AR5 and that extend beyond
 - CMIP5 is not meant to be comprehensive
 - cannot possibly include all the different model intercomparison activities that might be of value
 - expected that various groups will develop additional experiments
- **CMIP5 promotes a standard set of model simulations**
 - To evaluate how realistic the models are in simulating the recent past
 - provide projections of future climate change on two time scales, near term (out to ~2035) and long term (out to 2100+)
 - understand some of the factors responsible for differences in model projections (feedbacks, clouds, the carbon cycle, etc.)

- **EC-Earth model will produce ~200TB data at ICHEC**
 - Approx. 100,000 netCDF files
 - Runs from late-2009 to end 2010
- **to be used by consortium members and the wider climate community**
- **Read-write access to this data is available to EC-Earth scientists from ICHEC**
 - 1 Gbps ethernet, with upgrade to 10 Gbps light-path in 2010 planned.
- **Public access to be offered using ESG data node**
 - OpenDAP/HTTP allowing access to specific parts of files and metadata
 - Central gateways manage data nodes

- **Irish National e-Infrastructure**
- **Aims to provide a national e-infrastructure to academic researchers**
 - access to ICHEC capacity and capability computing facilities
 - specialist expert user support and training
 - secure HEAnet network and Grid-Ireland grid services
- **Also includes pilot national datastore**
 - For any discipline (not just science)
 - DHO archive of Irish Language <http://dho.ie/doegen/>
 - HELIO solar physics project (poster at forum)
 - ATLAS
 - CMIP5...
 - Initially federated across 3 sites (TCD, DIAS, UCC)
 - Currently <1PB but expansion soon!

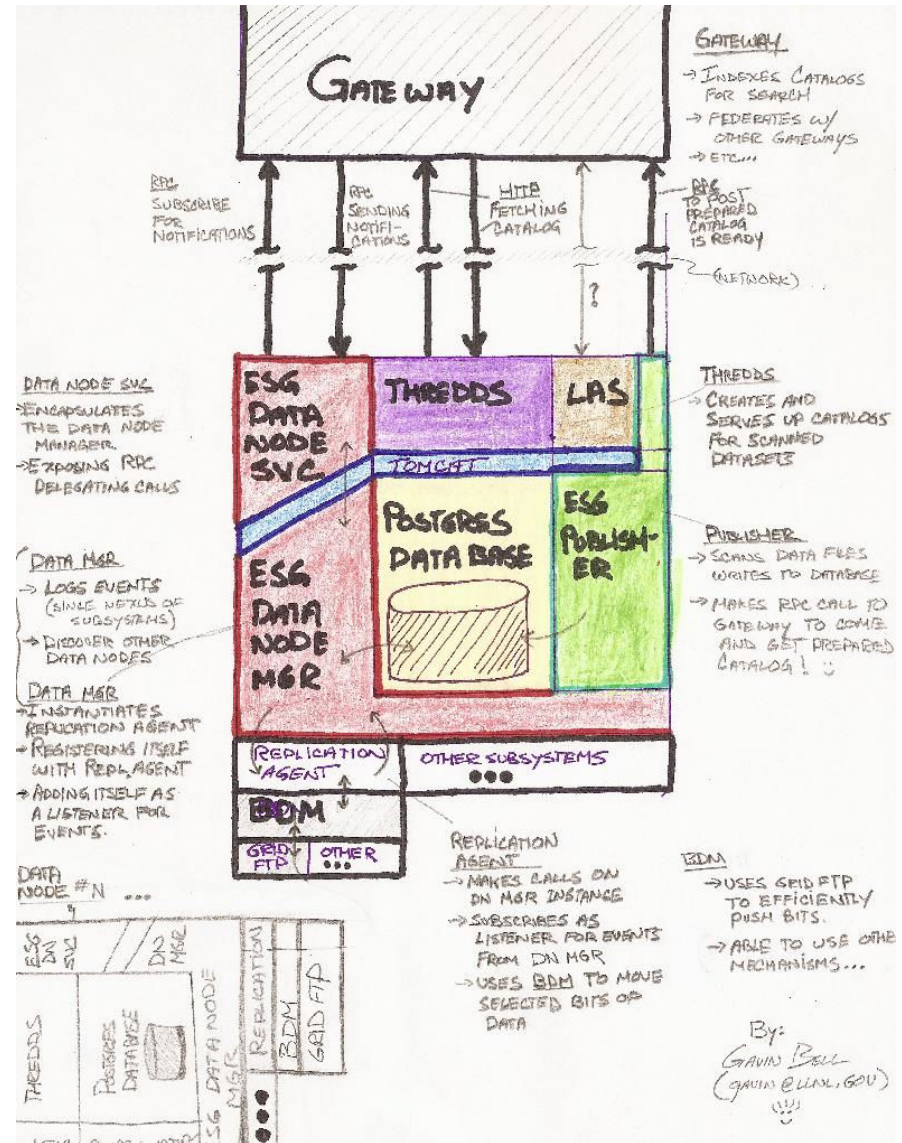
- **Sites act as federation, significant autonomy**
- **Some centralisation may be kept, e.g. for catalogs**
- **Two main back-end technologies**
 - EGEE gLite middleware (DPM, LFC) on Grid-Ireland
 - iRODS (Integrated Rule-Oriented Data System)
- **Very varied user requirements**
 - Can target data to most cost-effective storage that meets requirements
 - Need to ‘bridge’ between back-end protocols and user protocols
 - Insulate users from each other’s activities
- **Grid-Ireland OpsCentre (at TCD)**
 - Runs largest slice of national datastore
 - 600TB disk-based storage

- Bridge layer created
- Blades running virtual machines
- Current iteration uses Open Nebula cloud infrastructure
- VMs presented subset of storage
- Re-export using community protocol
- Community host own front-end
- Write-access always secure (policy)
- Read access can be open if desired



- **Data ingest and some access via gLite**
- **Read-only access required for publishing**
- **ESG data node software used to publish appears to expect local filesystem**
- **FUSE filesystem written to mount Grid storage**
 - Storage is at local site
 - GFAL API calls used
 - Only functions for listing, reading files implemented in FUSE
 - This works but is slow
 - Future iterations may prove better?
 - May need to implement more at RFIO level
 - Hides this detail from the ESG software!

- Suite of software components
- Locally creates catalog
- Publishes to gateway
- Uses own gridftp server for replication
- Separate (non-EGEE) grid stack



- **Data being generated at ICHEC and uploaded**
 - Approximately 37TB uploaded so far
- **gLite middleware working**
- **OpenNebula private cloud (bridge)**
- **Working prototype FUSE filesystem**
- **ESG data node install attempted**
- **Still debugging the data node install**
 - Problems caused by firewall/proxy servers

- **The Irish involvement in CMIP5 producing 200TB data**
- **Needs to be shared with a (non-Grid) distributed community**
- **Storage currently available uses gLite**
- **'Bridge' used to publish this data using community software**
- **Current problems are with inter-site networking**

- Grid-Ireland:** <http://www.grid.ie>
- ICHEC:** <http://www.ichec.ie>
- e-INIS:** <http://www.e-inis.ie/>
- CMIP5:** <http://cmip-pcmdi.llnl.gov/cmip5>
http://www.ichec.ie/research/met_eireann#cmip5

National datastore paper:

“The Back-end of a 2-Layer Model for a Federated National Datastore for Academic Research VOs that Integrates EGEE Data Management”, to be published Journal of Grid Computing, EGEE special edition, 2010. [DOI: 10.1007/s10723-010-9149-9](https://doi.org/10.1007/s10723-010-9149-9)