

# Intensity Frontier Computing Model

July 26, 2010

September 20, 2010

Lee Lueking for REX/IFront

## Background

For the Intensity Frontier (IF), computing resources have been provided piecemeal over the last decade. The experiments either had dedicated resources of their own or made use of FNALU nodes. MiniBooNE installed and administrated their own user cluster and storage servers, and MINOS utilized a CD supported cluster that provided user login and Condor batch processing. One group was able to leverage hardware and manpower from other projects such as ILC. In recent years some of the concerns about central storage mounts on FNALU have eased and specific nodes were associated with particular experimental efforts. Some experiments made effective use of FermiGrid resources, and have even used opportunistic computing cycles on Run II and CMS resources. Much of the simulation for MINOS was produced at sites beyond Fermilab, but using non-GRID approaches.

## Emerging Model

As the number of experiments increased a more coherent approach was warranted. In the last year and a half CD has helped experiments commission 136 CPU cores (17 dual quad core machines) in four experiment clusters for interactive login and Condor batch. Job submission scripts have been provided, based on MINOS's "minos\_jobsub", which enable users to submit batch jobs to experiment specific Condor pools on each system. These scripts also provide job submission to FermiGrid through glideinWMS services specific to each of the four experiments.

A centrally served disk storage model has been very successful and is provided via high performance NFS served mounts [1 BlueArc]. Each experiment has a storage allocation on this system appropriate to their requirements. At present there are over 200TB of disk served by this system to the IF experiments. Contention for this storage resource as hundreds, or thousands, of clients open and read/write files simultaneously is a potential bottleneck. This issue has been circumvented by a simple traffic control mechanism called "cpn" [2 Art's cpn utility] which manages the number connections. Critical data is archived to ENSTORE tape through either the direct encp or via dCache. Additional storage and caching solutions needed to meet the growing demand are outlined in a subsequent section.

The next step to higher scalability and maintainability is to separate the interactive login nodes and batch nodes. The details of the architecture for this system, called General Physics Computing

Facility or GPCF are described in [3 gpcf]. The interactive login component comprises multiple virtual machines on each physical box and shared “scratch” disk via NFS. The batch resources are non-virtualized boxes running a Condor pool shared among all the experiments. The general pool will be configured to provide dedicated access to specific nodes for each experiment, as well as general access when resources are available.

These two parts of the system, interactive login and local batch, provide the development, debugging and testing components of the system. The IF experiments have agreed to not include desktop clusters in their requirements, unlike Run II where such systems were administered for the experiments by CD. The tasks normally carried out on desktops are being done on the interactive login machines. Since these machines are uniformly procured and maintained, as well as centrally located in CD computer rooms their administrative load is manageable by FEF. It is assumed that users will bring desktop and laptop hardware, but the CD responsibility is limited to consulting support in OS installation, Kerberos, and security features.

The local batch is provided by an IF Condor pool to which users in each experiment submit jobs. To this end, each node in this pool has accounts for the union of all users, and mount points for the union of all mounts needed by all experiments. Job submission will be via a general tool “if\_jobsub” similar to the familiar “minos\_jobsub” mentioned above. The tool will determine the experiment information based on the node from which it is submitted, and the job will set the appropriate group when it starts on the general cluster. Ensuring that each group has both dedicated and shared resources on this Condor pool are managed through Condors Ranking provisions. User priorities within each group can be managed by setting priorities as determined by the experiment. Condor provides a tool called “condor\_ssh\_to\_job” that provides interactive debugging, ps, top, gdb, strace, lsof and even allows forwarding ports, using X, transferring files and other useful features [4 Condor user manual].

It is understood that the local batch represents an important part of IF processing in its own right, however most of the heavy duty processing will be performed on GRID resources. GRID resources include two important categories for Fermilab users: 1) the General Purpose GRID Cluster, and 2) the GRID at large. The GP Grid resources are characterized by their access to the centrally served file systems also available on the interactive and local batch systems. This makes configuring and running jobs on all three resources quite similar and straightforward for users. The second GRID category is potentially much larger and includes CDF, CMS and OSG computing sites beyond Fermilab. These nodes do not have the centrally served disk mounts and therefore have no immediate access to experiment software and data. An overview of the configuration for each of these computing platforms is summarized in Table 1.

