# Comprehensive Grid and Job Monitoring with Fifemon

• • •

Kevin Retzke
User Support for Distributed Computing @ Fermilab
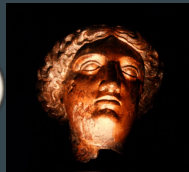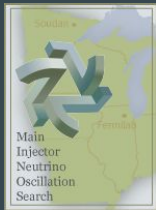HTCondor Week 2016

Landscape

Fermilab

# FIFE Project

FabrIc for Frontier Experiments:

Common computing for "not CMS" experiments at Fermilab

- O(10) experiments
- O(100) users
- O(10 000) simultaneous jobs
- O(1 000 000) jobs per week
- O(1 PB) data collected per month
- One global HTCondor pool (via GlideinWMS)
  - ~ ⅔ jobs run on dedicated local cluster
  - ~ ⅓ opportunistic through Open Science Grid

# Why Do We Need Monitoring?

Grid admins want to know:

- Overall health of the batch system
- Worker node status and availability
- Efficiency in matching jobs to resources
- Identify and fix problems quickly (before users and stakeholders notice… and open tickets)

Users want to know:

- State of their jobs
- Availability of resources
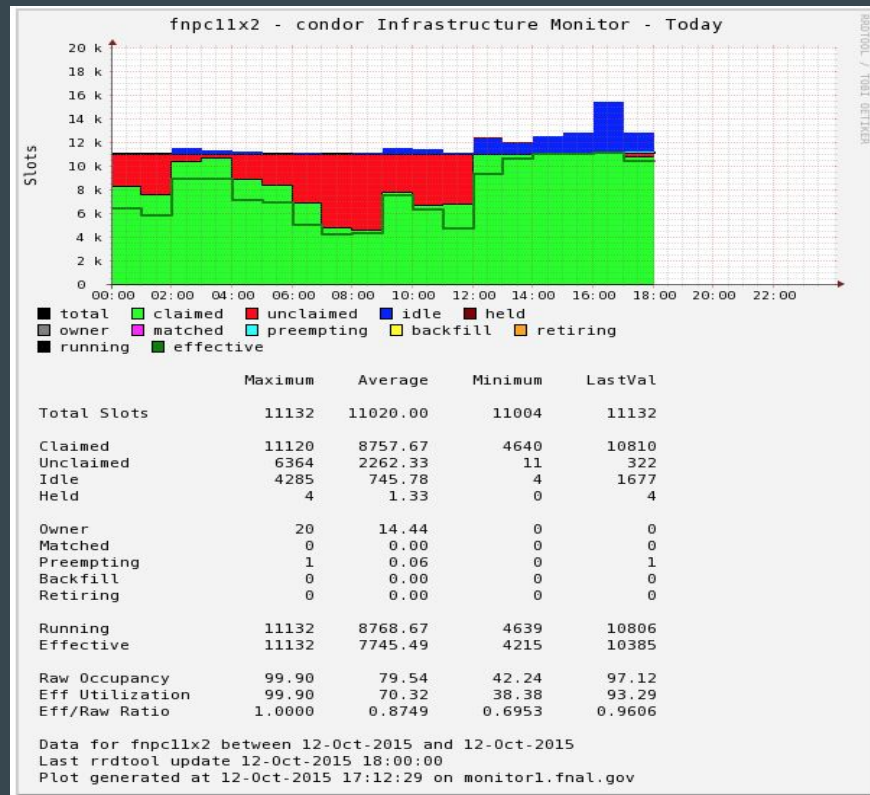- WHY ISN'T MY JOB RUNNING?

Stakeholders want to know:

- Each group is getting the resources it needs
- Resources are being used effectively

# Fermigrid Monitor (ca. 2004)

Monitoring for local HTCondor cluster (GPGrid).

- Aggregate metrics for grid and VOs.
- No offsite information, no user job information.
- Difficult to alter or expand.

OK for grid admins, good for stakeholders, bad for users.



fnpc11x2 - condor Infrastructure Monitor - Today

Slots

20 k
18 k
16 k
14 k
12 k
10 k
8 k
6 k
4 k
2 k
0

00:00  02:00  04:00  06:00  08:00  10:00  12:00  14:00  16:00  18:00  20:00  22:00

RRDTOOL / TOBI OETIKER

■ total        ■ claimed      ■ unclaimed    ■ idle        ■ held
■ owner        ■ matched      ■ preempting   ■ backfill    ■ retiring
■ running      ■ effective

|  | Maximum | Average | Minimum | LastVal |
|---|---|---|---|---|
| Total Slots | 11132 | 11020.00 | 11004 | 11132 |
| Claimed | 11120 | 8757.67 | 4640 | 10810 |
| Unclaimed | 6364 | 2262.33 | 11 | 322 |
| Idle | 4285 | 745.78 | 4 | 1677 |
| Held | 4 | 1.33 | 0 | 4 |
| Owner | 20 | 14.44 | 0 | 0 |
| Matched | 0 | 0.00 | 0 | 0 |
| Preempting | 1 | 0.06 | 0 | 1 |
| Backfill | 0 | 0.00 | 0 | 0 |
| Retiring | 0 | 0.00 | 0 | 0 |
| Running | 11132 | 8768.67 | 4639 | 10806 |
| Effective | 11132 | 7745.49 | 4215 | 10385 |
| Raw Occupancy | 99.90 | 79.54 | 42.24 | 97.12 |
| Eff Utilization | 99.90 | 70.32 | 38.38 | 93.29 |
| Eff/Raw Ratio | 1.0000 | 0.8749 | 0.6953 | 0.9606 |

