

Tier3 Setup at Univ. Wisconsin-Madison - With PROOF, PQ2, Condor

Neng Xu, Wen Guan, German Montoya, Sau Lan Wu

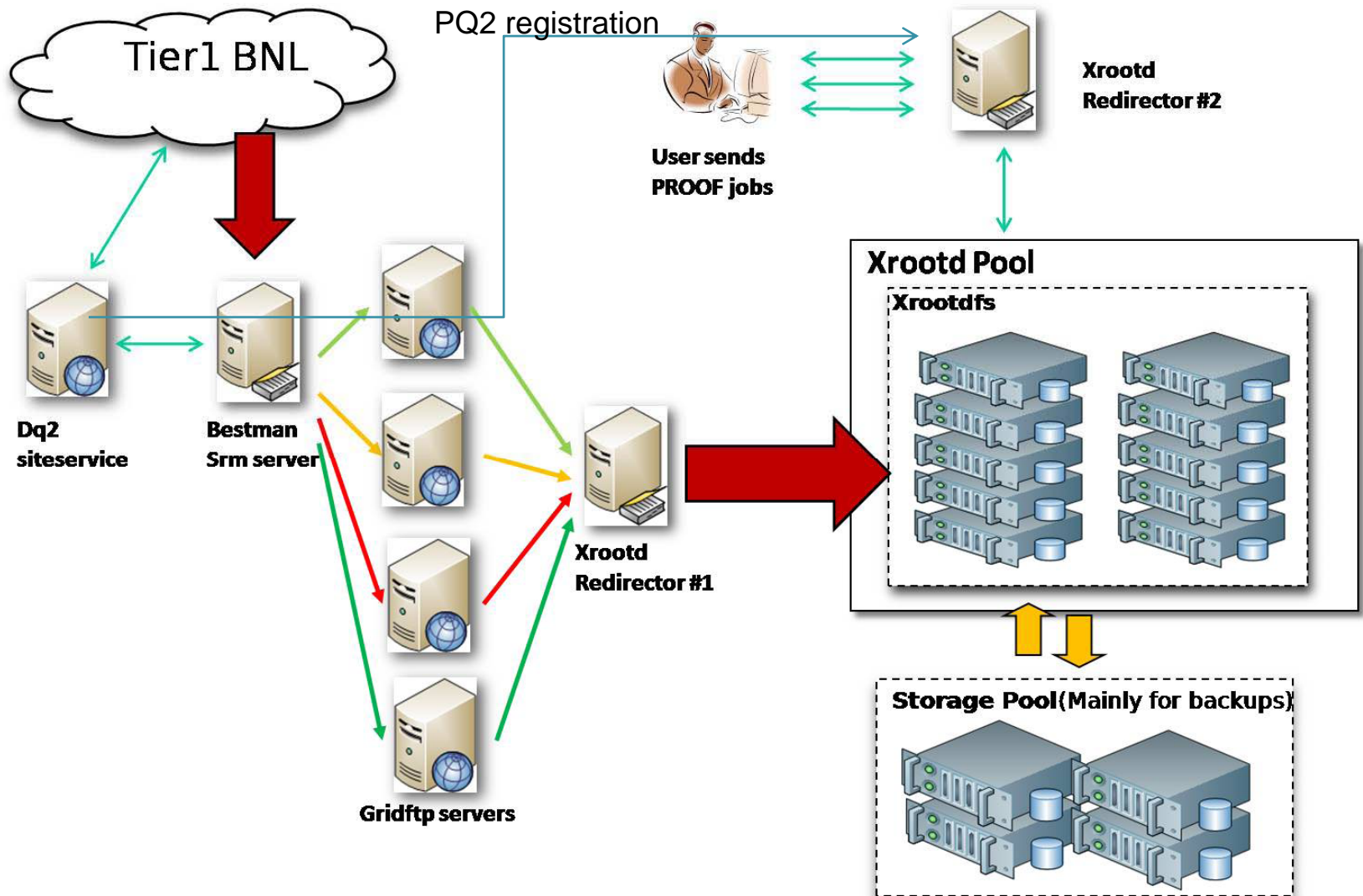
Outline

- Simple introduction of our Tier3 Facility.
- Data management with PQ2 tools.
- Our experience with PROOF
- Our experience with Condor-dagman.
- The use of ““opportunistic” CPU resources.
(Run ToyMC and Analysis with Windows machine or other OS cluster.)