Table 1. Overview of resources and configurations for the Intensity Frontier computing model.

	Interactive Login	Local Batch	"Local" GRID (FermiGRID)	GRID @ Large (OSG)
Description	Interactive user login	Condor batch pool	GRID Condor pool	Opportunistic facilities
Facility names or examples	GPCF VM	GPCF Worker Nodes	GP GRID CDF GRID	DO Grid CMS GRID GRID facilities beyond Fermilab
Machine config.	Virtual Machines	Bare Metal	Bare Metal	Unknown
Batch job submission and access	Job submission to Loc. Batch, Loc. Grid, or Grid@Large	Access to running jobs for evaluation and debugging condor_ssh_to_job	No access to running jobs.	No access to running jobs.
User home areas	/afs	/afs	none	none
Grid provided mounts of centrally served disk	/grid/data /grid/fermiapp /grid/app	/grid/data /grid/fermiapp /grid/app	/grid/data /grid/fermiapp /grid/app	BestMan gateway or via Parrot. (/grid/* are on DO clus.)
Experiment mounts of centrally served disk	/ <code>&lt;exp&gt;</code> /data / <code>&lt;exp&gt;</code> /app	All / <code>&lt;exp&gt;</code> /data, / <code>&lt;exp&gt;</code> /app mounts	All / <code>&lt;exp&gt;</code> /data, / <code>&lt;exp&gt;</code> /app mounts	no central served mounts
Enstore archive	/pnfs/ <code>&lt;exp&gt;</code> and dCache	/pnfs/ <code>&lt;exp&gt;</code> and dCache	?	SRM
SRMCP				
Users/Jobs run as.	<code>&lt;exp&gt;</code> user list	Union of all <code>&lt;exp&gt;</code> user yp lists.	<code>&lt;exp&gt;</code> ani <code>&lt;exp&gt;</code> pro <code>&lt;exp&gt;</code> cal <code>&lt;exp&gt;</code> gli	<code>&lt;exp&gt;</code> ani <code>&lt;exp&gt;</code> pro <code>&lt;exp&gt;</code> cal <code>&lt;exp&gt;</code> gli
Software tools	User products	User products	User products	
Experiment software	/grid/fermiapp / <code>&lt;exp&gt;</code> /app	/grid/fermiapp / <code>&lt;exp&gt;</code> /app	/grid/fermiapp / <code>&lt;exp&gt;</code> /app	

### Storage Evolution

Although the current central data storage solution has been very successful, demand will soon exceed supply and it may not be best suited for all applications. Alternative disk storage solutions are being explored that are anticipated to be cheaper, while supplying adequate performance at manageable support levels. Work in this area is described at [5 disk storage research]. On-demand caching from a tape archive can be provided via a SAM station and this is being set up as a trial for MINERvA on their cluster [6 minerva SAM caching experience]. If this provides value added for them this model may be useful for some kinds of data access.

## **Moving to the GRID at large**

Moving jobs to the GRID at large will involve additional steps to make the release software and experiment data available at the GRID facility. In some cases certain GRID sites have close associations with specific experiments and the software release might be installed at the site. This may be true for data as well, although access through SRMCP should make moving data to and from Fermilab straightforward. It is believed that SRMCP could replace cpn for local jobs as well (this needs to be confirmed), if so this would provide a useful generalization for interactive, local batch and GRID jobs. Pre-staging input data and managing output data are tasks that need to be accomplished and work is needed to explore existing solutions used by Run II, LHC, and others. Such tools will be adopted if possible, or alternative approaches developed if needed.

To enable truly opportunistic computing on any given GRID site, experiments will be required to provide some form of packaging of their software that will provide portability. There are many examples of tarball-like products that do this [7 MiniBooNE, D0?, CDF?, CMS DARball]. Of course validation is an important activity for this kind of data processing effort.

