Data Analysis with PROOF, PQ2, Condor

Neng Xu, Wen Guan, Sau Lan Wu University of Wisconsin-Madison

30-October-09 ATLAS Tier 3 Meeting at ANL

Outline

- Dataset management with PQ2 tools.
- Data analysis at Wisconsin.
- Running PROOF and Condor at BNL.

Introduction of PQ2 tools

- It's a local data management system using PROOF dataset management function.
- Users only deal with datasets instead of file lists.
- File information is stored locally on PROOF master. No database is needed.
- It can be used for both PROOF and Condor data analysis 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. (One of our users was running a PROOF session with 60000 files, validation time is almost 15 minutes.)

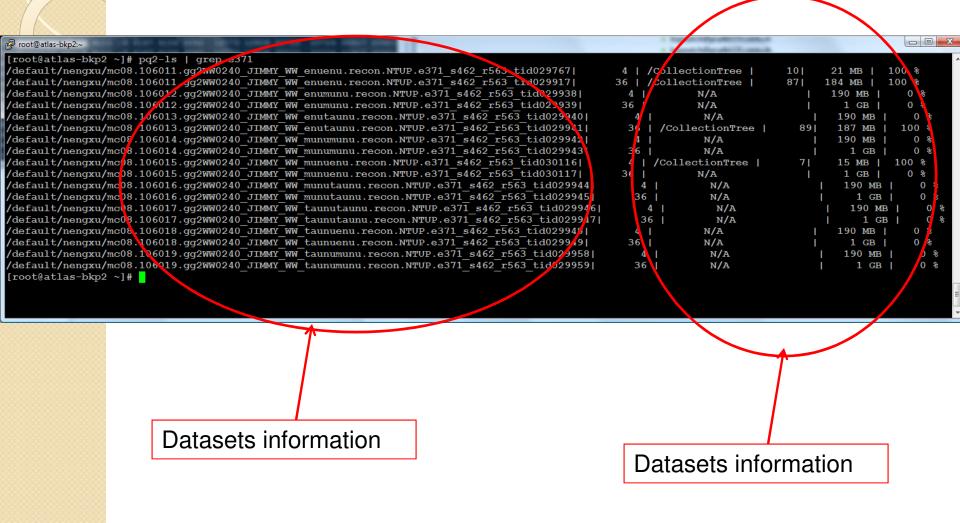
How to PQ2 tools

- PQ2 is part of ROOT(newer than 5.24) now.
- http://root.cern.ch/drupal/content/pq2-tools
- It's in \$ROOTSYS/etc/proof/util/pq2.
- To setup the environment:
 - I. Setup the ROOT environment.
 - 2. export PATH= \$ROOTSYS/etc/proof/util/pq2:\$PATH
 - 3. export PROOFURL="userid@proof-redirector.xxx.xxx:1093"

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.)

