



FermiGrid (The Fermilab Campus Grid)

Keith Chadwick
Fermilab
chadwick@fnal.gov

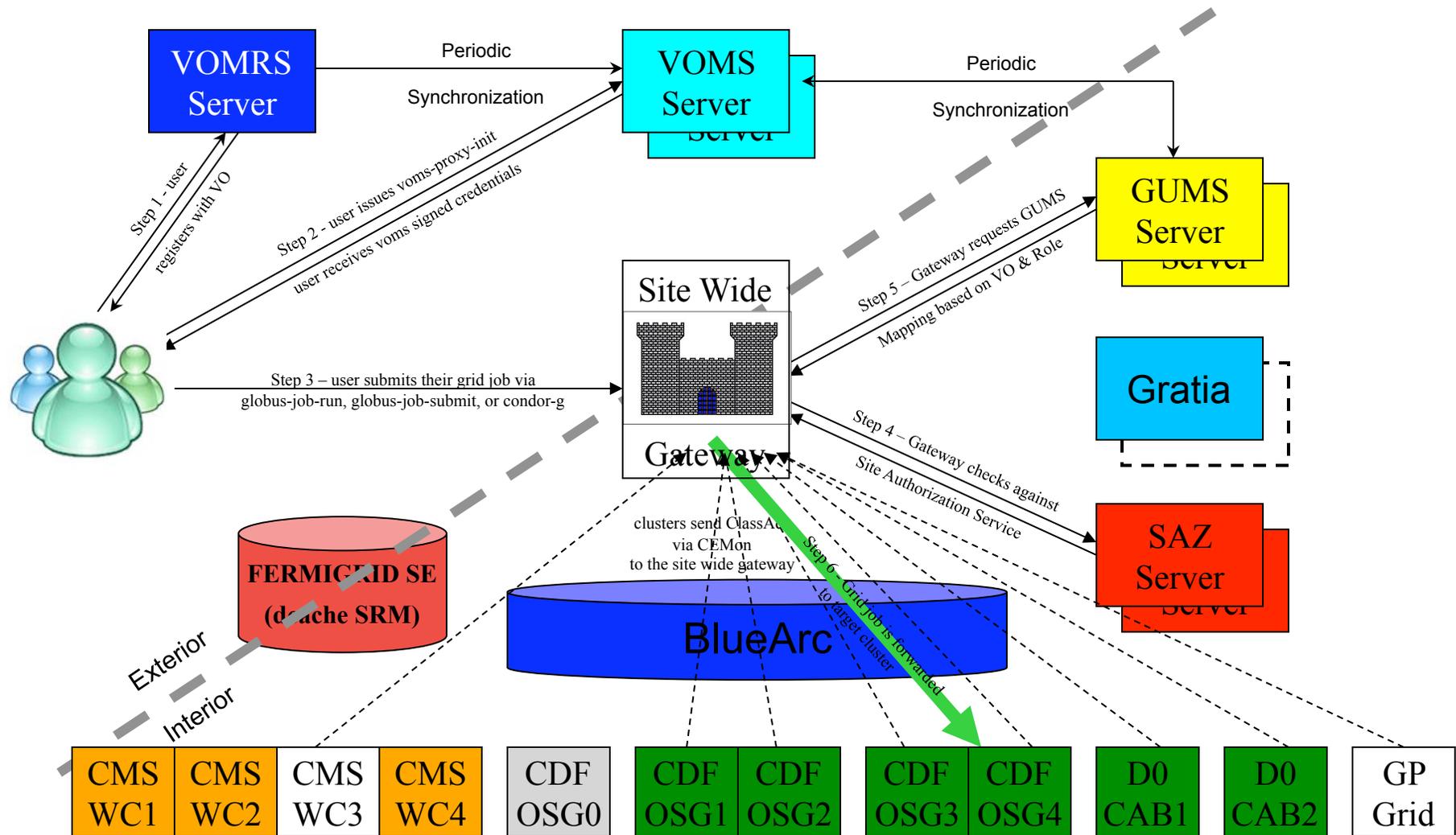


What is FermiGrid?

FermiGrid is the Fermilab Campus Grid.

- We operate the central Fermilab Grid services cyberinfrastructure.
 - VOMRS, VOMS, GUMS, SAZ, Squid, Site Gatekeeper, etc.;
 - Services offered to CDF, D0, CMS and other Fermilab consumers;
 - The services are deployed in a Highly Available (HA) infrastructure with automatic failover (FermiGrid-HA).
- We operate multiple Grid resources for a variety of client communities.
 - CDF, CMS & D0 experiments;
 - Other smaller Fermilab experiments.
- We coordinate and interoperate with other campus Grids (Purdue), regional Grids (NYSGrid, SuraGrid), and national cyberinfrastructures (Open Science Grid, TeraGrid) [FermiGrid is an OSG site].
- <http://fermigrid.fnal.gov>

Current Architecture





Communities and Job Slots

Community	# Clusters	# Gatekeeper	# Slots	Batch	VO Location
_____	_____	_____	_____	_____	_____
CDF Experiment	3	5*	5,315	Condor	Fermilab
CMS Experiment	1	4*	5,144	Condor	CERN
D0 Experiment	2	2	5,597	PBS	Fermilab
“Other” Fermilab	1	3	965	Condor	Fermilab
OSG VO's	--	--	--	--	National
TeraGrid	--	--	--	--	National
SuraGrid	--	--	--	--	Regional
_____	_____	_____	_____	_____	_____
Total	7	14	17,021		

Table Data Source = <http://fermigrid.fnal.gov/fermigrid-metrics.html>

- Four of the five CDF gatekeepers are open for opportunistic use.
- Only one of the four CMS gatekeepers is open for opportunistic use.

