# Fermilab

# Development of Remote Concatenation & Monitoring Tools for CDFII Production Farm

**Solomon Mikael**

*The Collider Detector at Fermi Laboratory*

## Abstract

The CDF detector is a vital part of the physics program at Fermi Laboratory. Monitoring the farm is the web based Tikiwiki pages. The tiki pages provide vital information to the users allowing the farm to operate smoothly. The concatenation of reconstructed files depends on farm. Modifications were made to both the monitoring and concatenation system to make both systems perform efficiently. This paper constitutes the final project for he SIST program.

## Introduction

Higher energy, higher intensity particle beams, and advanced detectors in high energy physics (HEP[1]) have led to the collection of larger volumes of data. The Collider Detector at Fermilab (CDF[2]) has improved its data acquisition capacity in the Run II program. The computing facility was also upgraded for processing larger volumes of data. The CDF II experiment began collecting data in 2000. The peak rate of data recording is 40MB/s. The total volume of raw data collected by the experiment is roughly 0.5 petabytes. This number is expected to increase dramatically in the future. As a result of this, the data management becomes a very important part of the experiment. This paper will describe the main hardware, software, monitoring, and control components of the CDF production farm.

The primary objective is to implement a method that will improve the disk write speed. Additionally modifications were made to the monitoring and concatenation systems of the farms.

## CDF Experiment

The CDF (collider detector experiment at FermiLab) is an

international collaboration involving many universities and national laboratories. The 100 ton CDF detector is a large general purpose solenoidal detector which combines precision charged particle tracking with fast projective calorimetry and fine grain muon detection. CDF must manage is made by the collisions or the protons and anti-protons. Of the collisions on a few are hard collisions, those which make energetic particles which go off at large angels from the beam allowing it to pass through the detector. A charged particle will be observed the rest can be tracked by taking into account the transverse momentum based on the conservation of momentum. Also particles that are not detected in the initial layers are detected in the calorimeter. Not all the energy can be accounted for some particles will inevitable escape detection.
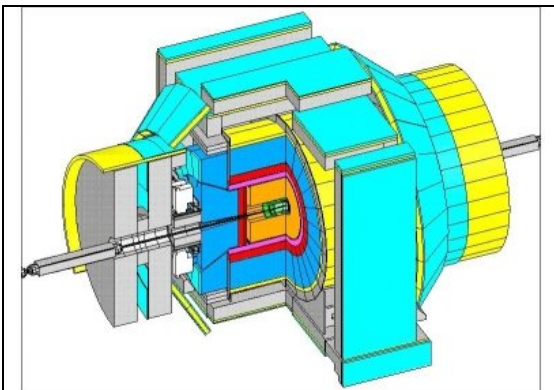


Figure 1 ( this is a diagram of the CDF detector )

To observe the particles generated by the collision there are several detector systems. Each detector system is best at detecting or interacting with a certain particle. Of the millions of collisions per second only a few of them give large angles. The detectors record electronic information on the 1000000 channels. Based on a few of the channels

preliminary information is provided. Triggers make a decision if a collision is interesting or not. As the precision and efficiency of the detector increases so is the amount of information it's able to draw from these collisions. As a result a data management system is imperative for CDF to operate.

**The CDF Farm**

The primary goal of CDF Production Farm is to perform computing and network intensive tasks in a cost effective manner. The data processing model for the CDF experiment is based on reconstruction of parallel data streams taken with different combinations of physics event triggers. Data processing consists of large clusters of Linux computers with data movement managed by CDF data handling system. The high energy physics experiments (HEP) being done at CDF generated enormous amounts of data and is expected to increase in the future. CDF's challenge was in finding a way to manage the large flow of data through the computing nodes. Condor is the software that manages, controls, and monitors the CDF production farm. Presently the production farm consists of 150 dual CPU PC's with a total computing power of about 800 GHz. The main software components of the farm architecture are SAM and CAF. $SAM_3$ stores and retrieves the metadata associated with files. CAF is control system used by all farms inside the experiment for batch job submission and access to the CDF data management system and databases. The new SAM data production system is suitable for job submission to any computing facility in the world that uses the CAF interface

with direct access to the SAM data handling and database connection. The production farm is setup to take advantage to grid computing was a CAF headnode and SAM stations. The mass storage system has three major parts: the Enstore system which is used to access data on the tapes, dCache which provides a system for storing and retrieving huge amounts of data which is distributed among a large number of nodes under a single virtual filesystem, and PNFS which is a software used to map all files stored in the root to a Unix like namespace.

**Software Products**

For CDF farms to accomplish its goals it needs many individual programs to operate. Each program contributes a job to the entire system of the farm. In the end it is this system that programmers and users are familiar with.