Example of pq2-ls



Example of pq2-ls-files

```
_ D X

√ root@atlas-bkp2:~

[root@atlas-bkp2 ~]# pq2/1s-files /defamlt/nengxu/mc08.106012.gg2WW0240 JIMMY WW enumunu.recon.NTUP.e371 s462 r563 tid029939
                        default/nengxu/mc08.106012.gg2WW0240 JIMMY WW enumunu.recon.NTUP.e371 s462 r563 tid029939' has 36 files/
pg2-ls-files: dataset
pq2-ls-files:
                                                                                                                         #Objs Obj|Type|Entries, ...
pq2-ls-files:
                       root://c136.chtc.wisd.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:
                       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:
                       root://c116.chtc.wisc.etu//atlas/xrootd/users/montoya/Ngor/mc08.106012.gg2WW0240 JIMMY WW enumunu.recon.NTUP.e371
tid029939/NTUP.029
                  939.
                       00007.pool.root.1
                                                           1 CollectionTree Tree 250
                       root://c127.chtc.wisc.edu//atlas/xrootd/users/montoya/NTUP/mc0
pg2-ls-files:
                                                                                        .106012.gg2WW0240 JIMMY WW enumunu.recon.NTUP.e371 s462 r563
tid029939/NTUP.029939.
                       00008.pool.root.1
                                                 5 MB
                                                           1 CollectionTree TTree 250
                       root://c129.chtc.wisc.edu//atlas/xrootd/users/montoya/NTUP/mc08.
pq2-ls-files:
                                                                                          .06012.gg2WW0240 JIMMY WW enumunu.recon.NTUP.e371 s462 r563
tid029939/NTUP.02
                 9939. 00009.pool.root.2
                                                           1 CollectionTree | TTree | 250
pq2-ls-files:
                       root://c097.chtc.wisc.edu
                                                 //atlas/xrootd/users/montoya/NTUP/mc08.1
                                                                                           06012.gg2WW0240 JIMMY WW enumunu.recon.NTUP.e371 s462 r563
tid029939/NTUP.029939. 00010.pool.root.2
                                                 5 MB
                                                           1 Collection ree TTree 250
                                                 //atlas/xrootd/users/mon
pq2-ls-files:
                       root://c125.chtc.wisc.edu
                                                                                           6012.gg2WW0240 JIMMY WW enumunu.recon.NTUP.e371 s462 r563
                                                                          toya/NTUP/mc08.10
tid029939/NTUP.02
                 9939.
                       00011.pool.root.2
                                                 4 MB
                                                           1 Collection ree | TTree | 250
                       root://c115.chtc.wisc.edh//atlas/xrootd/users/mo<mark>n</mark>toya/NTUP/mc08.1<mark>0</mark>6012.gg2WW0240 JIMMY WW enumunu.recon.NTUP.e371 s462 r563
pq2-ls-files:
                   8
tid029939/NTUP.02
                  9939.
                       00012.pool.root.1
                                                           1 Collection
                                                                          ree|TTree|250
                       root://c114.chtc.wisc.edu//atlas/xrootd/users/mor
pg2-ls-files:
                                                                          toya/NTUP/mc08.1
                                                                                            6012.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:
                       root://c128.chtc.wisc.edu//atlas/xrootd/users/montoya/NTUP/mc08.
                                                                                           06012.gg2WW0240 JIMMY WW enumunu.recon.NTUP.e371 s462 r563
tid029939/NTUP.029939. 00014.pool.root.2
                                                           1 CollectionTree | TTree | 250
pq2-ls-files:
                       root://c093.chtc.wisc
                                              edu//atlas/xrootd/users/montoya/NTUP/mc08
                                                                                          106012.gg2WW0240 JIMMY WW enumunu.recon.NTUP.e371 s462 r563
tid029939/NTUP.02993
                       00015.pool.root.1
                                                 5 kB
                                                           1 CollectionTree | TTree | -1
pg2-ls-files:
                  12
                       root://c115.chtc.wi/c.edu//atlas/xrootd/users/montoya/NTUP/mc/8.106012.gg2WW0240 JIMMY WW enumunu.recon.NTUP.e371 s462 r563
tid029939/NTUP.029939.
                       00016.pool.root.1
                                                 5 MB
                                                           1 CollectionTree Tree 250
                       root://c138.chtc/wisc.edu//atlas/xrootd/users/montoya/NTUP
pq2-ls-files:
                                                                                    mc08.106012.gg2WW0240 JIMMY WW enumunu.recon.NTUP.e371 s462 r563
                       00017.pool.root.1
tid029939/NTUP.029939.
                                                           1 CollectionTree | TTree
                       root://cloc.chtc.wisc.edu//atlas/xrootd/users/montoya/NTUP
pg2-ls-files:
                                                                                    mc08.106012.gg2WW0240 JIMMY WW enumunu.recon.NTUP.e371 s462 r563
                       00018.poel.root.1
tid029939/NTUP.029939.
                                                 5 MB
                                                           1 CollectionTree | TTree
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
                       root://cl17.chtc.wisc.edu//atlas/xrootd/users/montoya/NTUP
pq2-ls-files:
                  16
                                                                                    mc08.106012.gg2WW0240 JIMMY WW enumunu.recon.NTUP.e371 s462 r563
tid029939/NTUP.029939.
                       00020.pobl.root.1
                                                 5 MB
                                                           1 CollectionTree TTree
                       root://c117.chtc.wisc.edu//atlas/xrootd/users/montoya/NTUP
pg2-ls-files:
                  17
                                                                                    mc08.106012.gg2WW0240 JIMMY WW enumunu.recon.NTUP.e371 s462 r563
tid029939/NTUP.029939
                       00021.pool.root.1
                                                 5 MB
                                                           1 CollectionTree TTree
                       root://c104.chtc.wisc.edu//atlas/xrootd/users/montoya/NTUP
pq2-ls-files:
                  18
                                                                                    mc08.106012.gg2WW0240 JIMMY WW enumunu.recon.NTUP.e371 s462 r563
                       00022.pcol.root.1
tid029939/NTUP.029939.
                                                           1 CollectionTree|TTree
                                                 5 MB
pg2-ls-files:
                       root://d104.chtc.wisc.edu//atlas/xrootd/users/montoya/NTUP
                                                                                    mc08.106012.gg2WW0240 JIMMY WW enumunu.recon.NTUP.e371 s462 r563
tid029939/NTUP.029939.
                       00023.pc ol.root.1
                                                 5 MB
                                                           1 CollectionTree | TTree
pq2-ls-files:
                  20
                       root://d109.chtc.wisc.edu//atlas/xrootd/users/montoya/NTUP
                                                                                    <u>mc08.106012.qq2W</u>W0240 JIMMY WW enumunu.recon.NTUP.e371 s462 r563
tid029939/NTUP.029939.
                       00024.pool.root.1
                                                 4 MB
                                                           1 CollectionTree | TTree
                       root://c115.chtc.wisc.edu//atlas/xrootd/users/montoya/NTUP
pg2-ls-files:
                                                                                    mc08.106012.gg2WW0240 JIMMY WW enumunu.recon.NTUP.e371 s462 r563
                  21
tid029939/NTUP.029939. 00025.pool.root.1
                                                              CollectionTree|TTree
```