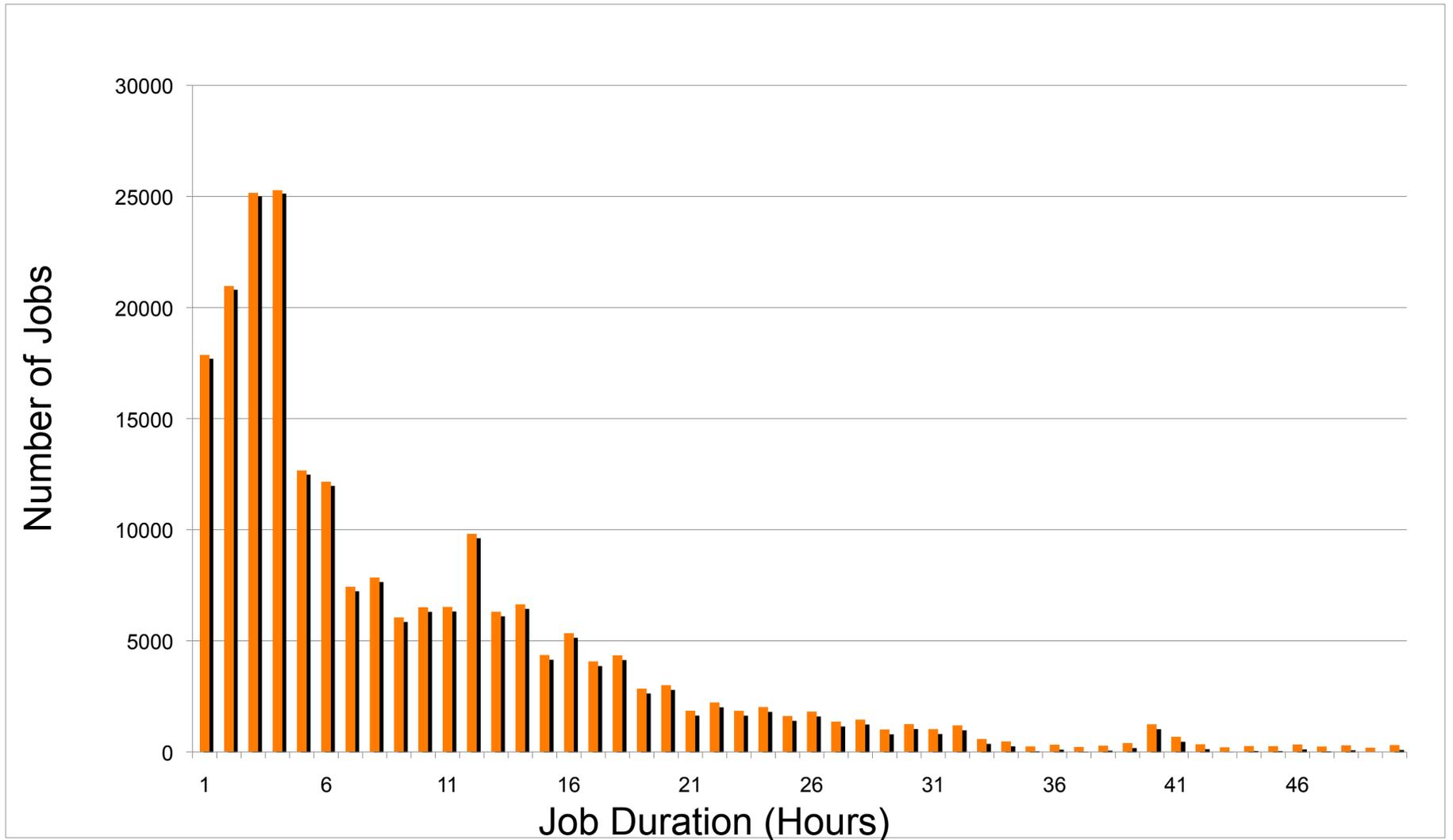


Scheduling within FermiGrid

1. The "owner" of the resource has first call on the resource.
2. Other Fermilab experiments have second call on the (remaining) resources.
3. Outside collaborators (OSG VO's) have third call on any remaining resources.
4. If the "owner" has a job waiting for the job slot, then "opportunistic" jobs are given 24 hours to finish:
 - Gratia accounting data shows that ~95% of *all* jobs finish in less than 24 hours of elapsed time;
 - Gratia accounting data shows that >99% of *opportunistic* jobs finish in less than 24 hours of elapsed time.



Job Duration Distribution





Storage Resources

FermiGrid provides the three canonical storage areas: \$APP, \$DATA, \$WN_TMP:

- \$APP is an area for software installation shared across all worker nodes. It is read/write from the head node via the jobmanager-fork queue, and read-only from worker nodes. The default quota is set to 30 GB for all the users in a VO. VO's that have more than one user may request that only a single privileged user be allowed quota to write files. By default, the directory is mode 1777. The APP directory seen on fermigrd1 is shared with all the member clusters except for the USCMS Tier 1 cluster. This space is backed up to tape on a nightly basis.
- \$DATA is a NFS-mounted area which is accessible read-write from the head node and all the workers. The default quota is set to 400 GB for all the users in a VO. By default the directory is mode 1777. The DATA directory seen on fermigrd1 is shared with all the member clusters except for the USCMS Tier 1 cluster. **This space is not backed up to tape.**
- \$WN_TMP is a read/write area for temporary files local to each worker node.

Dcache storage element FNAL_FERMIGRID_SE:

- The \$SITE_READ and \$SITE_WRITE environment variables are URL's that point to a local storage element accessible by SRM and dCache for reading, and by SRM for writing. This is declared to the OSG as FNAL_FERMIGRID_SE. This is volatile storage with a total available space of 5TB. Files are deleted, if necessary, on a least recently used basis and there are no quotas. This space is not backed up to tape.

Fermilab may delete or purge files in any of the above storage areas as necessary to assure operations.

- Under normal running conditions, we will attempt to contact the affected VOs and/or users through the OSG Operations Center and attempt to provide up to 4 hours of warning before deleting or purging files.
- During some running conditions we may be forced to delete or purge files without warning.



FermiGrid-HA - Highly Available Grid Services (e-infrastructure)

The majority of the services listed in the FermiGrid service catalog are deployed in high availability (HA) configuration that is collectively known as "FermiGrid-HA".

FermiGrid-HA utilizes three key technologies:

- Linux Virtual Server (LVS).
- Scientific Linux (Fermi) 5.3 + Xen Hypervisor.
- MySQL Circular Replication.