**I. Condor**

The Condor Project[4] is one of the systems used in the farm. The goal of the Condor Project is to develop, implement, deploy, and evaluate mechanisms and policies that support High Throughput Computing (HTC) on large collections of disruptively owned computing resources. This system enables scientists and engineers to increase their computing throughput. Condor achieves its goal by implementing a specialized work load management system for computer intensive jobs. The batch system in Condor provides a job queuing mechanism, scheduling policy, priority scheme, resource monitoring, and resource management. After submission of a job, in series or in parallel, Condor places them into a queue and chooses

when and where to run the jobs based upon a policy, monitors the jobs progress, and finally informs the user when the job is complete. Condor can be used to manage a cluster of dedicated computer nodes and harness wasted CPU power from idle desktop workstations. A big advantage of using Condor is that it does not require a shared file system across machines.

**II. CVS**

The concurrent version system also known as CVS[5] is another software tool used by the farms. CVS is a version control system. It allows one to record the history of source files. For example is there is a bug in a program after it has been modified. CVS allows for retrieval of the old version to see exactly what change caused the bug. CVS stores all the versions of a file in a single file which only stores the difference between versions. The software helps insulate different developers from each other so all developers work in their own directory avoiding overwriting others changes. The files are stored in a centralized repository which stores a complete copy of all the files and directories which are under version control.

**III. SAM**

For the CDF Farms to operate efficiently there is a need for a system that delivers files and keeps track of the files that have been analyzed. A system that does just this is called SAM or sequential data access via metadata. All the files accessible in a UNIX system area arranged in one big tree the file hierarchy is rooted at /. These files can be spread out over several devices. SAM

mitigates the problem of one person hogging the tape drive and/or flooding the tape system. SAM is a data handling system organized as a set of servers which work together to store and retrieve associated metadata ( data that is used to describe other data which may describe how, when, and by whom the data was received, created, accessed, and modified ) including a complete record of the processing which has used the files. SAM is able: to track locations and comprehensive metadata for each file in the system, provide storage utilities to add a file to a permanent storage location, deliver files on request to systems that have deployed the file delivery components of the SAM system, provide methods for job submission which are coupled to the file delivery mechanism and can utilize location and system information for performance optimization, and track processing information permitting processing to be organized on the basis of the consumption history  of a particular data set and allows the construction of processing jobs which use information which other files have processed successfully.

## IV. Data Base

Structured Query Language (SQL) implements a database and interfaces on a computer called a database server. A database is a hierarchy of increasingly complex data structures. MySQL$_7$ has an ability to handle an unlimited number of simultaneous users and handle millions or records making it essential to the CDF Farms.

A few of the requirements of the CDF production farm are high throughput rate, I/O capability, easily configurable system, and efficient error handling. To make these happen the farm must monitor the hardware performance including status reporting and tuning the software of the production farm. The production farm web interface (PFWI) will parse, calculate, and display all major characteristics of the farm. The farm's performance indicators change in real time so the interface reports online results. All steps of processing heavily rely on information from the Oracle database. The worker nodes retrieve: calibration, geometric constants, and reconstructed files. The intermediate output on the worker nodes is declared to the database in order to avoid lost or duplicated events.



Figure 2 ( the basic schematic of the CDF Farm )

The CDF farm begins raw data processing by accessing the data of the tape storage device. The collision information is stored on the tapes until this point. Then next step is reconstruction of the data. Reconstruction interprets the data that was collected to make it useful for physics analysis. The reconstructed data is then sent to the stagers. In some cases if the files are already 1-2 GB is size there is no need for concatenation. For

the instances where this is not true the concatenation begins on the stagers. The stagers then send the reconstructed and concatenated data back to the tape drives for storage and later access by physicists.

**Storage Speed**

The issue of increasing the rate the data is written to the tape storage has many factors. In the configuration when the concatenated and the tape upload are running simultaneously the tape upload speed is affected dramatically.