Simple introduction of our Tier3 Facility

- We are using BestmanSrm + Xrootd for data transfer and storage.
- Our hardware is 50x(8core+16GB+8x750GB).
- Our network is 10Gb uplink and 1Gb to each node.
- We are using PATHENA, PROOF and Condor- dagman for data analysis.
- We run Xrootd and PROOF separately but on same storage pool.
- We were using MySQL for local data management but now we are using PQ2 tools.

Our system design



Introduction of PQ2 tools

- It's a DQ2 like local data management system by using PROOF dataset management function.
- It can be used to manage local users' data and it can also synchronize with LFC and DQ2.
- File information is stored locally on PROOF master. It can easily replicate to different PROOF master.
- Users only deal with datasets for both PROOF and Condor dagman job submission.
- The files are pre-located and pre-validated. This will save lots of overhead for the PROOF sessions, especially for the large datasets.

Basic Commands

- For system admin:
 - **pq2-put** (Register a dataset.)
 - **pq2-verify** (File location check and pre-validate.)
 - **pq2-info-server** (list all the storage nodes' information.)
 - **pq2-ls-files-server** (list all the files on one of the storage nodes.)
 - **pq2-rm** (remove datasets)
- For users:
 - **pq2-ls** (List the datasets.)
 - **pq2-ls-files** (List the file information of a dataset.)

Where to find PQ2 tools

- PQ2 is part of ROOT now. Check out the newest SVN version.
- It's in `$ROOTSYS/etc/proof/util/pq2`.
- To setup the environment:
 1. Setup the ROOT.
 2. `export PATH= $ROOTSYS/etc/proof/util/pq2:$PATH`
 3. `export PROOFURL="nengxu@atlas-bkp2.cs.wisc.edu:2093"`

What we had before PQ2

- We built our own MySQL database and register all the file information to this database.
- We ran daemons on each Xrootd data node to register new files to the database.
- Users run the scripts to get the file location and information from the database before submitting analysis jobs.

PQ2 vs MySQL

- We don't need to run separate cron daemons on each storage node. We have to always worry if those scripts are running correctly or not.
- Add/remove datasets is much simpler with PQ2.
- We don't need to maintain our own scripts and don't need to maintain a MySQL database. PQ2 is part of PROOF now. Saved 1 FTE.
- It's easy to make load balance for the PROOF master.
- Users can even synchronize the file information to their desktops/laptops.

Use PQ2 tools with DQ2/LFC in Wisconsin

- Since we are running DQ2 site service there. we wrote a script which reads out all the file information from DQ2/LFC and converts srm url of PFN to xrootd url.
- For each dataset, we create a simple text file.
- pq2-put command takes those text files as inputs and registers them into PROOF.

Use PQ2 tools with dq2-get at CERN

- At CERN, we have a small Xrootd/PROOF storage pool with XrootdFS setup.
- We can directly run dq2-get command over the XrootdFS system.
- For the finished datasets, we just simply create a text file with the xrootd url.
- pq2-put takes those text files to register the files into the PROOF.