Physical Hardware, Virtual Systems and Services

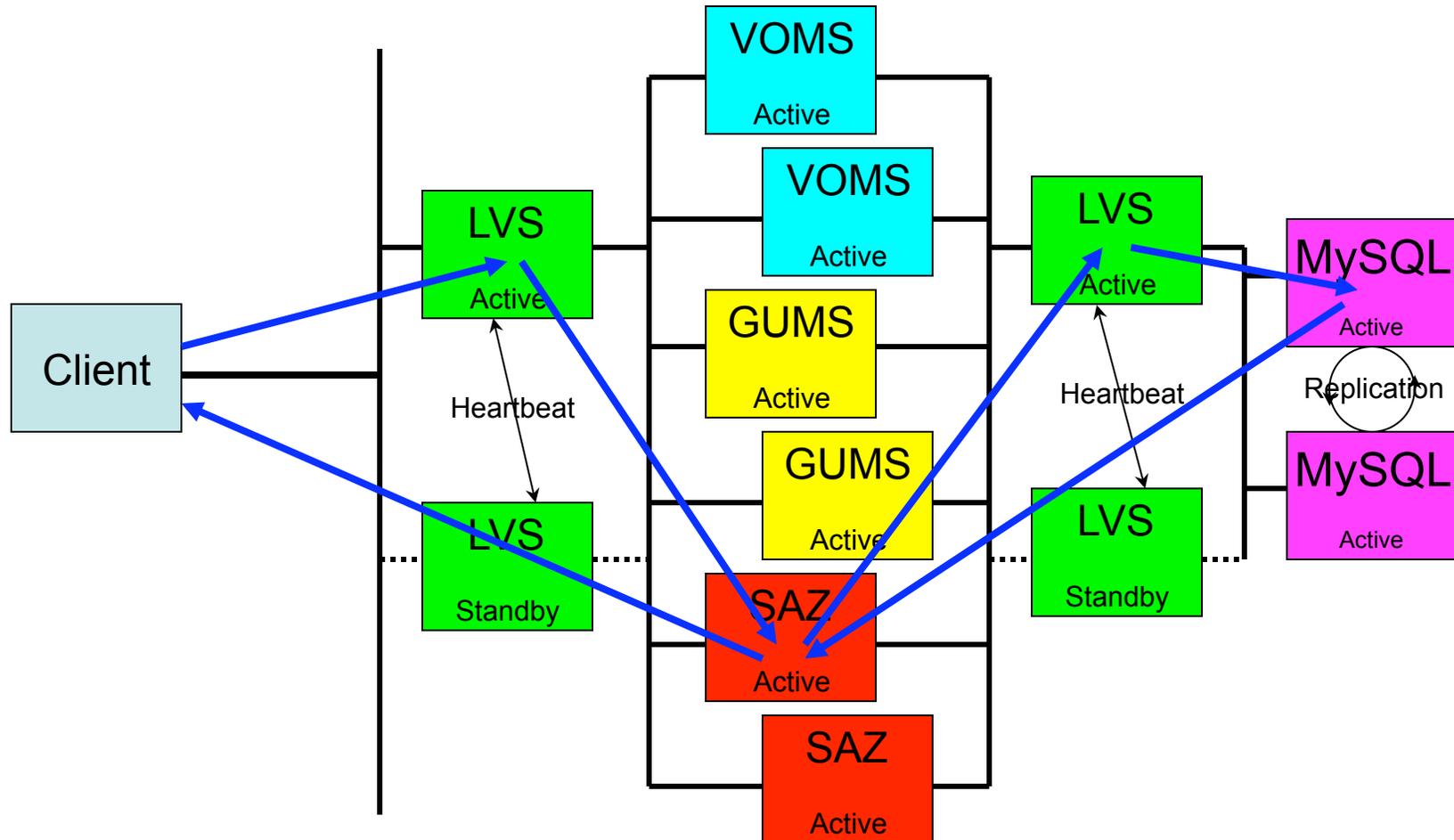
	Physical Systems	Virtual Systems	Virtualization Technology	Service Count
FermiGrid-HA Services	6	34	Xen	17
CDF, D0, GP Gatekeepers	9	28	Xen	9+6
Fermi & OSG Gratia	4	10	Xen	12
OSG ReSS	2	8	Xen	2
Integration Test Bed (ITB)	2+8	14+32	Xen	14
Grid "Access" Services	2	4	Xen	4
"FermiCloud"	8 (+16)	64 (+128)	Xen	--
"Fgtest" Systems	7	51	Xen	varies
"Cdf Sleeper Pool"	3	9	Xen	1+1
"GridWorks"	11	~20	Kvm	1



HA Services Deployment

FermiGrid employs several strategies to deploy HA services:

- Trivial monitoring or information services (examples: Ganglia and Zabbix) are deployed on two independent virtual machines.
- Services that natively support HA operation (examples: OSG ReSS, Condor Information Gatherer, FermiGrid internal ReSS deployment) are deployed in the standard service HA configuration on two independent virtual machines.
- Services that maintain intermediate routing information (example: Linux Virtual Server) are deployed in an active/standby configuration on two independent virtual machines. A periodic heartbeat process is used to perform any necessary service failover.
- Services that are pure request/response services and do not maintain intermediate context (examples: GUMS and SAZ) are deployed using a Linux Virtual Server (LVS) front end to active/active servers on two independent virtual machines.
- Services that support active-active database functions (example: circularly replicating MySQL servers) are deployed on two independent virtual machines.





Virtualized Non-HA Services

The following services are virtualized, but not (yet) currently implemented as HA services:

- Globus gatekeeper services (such as the CDF and D0 experiment globus gatekeeper services) are deployed in segmented "pools":
 - Loss of any single pool will reduce the available resources by approximately 50%;
 - Expect to segment the GP Grid cluster in late FY10;
 - We need a secure block level replication solution to allow us to implement this in an active/standby HA configuration;
 - DRBD may be the answer, and we have an active project to see if we can successfully incorporate the DRBD Kernel modifications into the Xen Kernel.
- MyProxy:
 - See comments about DRBD under Globus gatekeeper above.
- Fermi & OSG Gratia Accounting service [Gratia]:
 - Not currently implemented as an HA service;
 - If the service fails, then the service will not be available until appropriate manual intervention is performed to restart the service;
 - Equipment is has just been delivered to "HA" the Gratia services.



Measured Service Availability

FermiGrid actively measures the service availability of the services in the FermiGrid service catalog:

- <http://fermigrid.fnal.gov/fermigrid-metrics.html>
- <http://fermigrid.fnal.gov/monitor/fermigrid-metrics-report.html>

The above URLs are updated on an hourly basis.

The goal for FermiGrid-HA is > 99.999% service availability.

