# CCJ Operation in 2008-2009

T. Nakamura, H. En'yo, Y. Goto, T. Ichihara, Y. Watanabe and S. Yokkaichi

## 1  Overview

The RIKEN Computing Center in Japan (CCJ)[1] has been developed since April 1998 for analyzing the huge amounts of data collected in the PHENIX experiment at RHIC. Thus far, CCJ has been providing numerous services as Asia's regional computing center. For instance, CCJ maintains sufficient computing power by the PC cluster operated by Scientific Linux[2] for simulation and individual data analysis. The collected data are transferred through SINET3 with a 10 Gbps bandwidth, maintained by NII[3], from Brookheaven National Laboratory (BNL) by using GridFTP[4]. The transferred data are once stored in High Performance Storage System (HPSS)[5] before starting the analysis. This HPSS is one of the joint operations with the RIKEN Integrated Cluster of Clusters (RICC) project. One of the most successful achievements of the CCJ operation up to now is the completion of more than 30 doctoral dissertations with the analysis results obtained using the computing resources at CCJ.

A summary of the basic configuration of the CCJ system has been published elsewhere[6]. Since last year, several major hardware-level upgrades have been carried out. The center network switch has been replaced with Catalyst4900M, which has eight 10 Gbps ports, in July 2009. The operation of the HPSS was started with a new version 7.1 in January 2010. In February 2010, four sets of UPS were replaced by SANUPS ADS series (Sanyo Denki). The server (ccjnfs11) for Network File system (NFS), which managed a 6.8 TB FC-FC RAID box and a 8.9 TB FC-SATA RAID box, was ended the service to be upgraded in March 2010. Recently, effective use of the existing computing resources has become difficult owing to the rapid increase in the PHENIX data size. In this report, the details of several upgrades and developments made to the CCJ system in the 2008-2009 period are included.

## 2  PC cluster specification

Until early 2008, CCJ was operated by approximately 190 PC nodes, 166 nodes of which were used as calculation nodes for the simulation and PHENIX data analysis. In June 2008, 112 nodes were eliminated. The remaining 36 and 18 dual CPU nodes have Intel Pentium III (1.4 GHz) CPU and Intel Pentium 4 (2.0 GHz) CPU, respectively. Each calculation node has 1 GB memories.

In February 2009, 18 PC nodes (HP ProLiant DL180 G5) were newly introduced to compensate for the low total computing power. Each node has dual CPUs (Quad-Core Intel Xeon E5430 2.66 GHz) and 16 GB memories. These nodes have twelve 3.5 inch bays of
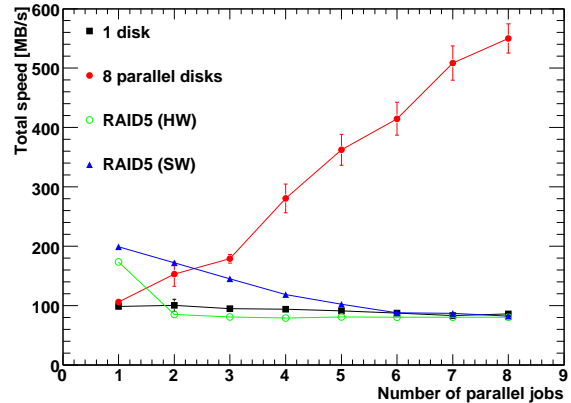


Fig. 1. Average total speed for reading 1 GB files in one disk (square), 8 parallel disks (filled circle), hardware RAID5 (open circle), and software RAID5 (triangle) as a function of number of parallel jobs.

HDD for each chassis. Two 146 GB SAS disks are mounted for the operating system and used by mirroring to reduce the down time originating from the troubles on the HDD. Initially, eight 1 TB SATA disks were mounted for local data storage. In December 2009, two additional 1 TB SATA disks were added to each node. Currently, local disks with a total capacity of 180 TB are available for the data storage. This specification is the key feature for I/O bound jobs, as will be described in the next section. Each node has a 1 Gbps network interface card, and all of the nodes are connected to a network switch (Dell PowerConnect 6224) mounted on the same rack. This network switch is up linked to the center network switch at CCJ (Catalyst 4900M) via a 10 Gbps connection. This new cluster has the capability to process a typical PHENIX analysis job approximately three times faster than the old PC nodes in a CPU core.

In October 2009, 20 PC nodes were setup as a cluster at RICC for the exclusive use of the CCJ users. Each node has dual CPUs (Quad-Core Intel Xeon X5570 2.93 GHz) and 12 GB memories and is operated by Scientific Linux 4.4 on VMware ESXi[7] so that it is dedicate to the PHENIX analysis environment. Condor[8] system is available as a batch job scheduler for this cluster.

Thus, a total of 412 CPU cores are presently available as calculation nodes for the CCJ users.