Example of pq2-ls

```
root@atlas-bkp2 ~# pq2-ls | grep e371
/default/nengxu/mc08.106011.gg2WW0240_JIMMY_WW_enuenu.recon.NTUP.e371_s462_r563_tid029767|      4 | /CollectionTree |      10|      21 MB | 100 %
/default/nengxu/mc08.106011.gg2WW0240_JIMMY_WW_enuenu.recon.NTUP.e371_s462_r563_tid029917|     36 | /CollectionTree |     87|     184 MB | 100 %
/default/nengxu/mc08.106012.gg2WW0240_JIMMY_WW_enuunu.recon.NTUP.e371_s462_r563_tid029938|      4 |           N/A      |           |      190 MB |  0 %
/default/nengxu/mc08.106012.gg2WW0240_JIMMY_WW_enuunu.recon.NTUP.e371_s462_r563_tid029939|     36 |           N/A      |           |      1 GB |  0 %
/default/nengxu/mc08.106013.gg2WW0240_JIMMY_WW_enuaunu.recon.NTUP.e371_s462_r563_tid029940|      4 |           N/A      |           |      190 MB |  0 %
/default/nengxu/mc08.106013.gg2WW0240_JIMMY_WW_enuaunu.recon.NTUP.e371_s462_r563_tid029941|     36 | /CollectionTree |     89|     187 MB | 100 %
/default/nengxu/mc08.106014.gg2WW0240_JIMMY_WW_munumuunu.recon.NTUP.e371_s462_r563_tid029942|      4 |           N/A      |           |      190 MB |  0 %
/default/nengxu/mc08.106014.gg2WW0240_JIMMY_WW_munumuunu.recon.NTUP.e371_s462_r563_tid029943|     36 |           N/A      |           |      1 GB |  0 %
/default/nengxu/mc08.106015.gg2WW0240_JIMMY_WW_munuenu.recon.NTUP.e371_s462_r563_tid030116|      4 | /CollectionTree |      7|      15 MB | 100 %
/default/nengxu/mc08.106015.gg2WW0240_JIMMY_WW_munuenu.recon.NTUP.e371_s462_r563_tid030117|     36 |           N/A      |           |      1 GB |  0 %
/default/nengxu/mc08.106016.gg2WW0240_JIMMY_WW_munutaunu.recon.NTUP.e371_s462_r563_tid029944|      4 |           N/A      |           |      190 MB |  0 %
/default/nengxu/mc08.106016.gg2WW0240_JIMMY_WW_munutaunu.recon.NTUP.e371_s462_r563_tid029945|     36 |           N/A      |           |      1 GB |  0 %
/default/nengxu/mc08.106017.gg2WW0240_JIMMY_WW_taanutaunu.recon.NTUP.e371_s462_r563_tid029946|      4 |           N/A      |           |      190 MB |  0 %
/default/nengxu/mc08.106017.gg2WW0240_JIMMY_WW_taanutaunu.recon.NTUP.e371_s462_r563_tid029947|     36 |           N/A      |           |      1 GB |  0 %
/default/nengxu/mc08.106018.gg2WW0240_JIMMY_WW_taanuenu.recon.NTUP.e371_s462_r563_tid029948|      4 |           N/A      |           |      190 MB |  0 %
/default/nengxu/mc08.106018.gg2WW0240_JIMMY_WW_taanuenu.recon.NTUP.e371_s462_r563_tid029949|     36 |           N/A      |           |      1 GB |  0 %
/default/nengxu/mc08.106019.gg2WW0240_JIMMY_WW_taanumuunu.recon.NTUP.e371_s462_r563_tid029958|      4 |           N/A      |           |      190 MB |  0 %
/default/nengxu/mc08.106019.gg2WW0240_JIMMY_WW_taanumuunu.recon.NTUP.e371_s462_r563_tid029959|     36 |           N/A      |           |      1 GB |  0 %
[root@atlas-bkp2 ~]#
```