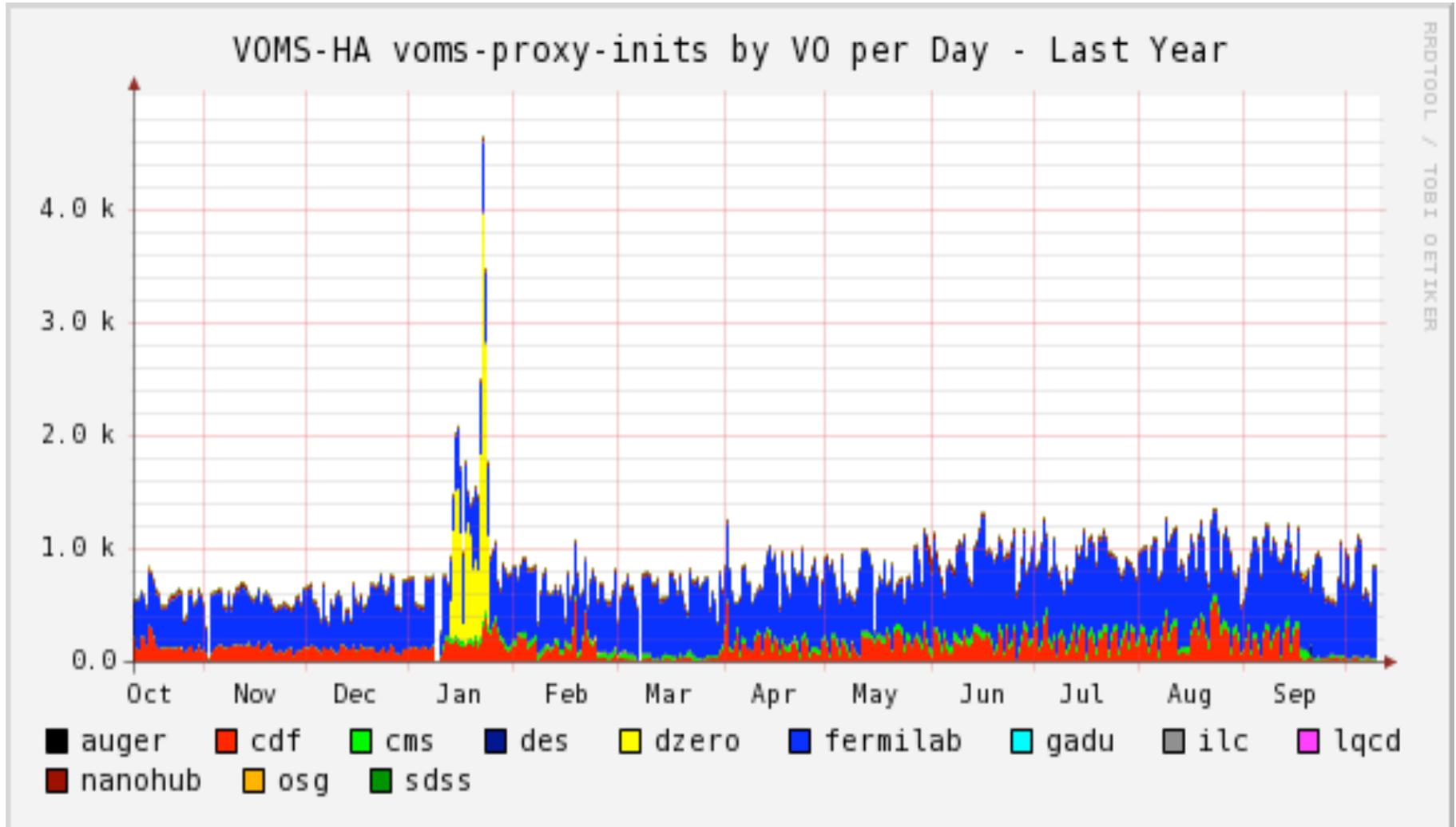
- Not including Building or Network failures.

For the period 01-Dec-2007 through 30-Jun-2008, we achieved a service availability of 99.9969% (10 minutes of downtime during 7 months of operation).

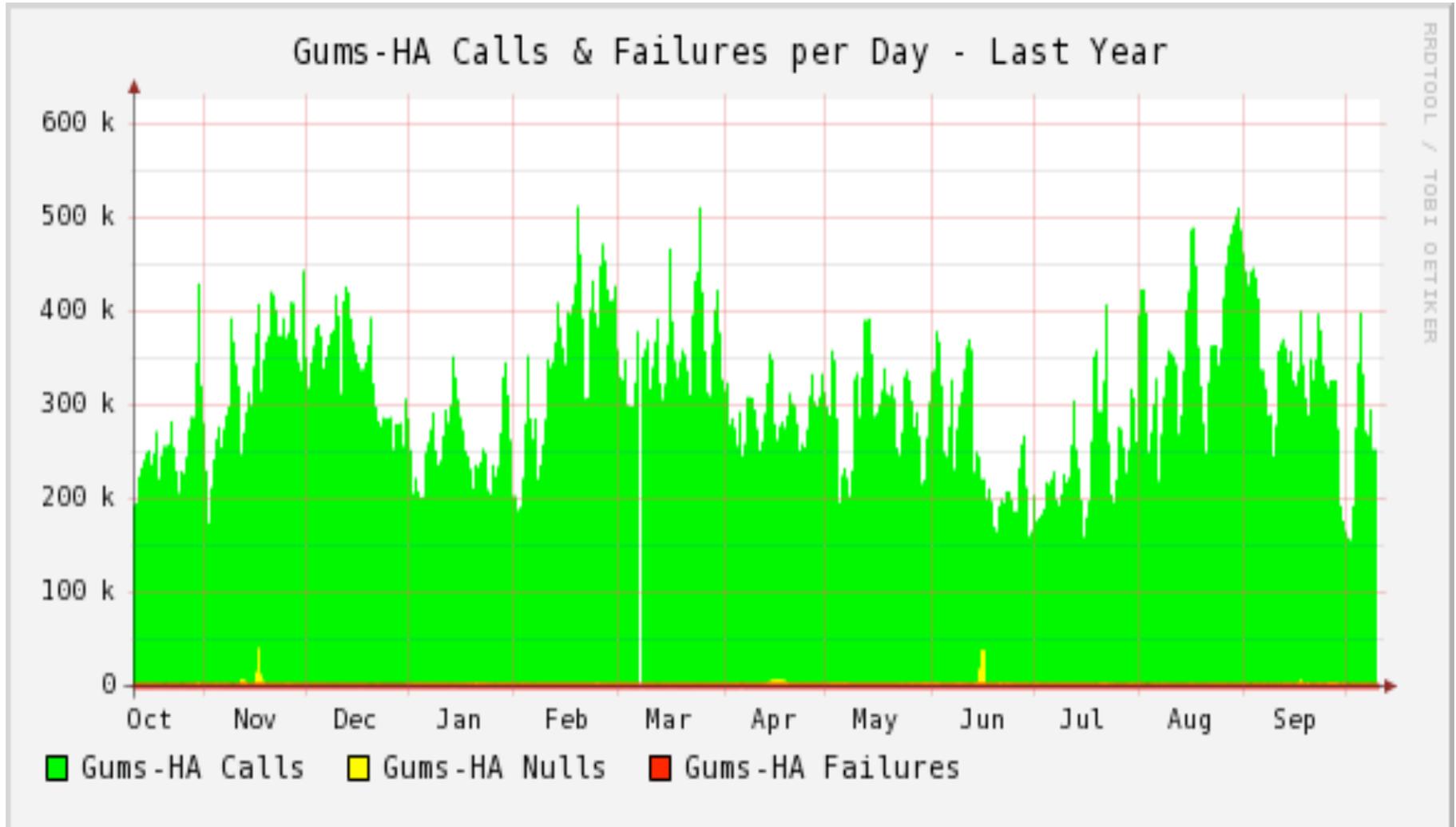
For the current year, we have only achieved a collective core service availability of 99.963% - Chiefly due to several user based denial of service attacks:

- Lets see what happens when a user attempts to download a 1.2 Gbyte file through the squid servers on 1,700 systems simultaneously...

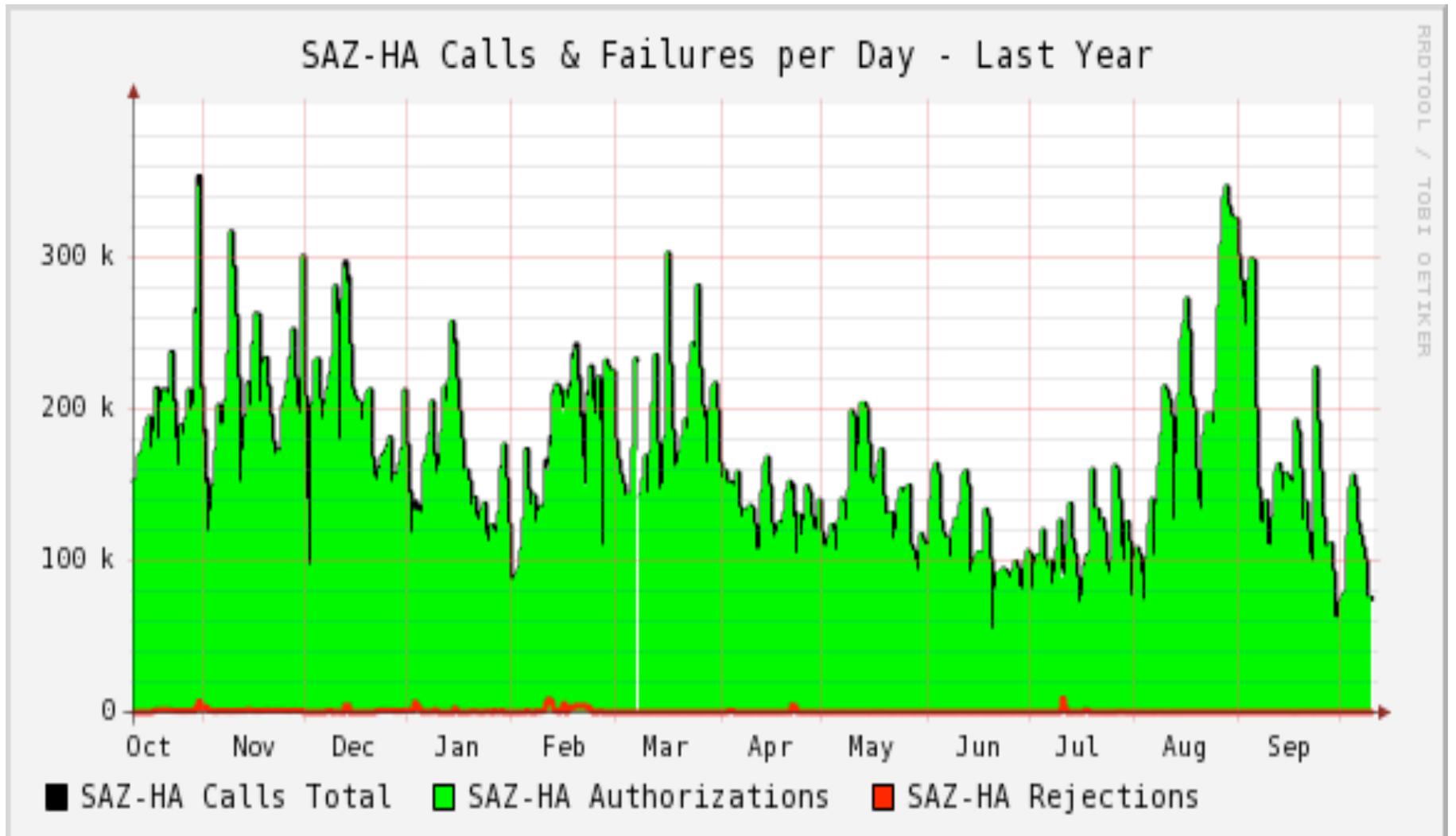
voms-proxy-init's by VO/day



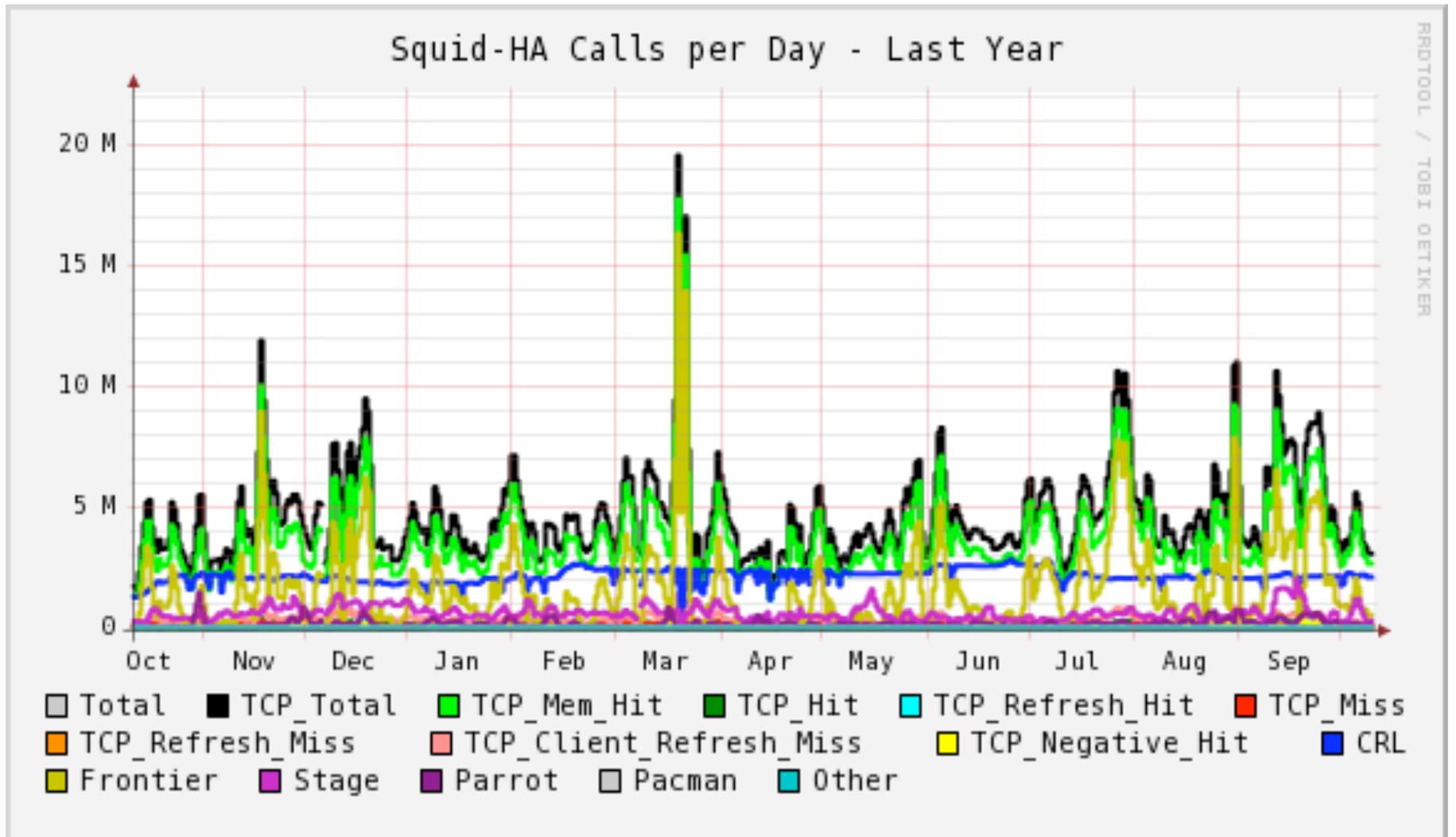
GUMS mappings/day



SAZ Decisions/day