Figure 3 ( the x –axis represents the MB/s while the y –axis represents the number of events that were sent at that particular speed the majority of the events seem to congregate around 12 MB/s )



Figure 4 ( a majority of the events congregate around 25 MB/s )

In the present scheme the concatenator and tape uploader are running at the same time resulting in limited I/O from the stagers. Disk access is directly related to the number of simultaneous I/O operations from the disk RAID 5. The two above figures demonstrate two different scenarios. The scenario in Figure 3 demonstrates when the stager is running concatenator as well as the tape uploader. The graphs x – axis represents the transfer rate in MB/s while thy y – axis represents the number of events that were transferred at a particular speed. The majority of the events that are being sent are on average being transferred at an average rate of about 12 MB/s. In Figure 4 the concatenation is removed from the stagers so the I/O to the disk RAID 5 is limited to the transfer of the data on the stagers. The difference is dramatic. The number of events that are send at a speed of 25 MB/s or higher constitute the majority of the events. This change nearly doubles the transfer rate of the information. The essence of my project at CDF entailed removing the concatenation to the workers to enable the stagers to achieve higher transfer rates to the disk storage.
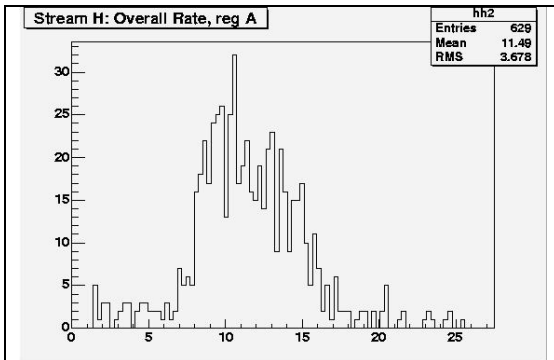
**Proposed Structure**

To improve the tape write speed on the CDF farm it was imperative the concatenation be moved to another location. The proposed structure is displayed in Figure 5.
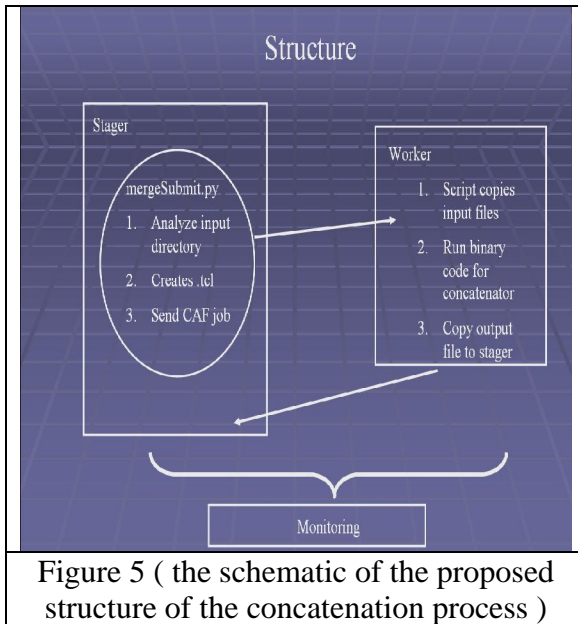
Figure 5 ( the schematic of the proposed structure of the concatenation process )

The stagers are responsible for running the python script mergeSubmit.py. This code will analyze the input directory. The reason for this is so it can make a list of the input files needed for each of the production datasets. The list of input files is inserted in a file called a Tool Command Langue (.tcl) file which drives a binary code for concatenation. Once this file has been composed a job can be submitted to CAF. CAF divides the job into sessions and assigns to the workers. The workers will be held with the responsibility of concatenation. Once on the workers the script copies the input files. It knows where the files are because of the constructed .tcl file. Additionally when the submission to the worker is made certain parameters are established to specify certain parameters of the specific submission. The binary code for the concatenation is then run. The goal is creating files that are 1-2 GB in size. The reason for this is to optimize the speed that the data is being written to the tape drives. The output files are then written back to the stager. An advantage of this process is the output file generated is on the worker. This allows

for it to be accessed by other mean instead of it being localized on a stager. For this configuration the output file is then written back to the stager. This system is intended to circumvent the problem of the RAID 5 getting numerous I/O calls which affect the rate data can be transmitted to the tape drives.

```bash
#!/bin/bash

name=`basename $0`

. ./cdfopr/scripts/common_procedures
. ./cdfopr/scripts/parse_parameters $* > temp_parse_log

echo $TCL_FILE
echo $PARAM_USER
echo $PARAM_HOST
echo $PARAM_PATH


cmd="fcp -c ${RCP}
${FCP_USER}@${FCP_HOST}:${PARAM_PATH}/${TCL_FILE}
.";
$cmd
STATUS=$?;

if[ STATUS -ne 0 ];then
    echo "$TCL_FILE was not able to be copied"
    exit 1;
fi

for file_loc in ` temp.awk –v num=$SEGMENT_NUMBER `;do

    cmd="fcp -c ${RCP}
${FCP_USER}@${FCP_HOST}:${file_loc} . ";
    $cmd
    STATUS=$?;

    if [ STATUS -ne 0 ];then
    echo "$file_loc was not able to be copide
    fi

done
```