Example of pq2-ls-files

```
root@atlas-bkp2:~  
[root@atlas-bkp2 ~]# pq2-ls-files /default/nengxu/mc08.106012.gg2WW0240_JIMMY_WW_enumunu.recon.NTUP.e371_s462_r563_tid029939  
pq2-ls-files: dataset '/default/nengxu/mc08.106012.gg2WW0240_JIMMY_WW_enumunu.recon.NTUP.e371_s462_r563_tid029939' has 36 files  
pq2-ls-files: # File Size #Objs ObjType|Entries, ...  
pq2-ls-files: 1 root://c136.chtc.wisc.edu//atlas/xrootd/users/montoya/NTUP/mc08.106012.gg2WW0240_JIMMY_WW_enumunu.recon.NTUP.e371_s462_r563  
tid029939/NTUP.029939.00005.pool.root.2 5 MB 1 CollectionTree|TTree|250  
pq2-ls-files: 2 root://c116.chtc.wisc.edu//atlas/xrootd/users/montoya/NTUP/mc08.106012.gg2WW0240_JIMMY_WW_enumunu.recon.NTUP.e371_s462_r563  
tid029939/NTUP.029939.00006.pool.root.2 5 MB 1 CollectionTree|TTree|250  
pq2-ls-files: 3 root://c116.chtc.wisc.edu//atlas/xrootd/users/montoya/NTUP/mc08.106012.gg2WW0240_JIMMY_WW_enumunu.recon.NTUP.e371_s462_r563  
tid029939/NTUP.029939.00007.pool.root.1 4 MB 1 CollectionTree|TTree|250  
pq2-ls-files: 4 root://c127.chtc.wisc.edu//atlas/xrootd/users/montoya/NTUP/mc08.106012.gg2WW0240_JIMMY_WW_enumunu.recon.NTUP.e371_s462_r563  
tid029939/NTUP.029939.00008.pool.root.1 5 MB 1 CollectionTree|TTree|250  
pq2-ls-files: 5 root://c129.chtc.wisc.edu//atlas/xrootd/users/montoya/NTUP/mc08.106012.gg2WW0240_JIMMY_WW_enumunu.recon.NTUP.e371_s462_r563  
tid029939/NTUP.029939.00009.pool.root.2 4 MB 1 CollectionTree|TTree|250  
pq2-ls-files: 6 root://c097.chtc.wisc.edu//atlas/xrootd/users/montoya/NTUP/mc08.106012.gg2WW0240_JIMMY_WW_enumunu.recon.NTUP.e371_s462_r563  
tid029939/NTUP.029939.00010.pool.root.2 5 MB 1 CollectionTree|TTree|250  
pq2-ls-files: 7 root://c125.chtc.wisc.edu//atlas/xrootd/users/montoya/NTUP/mc08.106012.gg2WW0240_JIMMY_WW_enumunu.recon.NTUP.e371_s462_r563  
tid029939/NTUP.029939.00011.pool.root.2 4 MB 1 CollectionTree|TTree|250  
pq2-ls-files: 8 root://c115.chtc.wisc.edu//atlas/xrootd/users/montoya/NTUP/mc08.106012.gg2WW0240_JIMMY_WW_enumunu.recon.NTUP.e371_s462_r563  
tid029939/NTUP.029939.00012.pool.root.1 4 MB 1 CollectionTree|TTree|250  
pq2-ls-files: 9 root://c114.chtc.wisc.edu//atlas/xrootd/users/montoya/NTUP/mc08.106012.gg2WW0240_JIMMY_WW_enumunu.recon.NTUP.e371_s462_r563  
tid029939/NTUP.029939.00013.pool.root.2 4 MB 1 CollectionTree|TTree|250  
pq2-ls-files: 10 root://c128.chtc.wisc.edu//atlas/xrootd/users/montoya/NTUP/mc08.106012.gg2WW0240_JIMMY_WW_enumunu.recon.NTUP.e371_s462_r563  
tid029939/NTUP.029939.00014.pool.root.2 5 MB 1 CollectionTree|TTree|250  
pq2-ls-files: 11 root://c093.chtc.wisc.edu//atlas/xrootd/users/montoya/NTUP/mc08.106012.gg2WW0240_JIMMY_WW_enumunu.recon.NTUP.e371_s462_r563  
tid029939/NTUP.029939.00015.pool.root.1 5 kB 1 CollectionTree|TTree|-1  
pq2-ls-files: 12 root://c115.chtc.wisc.edu//atlas/xrootd/users/montoya/NTUP/mc08.106012.gg2WW0240_JIMMY_WW_enumunu.recon.NTUP.e371_s462_r563  
tid029939/NTUP.029939.00016.pool.root.1 5 MB 1 CollectionTree|TTree|250  
pq2-ls-files: 13 root://c138.chtc.wisc.edu//atlas/xrootd/users/montoya/NTUP/mc08.106012.gg2WW0240_JIMMY_WW_enumunu.recon.NTUP.e371_s462_r563  
tid029939/NTUP.029939.00017.pool.root.1 4 MB 1 CollectionTree|TTree|250  
pq2-ls-files: 14 root://c102.chtc.wisc.edu//atlas/xrootd/users/montoya/NTUP/mc08.106012.gg2WW0240_JIMMY_WW_enumunu.recon.NTUP.e371_s462_r563  
tid029939/NTUP.029939.00018.pool.root.1 5 MB 1 CollectionTree|TTree|250  
pq2-ls-files: 15 root://c091.chtc.wisc.edu//atlas/xrootd/users/montoya/NTUP/mc08.106012.gg2WW0240_JIMMY_WW_enumunu.recon.NTUP.e371_s462_r563  
tid029939/NTUP.029939.00019.pool.root.1 5 MB 1 CollectionTree|TTree|250  
pq2-ls-files: 16 root://c117.chtc.wisc.edu//atlas/xrootd/users/montoya/NTUP/mc08.106012.gg2WW0240_JIMMY_WW_enumunu.recon.NTUP.e371_s462_r563  
tid029939/NTUP.029939.00020.pool.root.1 5 MB 1 CollectionTree|TTree|250  
pq2-ls-files: 17 root://c117.chtc.wisc.edu//atlas/xrootd/users/montoya/NTUP/mc08.106012.gg2WW0240_JIMMY_WW_enumunu.recon.NTUP.e371_s462_r563  
tid029939/NTUP.029939.00021.pool.root.1 5 MB 1 CollectionTree|TTree|250  
pq2-ls-files: 18 root://c104.chtc.wisc.edu//atlas/xrootd/users/montoya/NTUP/mc08.106012.gg2WW0240_JIMMY_WW_enumunu.recon.NTUP.e371_s462_r563  
tid029939/NTUP.029939.00022.pool.root.1 5 MB 1 CollectionTree|TTree|250  
pq2-ls-files: 19 root://c104.chtc.wisc.edu//atlas/xrootd/users/montoya/NTUP/mc08.106012.gg2WW0240_JIMMY_WW_enumunu.recon.NTUP.e371_s462_r563  
tid029939/NTUP.029939.00023.pool.root.1 5 MB 1 CollectionTree|TTree|250  
pq2-ls-files: 20 root://c109.chtc.wisc.edu//atlas/xrootd/users/montoya/NTUP/mc08.106012.gg2WW0240_JIMMY_WW_enumunu.recon.NTUP.e371_s462_r563  
tid029939/NTUP.029939.00024.pool.root.1 4 MB 1 CollectionTree|TTree|250  
pq2-ls-files: 21 root://c115.chtc.wisc.edu//atlas/xrootd/users/montoya/NTUP/mc08.106012.gg2WW0240_JIMMY_WW_enumunu.recon.NTUP.e371_s462_r563  
tid029939/NTUP.029939.00025.pool.root.1 5 MB 1 CollectionTree|TTree|250
```