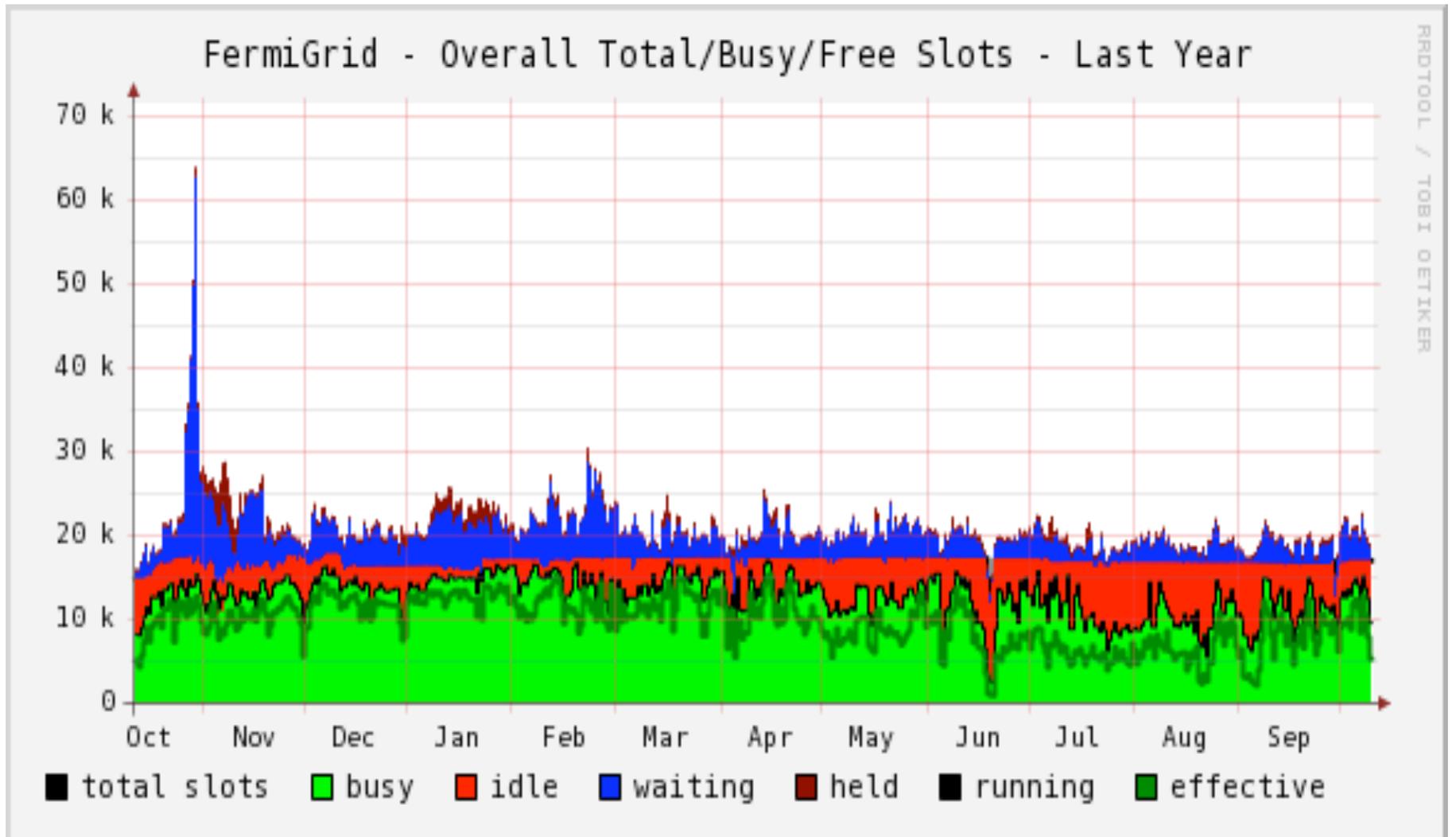


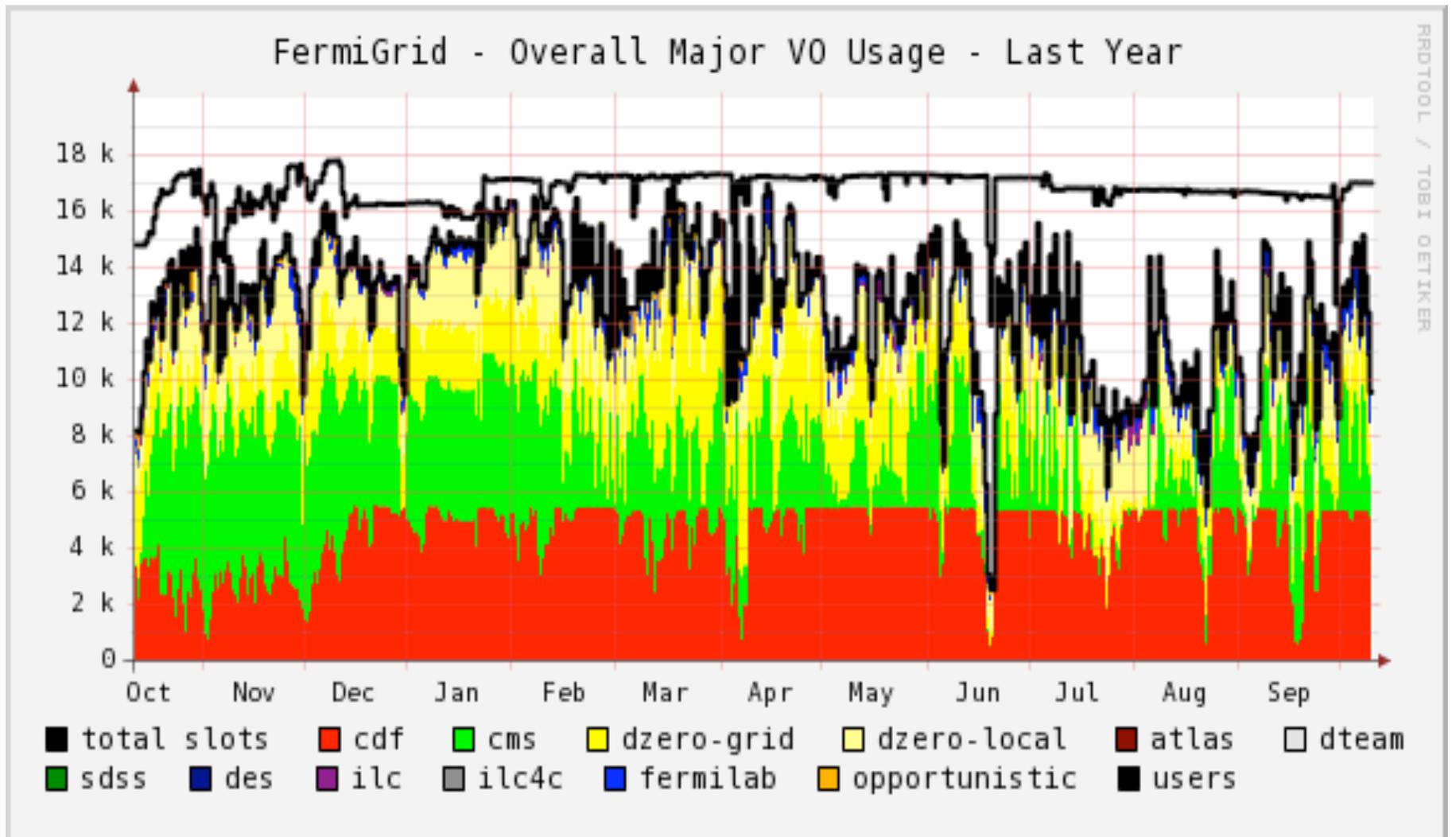
Squid cache accesses/day

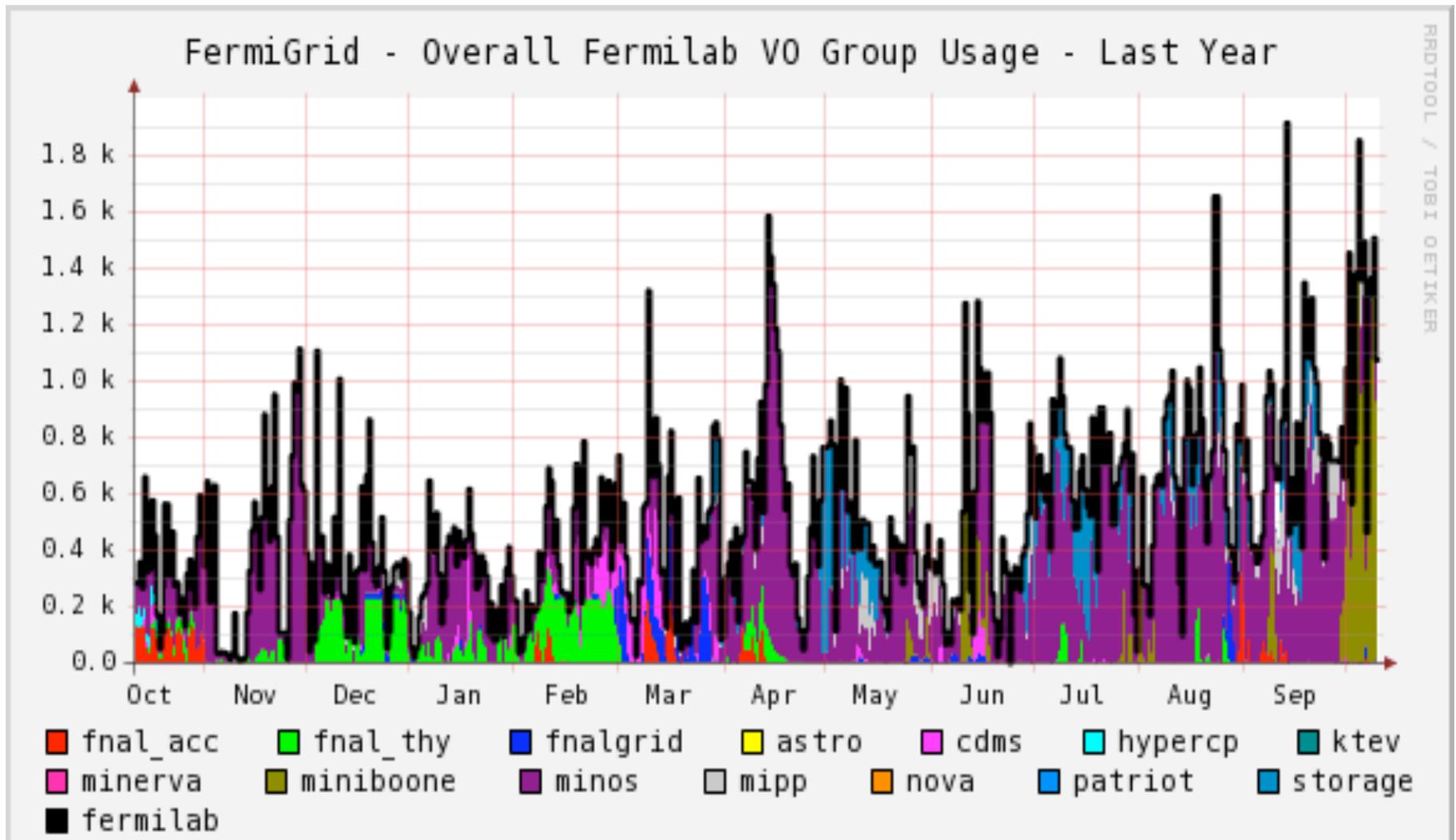


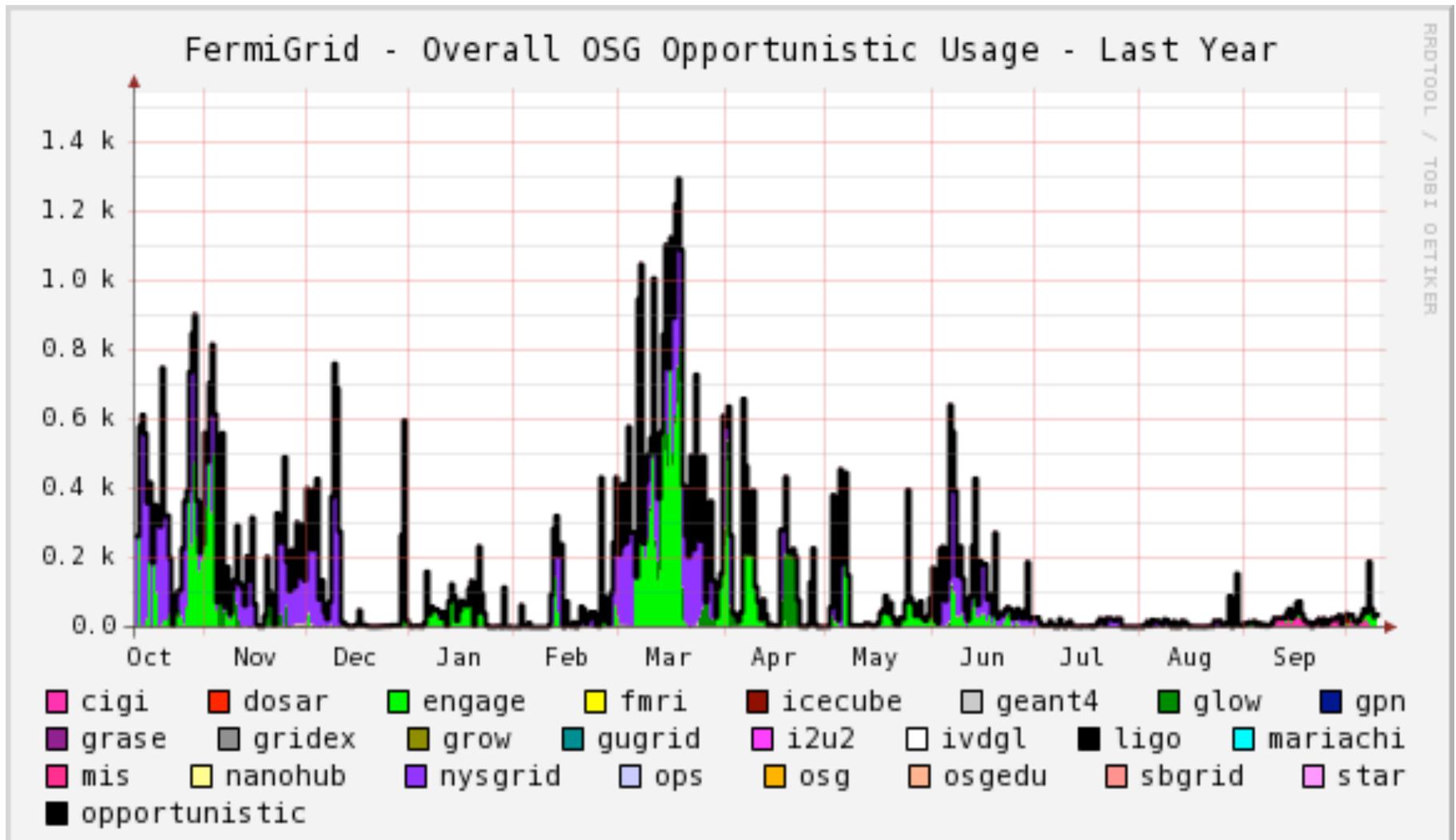


Jobs on FermiGrid - Last Year











FermiGrid - Summary

Provides the centralized Grid cyberinfrastructure for the Fermilab Campus Grid, and selected OSG services (ReSS, Gratia, VOMRS/VOMS).

Deployed cyberinfrastructure based on the OSG / VDT software stack.

Actively processing jobs from "on-site" Virtual Organizations (VO's) on their "dedicated" resources.

Actively processing opportunistic jobs from other "on-site" VO's.

Actively processing opportunistic jobs from VO's registered with OSG.

Operating a "transparent" OSG style GT2 and GT4 Gateway to selected TeraGrid resources.

Making plans to deploy a local Cloud Computing infrastructure "FermiCloud".

Any questions?