Data for fnpc11x2 between 12-Oct-2015 and 12-Oct-2015
Last rrdtool update 12-Oct-2015 18:00:00
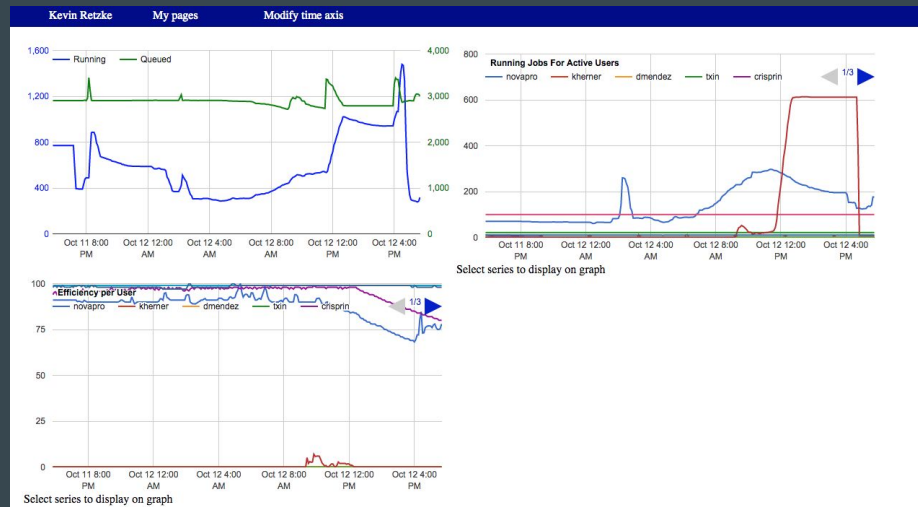Plot generated at 12-Oct-2015 17:12:29 on monitor1.fnal.gov

# Fifemon v1 (ca. 2014)

Growing usage of offsite resources through OSG; needed new monitoring.

- Aggregate metrics for users and VOs.
- No cluster information.
- Cumbersome to maintain and expand.

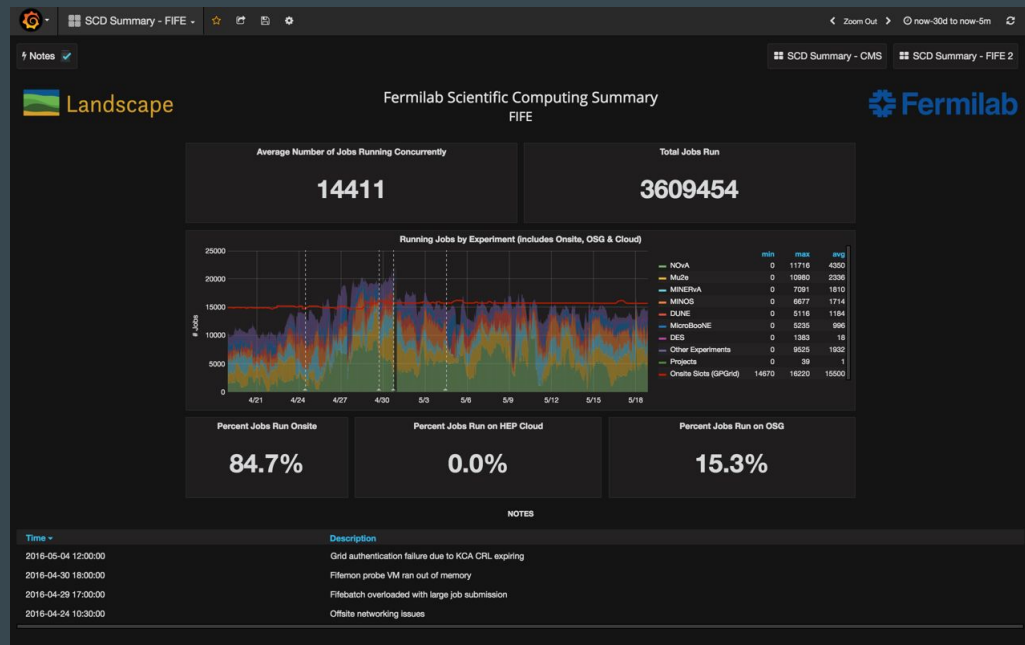OK for grid admins, bad for stakeholders, good for users.

# Fifemon v2+ (ca. 2015)

Landscape Program: develop comprehensive monitoring for FIFE, HEP Cloud, and beyond.