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.v14022363-transferInput=/home/wguan/dag_pq2/input.tgz --jobsPerMachine=2

Advantages of PQ2

- Easy to maintain.
- No database needed.
- Add/remove datasets is very simple.
- Fast access to the file information.
- Avoid some overhead of PROOF jobs.
- Easy to make load balance for the PROOF master.
- Not only work with Xrootd but also NFS, Dcache, CASTOR or even local files.
- . . .

Our Tier3 Facility at Wisconsin

- 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 are using PQ2 tools for local file management.
 Recently PQ2 need to work with XrootdFS to get the file list.
- All of our users are working from CERN Neng Xu, Univ. of Wisconsin-Madison

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.

Analysis with Condor-dagman

- The submit script talks with PQ2 server to get file location and define the Condordagman JDLs.
- Check the machine status before submit the jobs.
- Limit the number of jobs to each node.
- Multi-level collecting or merging the output files.
- Example script is here:

http://wisconsin.cern.ch/~wguan/dag_pq2_v2_bnl.tgz http://wisconsin.cern.ch/~wguan/dag_pq2_v2_wisc.tgz

Our Tier3 experience at BNL

- We are going to have some CPUs and storage at BNL. How to use them is our new project.
- We are working on two ways to run analysis jobs at **BNL**:
 - Using Condor to start a PROOF pool.
 - Using Condor-dagman to run analysis jobs.
- All the input files are in the Dcache. The PROOF jobs read via the Xrootd door and Condor-dagman jobs read directly from DCache.
- We registered the datasets we need to a PQ2.
- For PROOF jobs, we use SSH tunnel from CERN to BNL without login to BNL machines. Very convenient for users.

13

Use PQ2 tools at BNL

- We can run dq2 command to get the datasets information and file SURL.
- We just simply create a text file with the xrootd url converted from SURL.
- pq2-put takes those text files to register the files into the PROOF redirector.
- Users can access the file information via SSH tunnel between CERN and BNL.

Some performance problems at BNL

- Since PROOF reads the Dcache files via Xrootd door, the performance was limited by Xrootd door. The maximum speed is 8MB/s.
- PROOF developers are working on direct access Dcache files.
- We also don't want to overload the Dcache system at BNL because PROOF is very aggressive on I/O. There should be a way to limit maximum network traffic.

Summary

- PQ2 works quite good for our local data management.
- If your Tier3 want to put machines at BNL or SLAC, the way we setup PROOF over Condor could be good example. It also works with other batch system like LSF, PBS.
- Condor-dagman is good for those users who don't want to modify their code. Multi-level merging can reduce the load on submit nodes.

•

Thanks to

ROOT/PROOF team.

Condor team.

DDM team at BNL.

SLAC Xrootd team.