## 3  Development of job submission scheme

The PHENIX data stored once in the HPSS will be transferred to several RAID boxes for the analysis. Although users can access the data in the RAID boxes through the NFS servers, multiple access from numer-

Table 1. Summary of nDSTs in local disk.

| Dataset | nDST type | Data amount |
|---------|-----------|-------------|
| Run 9 $p + p$ 200 GeV | All type | 65.4 TB |
| Run 9 $p + p$ 500 GeV | All type | 31.2 TB |
| Run 8 $p + p$ 200 GeV | All type | 21.2 TB |
| Run 6 $p + p$ 200 GeV | w/o detector | 14.6 TB |
| Run 5 $p + p$ 200 GeV | w/o detector | 9.9 TB |

ous calculation nodes at the same time is not possible because of the decrease in the I/O speed. Therefore, users must transfer the data from the RAID boxes to the calculation node for each batch job. Since the size of PHENIX data is growing steadily, such data transfer becomes a bottleneck in data analysis. This problem is eliminated with the use of the newly introduced calculating nodes, which have large capacity local disks for storing the data a priori (see previous section). However, since these new calculating nodes have multi-core CPUs, which is predominant in the market, data analysis remains an I/O bound type job. Therefore, it is necessary to optimize the composition of the local HDD. We performed a benchmark test to evaluate the I/O performance for the new cluster. Figure 1 shows the average total speed for reading 1 GB files as a function of the number of parallel jobs. Originally, each HDD shows the I/O performance of 100 MB/s. However, the use of a multi-HDD is not advantageous with the RAID configuration, as shown by the open circles and triangles in Fig. 1. Since the RAID configuration gives us a single name space, maintaining the data location becomes easy. Nevertheless, we do not chose this configuration to maximize the I/O performance. In October 2009, approximately 96 TB Run 9 $p + p$ data, so-called nano-DST, were transferred from BNL as soon as data reconstruction was completed at the RHIC Computing Facility (RCF)[9]. They were stored in local disks along with the previously stored data. Table 1 shows a summary of the dataset in the local disks accessible to users by the batch queuing system.

In the calculation nodes, users can process their own analysis code via the batch queuing system (LSF version 7.0[10]) in the CCJ cluster. We have added some software modules to enable the user specify the nDSTs distributed to the local disks during job submission. Figure 2 shows a brief flowchart of the data-oriented job submission scheme. Since all the subsets of the PHENIX data have identical run numbers, the module first obtains the location of the nDST by the user-specified run number from a database. Then, the module submits the user jobs to the appropriate node by the LSF. To avoid multiple access for a local disk from several jobs, the module sets a lock file for exclusive access and grants permission only to the user by the Access Control List (ACL) method. As a result, each job dispatched to a calculating node exclusively handles a local disk. The advantage of this method is that the I/O performance is enhanced, as shown by
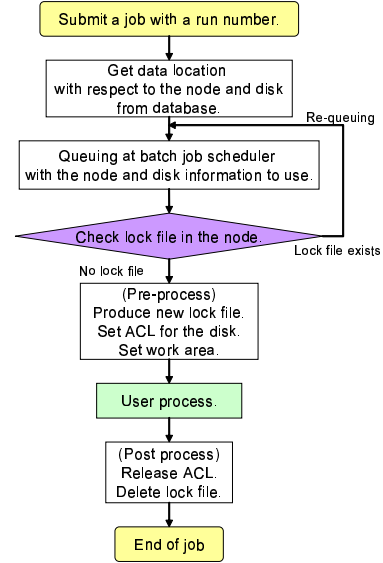


Fig. 2. Data-oriented job submission scheme.

the filled circles in Fig. 1. Further, a temporary work area for the job is assigned to the same disk with the data location. Therefore, this scheme is effective for eliminating the I/O bound problem in case of generic jobs as well *e.g.* simulation. Each job are able to save the processing time on the data transfer approximately 10 times as compared to the typical jobs in the preexisting calculating nodes. An instruction is provided in this URL[11] on how to access to the data stored in the local disks.

## 4 Summary

18 calculating nodes with 180 TB local disks were introduced for effectively analyzing huge amounts of PHENIX data. A data-oriented batch queuing system was developed as a wrapper of the LSF system to increase the total computing throughput. Indeed, the total throughput was improved by roughly 10 times as compared to that in the existing clusters; CPU power and I/O performance are increased threefold and tenfold, respectively. Thus, users can analyze data of several tens of terabytes within a few hours. This is one of the most significant developments made to the CCJ operation in 2009.

References
1) http://ccjsun.riken.jp/ccj/
2) http://www.scientificlinux.org/
3) http://www.nii.ac.jp/
4) http://www.globus.org/grid_software/data/gridftp.php
5) http://www.hpss-collaboration.org/
6) S. Yokkaichi et al., RIKEN Accel. Prog. Rep. **42**, 223 (2009).
7) http://www.vmware.com/products/esxi/
8) http://www.cs.wisc.edu/condor/description.html
9) http://www.rhic.bnl.gov/RCF/
10) http://www.platform.com/products/LSF/
11) http://ccjsun.riken.go.jp/ccj/doc/LSF/lsf-wrapper.html