Some data processing may require access to conditions data, i.e. calibration, alignment, et cetera. Generally remote access to a central database service at Fermilab does not scale well, and replicating databases is impractical. Solutions to this via conditions files moved with the data can be used and SQLite is very convenient for this. A scalable approach like FroNTier can be provided if the experiments conditions data Interval Of Validity (IOV) constraints meet certain criteria [8].

## **Monitoring**

Several tools are in place for monitoring, but additional work in this area will make the system more robust and manageable by a small operations team. This includes system, storage, condor batch, and GRID job monitoring. Some simple tools like Ganglia, Condor pool monitor, and user level Condor monitoring developed by MINOS will be employed. Additional areas need to be explored.

## **Workflow Management**

Processing can be quite complex as the simulation and analysis proceed. Providing a framework for this will empower the experiments and discourage each group from inventing their own specific tools. Many such frameworks already exist, such as Condor's DAGMan and those developed by Run II and LHC. Work needs to be done in this area to provide tools that meet the needs of the IFRONT groups without adding unnecessary complexities to their tasks.

## **Other Services**

There are several additional services needed by experiments including, for example Apache, Tomcat, and database. These are currently being provided in a piecemeal approach on servers specific to each experiment. One example of consolidating the services for many experiments is the SAM Oracle database and SAMDB (python) server. In this case, many instances of Oracle are installed on one machine (actually one for Development, and one for Production). Some of the other database services

now in use by various groups are described in Table 2. Some experiments have requested their own apache servers so they can develop and maintain web pages and applications. Services like FroNTier (TOMCAT+SQUID) also need server machines, and high availability is a requirement. Most of these applications require minimal CPU and storage, and are highly suited to deployment on Virtual Machines. The deployment and support details for these VM's need to be established soon.

Table 2. Existing database servers for IF experiments.

Exp.\DB Service	Conditions	Construction	File Catalog
MINERvA	Fnalmysqlprd, fnalmysqldev, Minervadbprod, If0ra2 (dev)		IFORA1,2
NOvA	Novadbprod, novadbdev	Novadbprod, novadbdev	IFORA1,2
MINOS	MINOS-MYSQL1		IFORA1,2
MIPP	mippdbsrv01, mippdbsrv02		

## References

- [1] BlueArc
- [2]cpn
- [3]gpcf
- [4]condor user manual
- [5]disk storage research
- [6] minerva SAM caching experience
- [7] MiniBooNE, D0?, CDF?, CMS DARball
- [8] Interval Of Validity criteria

Notes from Steve Timm

First of all, in the row "Experiment mounts of centrally served disk" everything that is available on the General Purpose Grid cluster is also available on the CDF Grid cluster, this includes all the /minos, /lbne, /argoneut, /minerva, etc. mounts. It is still FermiGrid's goal to get the other 2 clusters in FermiGrid (D0 and CMS) to mount all of those too. CMS refused to do it in 2009 but that was before Art's scripts came into being and now 1 1/2 years of stable operations have been done by the IF experiments. You might better characterize the final 2 columns of the table as "FermiGrid" (which is its name after all) and "OSG". the /grid/app, /grid/fermiapp, and /grid/data are also available on the D0 cluster. On the wider OSG (and on D0, CMS clusters here) all of these volumes can be accessed via the BestMan gateway or via Parrot.

Second on the row "enstore archive", access to public stken dCache is available from all nodes in FermiGrid and the right clients are installed. for that matter gsidcap access is available from off-site but it is not recommended to be used. We would be willing to put up a few nodes that could access the pnfs volumes directly via encp but there hasn't been demand for this up to now. On the wide-area OSG the files can be accessed via SRM.

Third on the "user jobs run as"--for the OSG it is the same as "local grid:

Fourth--for software tools and experimental products, again parrot is available. Also each OSG site must have the equivalent of /grid/app, although it is not shared with ours, it is possible to load your applications on it via gridftp and fork jobs.

Re Monitoring--a number of grid tools already have monitoring plugins available for Nagios, which is FEF's monitoring tool of choice. this needs to be explored.

In the .pdf file I got, the end note numbers came through but the references that they referred to didn't.