Figure 6 (find_tcl.sh code used to extract the corresponding location of the input files and send the files to the appropriate worker )

```awk
#!/bin/awk

BEGIN {
 flag = 0;
}

/SEGMENT_NUMBER/ {
 if ($5 == seg) {
   flag = 1
 }
}
```

```
/include/ {
  if (flag == 1) {
    print $0
  }
}

{
  if ($1 == "}") {
    flag = 0
  }
#  if (($4 == "==") && (ENVIRON["SEGMENT_NUMBER"] ==
$5)) {print $0}
}
```

Figure 7 (temp.awk code that would extract the information needed from the .tcl file )

Figure 6 displays a short excerpt of the code used to perform the proposed task. The code is a combination of shell scripting and the awk editor. The awk editor helped in parsing the .tcl file to extract the location of the input files associated with a specific project. This code would be implemented right after step 2 on the stager in Figure 5. While this entire process is taking place the farm is being monitored. Monitoring is an essential part of the CDF farms. It enables the farm members to maintain the entire farm system in an efficient manner.

**Farm Control & Monitoring**

The CDF Farm is a huge system that requires the attention on many individuals. To take the most advantage of those individuals time there needs to be an intuitive interface. Farm maintenance is significantly simplified by introducing wiki-based farm projects. Tikiwiki$_8$ software is an extremely efficient tool for webbased documentation. Additionally its underlying database serverkeeps a history of all changes to the Farm Projects. The tiki pages with Project configurations give the possibility to:

keep track of all existing projects, start or stop the projects, change resource sharing between projects. Several of the pages that the farms relied on needed slight modification to improve the overall presentation of the information on the tiki pages.



Figure 8 ( tiki page initially )



Figure 9 ( tiki page after modifications )

Figure 8 shows the initial orientation of the page. This information is all presented in an unclear manner and there isn't a way to contact the individual running the particular project. To make these edits possible a few modifications needed to be made to include columns and remove others. This was all done using the online tiki editor. Another group of pages were modified to make

the information more presentable. Instead of using the online tikiwiki editor python script was edited. The script itself was made to interface with html format. Also some adjustments were made including: removing columns, adding jobid information, and fixing hyperlinks. Figures 10 and 11 show the page before and after respectively.

respective result are displayed in Figures 12 and 13.



Figure 10 ( web page initially )



Figure 11 ( web page after modifications )

Another python script that was edited is df_disk.py. The intention with this script was to display the disk space on the 32 partitions on the different servers. This expedited the process of obtaining the server information. The code and the



Figure 12 ( info on the partitions which are updated with cron jobs )

```
p = string.find(output, '/export/data4')
usage = string.strip(output[p-24:p])
fp.write(""" fncdfsrv5 %20s /export/data4 """ % usage)
fp.write("\n")
percentage = string.strip(output[p-5:p-2])
if string.atoi(percentage) > 90 : IsFull=1
```

Figure 13 ( an excerpt of code that helps generate the partition information in Figure 12 )

The information about the servers are continuously updated using the cron jobs. The information's previous update was made Aug 5 as can be seen at the bottom of figure 12. All of these systems are interconnected resulting in one large cohesive farm unit. Enabling users to access and monitor the farms in an effective and efficient manner.

**Conclusion**

The monitoring and concatenation systems at the CDF Farms are an essential part to its function. The monitoring system obtained the needed modifications. The concatenation aspect of the project received several changes.

8

This system is nearly operational. The benefits for operating in the new configuration are significantly larger. Removing the concatenation from the stager and to the workers will help improve the rate at which information is written to the tape drives. The Monte Carlo events as well as ntuples will be able to benefit from these modifications. In the future when this project is fully implemented its true benefits can be seen.

**Acknowledgments**

I would like to thank my supervisor Elena Vataga for showing an indispensable amount of patience. I cannot recall when she did not make time to address my questions. I would also like to thank Murat Pavel. Murat would always find time to squeeze me into his busy schedule. I learned many things from the both of them. Last but not least the SIST committee and the individuals affiliated, for giving me this opportunity to explore my interests.

**References**

1. High Energy Physics Information Ceter Feb 14, 2002 , http://www.hep.net/

2. Collidewr Detector at FermiLab Aug 2006 http://www-cdf.fnal.gov/about/index.html

3. SAM Home Page March 2 2006 http://d0db.fnal.gov/sam/

4. Condor Project Homepage August 2006 http://www.cs.wisc.edu/condor/

5. Concurrent Versions System August 2006 http://en.wikipedia.org/wiki/Concurrent_Versions_System

6. MySQL AB August 2006 http://www.mysql.com/

7. Tikiwiki: The Tikiwiki Community August 2006 http://tikiwiki.org/