How do users use PQ2 tools

- For PROOF users, they just need to put the name of the dataset like this:

```
p-> Process("/default/nengxu/mc08.109067.PythiaH190zz4l.merge.NTUP.e384_s462_r635_t53_tid056686",  
           "CollectionTree.C+",options);
```

- For Condor-dagman users, they submit their analysis jobs like this:

```
python submit.py --initdir=/home/wguan/dag_pq2/test3 --parentScript=/home/wguan/dag_pq2/dump.sh --  
childScript=/home/wguan/dag_pq2/merge.sh --dataset=/default/nengxu/mc08.109999.PythiaQCDbb2l.fullsim.v14022303 --  
transferInput=/home/wguan/dag_pq2/input.tgz --jobsPerMachine=2
```

New features in PROOF

- Session queuing system. (Limit number of running sessions and waiting sessions.)
- PROOF-Batch. (Run PROOF as a batch system.)
- PROOF-Lite. (Make use of all your CPU cores on your desktop.)
- PQ2 tools. (New data management system with PROOF.)

Our development with Condor-dagman

- Check the machine status before submit the jobs.
- The way to setup the fast job slots.
- Setup the timeout for the child jobs.
- Limit the number of jobs to each node.
- Change the interface from MySQL to PQ2 tools.

The use of “opportunistic” resources

- Following certain steps, we found the ToyMC and Analysis code(ROOT) can be compiled and run on different operating systems.
- We can run those jobs on many Windows machines (with Condor and Cygwin installed) in the campus.
- This saves lots of CPU time for our limited SLC4 CPU resources.
- No additional software needed on the remote Windows nodes, like compilers. We just need to ship the binary file and some libraries.
- PROOF client can also run on those resources, too.

Summary

- There are many available methods to let users to run analysis jobs. PROOF, PROOF-batch, PROOF-Lite, Condor-dagman, PATHENA.
- PQ2 tools is a handy software package for Tier3s.
- ROOT works fine on many operating systems. Try to run your ToyMC and analysis jobs on Windows machines.
- Wiki pages will be available soon. If you are interested in what we are doing, you can contact me (neng.xu@cern.ch) or ROOT/PROOF team.

Thanks to

- ROOT/PROOF team.
- Condor team.
- DDM team at BNL.
- SLAC Xrootd team.