# FRESCO: The Open Failure Data Repository

**The** open failure data repository for computer site and machine failures for large scale clusters is an effort from the **D**ependable **C**omputing **S**ystems **L**aboratory (DCSL) research group of Purdue University with the support of NSF (grant CNS-1405906).

The document is consists of following parts
1. Overview of the data set
2. Overview of the Conte cluster
    a. Hardware specifications
    b. Software specifications
3. Data collection steps
4. Description of data set – Accounting Statistics
5. Description of data set – TACC Stats

# 1. Overview of the data set

The data is collected from Purdue University's latest computing cluster- Conte. In a duration of 6 months starting from October 2014, we have collected the different types of data from the cluster to for a User-centric workload study. Towards this we have collected information pertaining to the jobs that are submitted by the user and the information from all the nodes of the cluster used for scheduling jobs. Using the data set spanning six months, we analyzed over 489,000 jobs submitted by over 300 different users. In conjunction, we also analyzed roughly 3.8K user tickets (but they are not shared due to privacy concerns).

The data set is consists of the following components
1. Accounting statistics extracted from job scheduler
2. Node level statistics extracted from monitoring tool

| Data set Summary | |
|---|---|
| Duration over which data set is collected | October 2014 to March 2015 |
| Total number of jobs analyzed | 489,971 |
| Number of unique users | 306 |
| Number of user tickets analyzed | 3,873 |

# 2. Overview of the Conte cluster

Conte is homogenous cluster with all nodes running on same specifications (hardware and software). In this section we describe the hardware and software specifications of Conte.

| Hardware Specifications | |
|---|---|
| Total number of Nodes in the cluster | 580 |
| Number of cores per node | 16 |
| The processors in node | 2 x 8-Core Intel Xeon-E5 @ 2.6 GHz |
| Accelerator Cards | 2 x 60-Core Xeon Phi |
| Memory available per node | 64 GB of DDR3, 1,600 MHz |
| Interconnect used for inter-node communication | 40 Gbps FDR 10 Infiniband |
| Software Specifications | |
| Operating system | Red Hat Enterprise Linux v6.6 |
| Job Scheduler | TORQUE 4 |
| Resource Manager | Moab 7 |
| Scratch filesystem by jobs | Lustre 2.4 |
| Performance monitoring Tool | TACC Stats |

# 3. Data Collection steps

As mentioned earlier, the data collected from Conte has multiple parts. Each part of data is contributed by different parts of the cluster system. The scheduler, TORQUE, provides the information about every job that was submitted to the system. The performance monitoring tool, TACC Stats, provides data with respect to each node unlike TORQUE. TACC stats provided regular snapshots of the system state by extracting the system statistics for various components like network, file system, virtual memory, CPU, Lustre file system and many more. In Conte, the data collected by TACC stats is at interval of 10 minutes.

# 4. Accounting Statistics

This section provides detailed description of each field present in accounting statistics file. The Sample data file (SampleAccStats_Data.tsv) provided in the repository provides a glimpse of the full data (AccStats_Oct2014ToMar2015.tar.gz) that is available in the zipped format. Many fields in the provided data are self-explanatory. In the table below we explain few important fields of the data.

| Field Name | Field description |
|---|---|
| Job ID | The ID assigned by the scheduler for the submitted job |
| user | Username of the user submitting the job |
| Job Status | Value indicates status of the job. 'E' implies Job has exited (successfully or unsuccessfully). This are also called record markers. The different values and their descriptions are provided here |
| account | Account name used for submitting the job |
| ctime | Time the job was created |
| etime | Time the job became eligible to run |
| qtime | Time the job was queued |
| start | Time the job started |
| end | Time the job ended |
| group | Group name of the user who submitted the job |
| jobname | Name of job submitted |
| owner | A hostname that is the owner of the job |
| queue | Which queue that job is submitted |
| session | Session ID for the submitted job |
| Resource_List.mem | Memory requested by the job |
| Resource_List.naccesspolicy | Type of access policy. Default or no-value indicates that no sharing is allowed on node used by a job. A value of "Shared" indicates that the node can be shared between multiple jobs. |
| Resource_List.neednodes | This field is a tuple. First field indicates Number of nodes needed. The second field is a ppn field which indicates the number of processors per node needed. |
| Resource_List.ncpus | Number of cpus needed |
| Resource_List.pmem | Memory requested by the job |
| Resource_List.walltime | Wall time requested by the job |
| resources_used.cput | Amount of cpu time used by the job |
| resources_used.pmem | Amount of peak memory used by the job |
| resources_used.vmem | Amount of virtual memory used by the job |
| resources_used.walltime | Amount of wall time consumed by the job |
| exec_host | Node and Cores which the job is scheduled to use. Node000/01 indicates that 1$^{st}$ core of Node000 is allocated to be used by the job |

## 5. TACC Stats

\*\*Due to large number of fields present in the TACC Stats, the description is provided in a separate XL sheet. It is available in the *Documentation* folder of the repository. Link \*\*

As TACC stats information is available per node, any analysis related to jobs will have to be aligned with the information from accounting statistics. As there is no explicit key used to tie up records of TACC stats and accounting records, timestamps of the jobs (start and end attributes) are good choices. Using these attributes we can find out all the statistics pertaining to a job.

In the *.tsv* files available in the TACC stats folder, the column headers are not provided. The column headers for each *.tsv* file name is provided as semicolon (;) delimited list in the file *TACCStats_AttributeList.txt*. All the headers provided here are in correct order with the values in *.tsv* files. Hence in order to manually inspect a *.tsv* file, please follow these steps.
1. Pick the *.tsv* file you wish to inspect.
2. Pick the corresponding list of column headers provided for that file from *TACCStats_AttributeList.txt*
   a. Replace all semicolon with Tabs
   b. Now it'll be a tab delimited list
3. Copy the tab delimited list on the top of *.tsv* file
4. Now, the .tsv file is can be viewed (MS-XL or any other viewer)