- Leverage open-source monitoring technology
- Focus on incorporating new data sources and new dashboards
- Rapid development and iteration of tailored views for each target audience.

Good for grid admins, stakeholders, and users alike.

# Fifemon Backend

Data collection:

- Generic HTCondor probe collecting daemon, machine, and job status
- Logstash collecting live HTCondor Events
- Several other centrally-run probes querying other resources
- Some services directly reporting to Graphite

- Most probes report stats every five minutes
- Graphite:
  - 250K individual metrics
  - ~80GB
  - 10 year history
- Elasticsearch: ~8GB per day

Graphite:

- Time-series database, stores data in files similar to RRD with caching layer.
- Simple line protocol
- Powerful query manipulations and aggregations

Elasticsearch:

- "NoSQL" document database, powered by Apache Lucene.
- Store full details on current jobs, batch slots, and logs.

# Fifemon Frontend

Grafana:

- Time-series (primarily) visualization dashboard platform.
- Supports numerous data sources (Graphite, InfluxDB, Elasticsearch, etc).
- Several auth methods (LDAP, OAuth, proxy).
- Rich user interface for graphing metrics and composing dashboards.
- Scripted and templated dashboards and raw HTML panels allow extensive customization.
- V3 (released last week) introduces new plugin system to support custom datasources and panels.

Kibana:

- Elasticsearch data only
  - Current jobs and machine status
  - Event logs
- Explore data, create ad-hoc visualizations, combine into dashboards
- Used for analytics and troubleshooting
- Access limited to grid admins and power users

# Fifemon Architecture

# Next Steps

Fifemon is constantly evolving:

- Adding new data sources and metrics
- New dashboards:
    - Tailored views based on user request
    - Discovering new ways of looking at the data
- New Grafana panels
- Further leverage HTCondor event logs & gangliad/metricsd for true real-time monitoring

# Case Studies

# "There's a dashboard for that..."

# Case Study: Grid Admin

"Is the batch system healthy?"



Cinnamon

Woof, all green!

Grid: gpgrid

FIFE Onsite Summary | GPGrid | GPGrid Group | Why Are There Unused Slots on GPGrid?

PAGE HELP

RELATIVE UTILIZATION

## CPU Utilization

| | min | max | avg | current |
|---|---|---|---|---|
| CPU Claimed | 93.488% | 93.622% | 93.549% | 93.580% |
| CPU Utilized | 66.928% | 74.206% | 71.271% | 74.206% |

## Memory Utilization

| | min | max | avg | current |
|---|---|---|---|---|
| Memory Claimed | 80.853% | 81.224% | 81.063% | 81.180% |
| Memory Utilized | 53.703% | 55.235% | 54.930% | 54.102% |

## Disk Utilization

| | min | max | avg | current |
|---|---|---|---|---|
| Disk Claimed | | 5.301% | 5.438% | 5.386% | 5.415% |
| Disk Utilized | | 3.938% | 4.053% | 3.990% | 3.954% |

ABSOLUTE UTILIZATION

## CPU

| | min | max | avg | current |
|---|---|---|---|---|
| Claimed | | 14521 | 14606 | |
| Unclaimed | | 1001 | 1002 | |
| Unusable | | 0 | 0 | 0 |
| Total | | 15523 | 15608 | |
| | | 11063 | 11582 | |

Grid utilization is OK.

## Memory

| | min | max | avg | current |
|---|---|---|---|---|
| Claimed | 27.6 TiB | 28.0 TiB | 27.8 TiB | 27.9 TiB |
| Unclaimed | 4.1 TiB | 4.2 TiB | 4.1 TiB | 4.1 TiB |
| Unusable | 2.3 TiB | 2.4 TiB | 2.3 TiB | 2.4 TiB |
| Total | 34.1 TiB | 34.4 TiB | 34.2 TiB | 34.4 TiB |
| Effective | 18.4 TiB | 18.9 TiB | 18.8 TiB | 18.6 TiB |

## Disk

| | min | max | avg | current |
|---|---|---|---|---|
| Claimed | 63 TiB | 65 TiB | 64 TiB | 65 TiB |
| Unclaimed | 230 TiB | 232 TiB | 231 TiB | 230 TiB |
| Unusable | 887 TiB | 901 TiB | 892 TiB | 900 TiB |
| Total | 1.184 PiB | 1.193 PiB | 1.187 PiB | 1.193 PiB |
| Effective | 47 TiB | 48 TiB | 47 TiB | 47 TiB |

GROUP UTILIZATION

Probe Status - Dev

Zoom Out    Last 6 hours

probe: gpce01_status + gpce02_status

Update Time

| Metric ▲ | Min | Max | Avg | Current |
|---|---|---|---|---|
| awsmonitor | - | - | - | - |
| cmssrv14_status | 1.61 s | 8.98 s | 2.05 s | 1.84 s |
| cmssrv274_status | 0.32 s | 1.02 s | 0.39 s | 0.38 s |
| cmssrv39_status | 0.86 s | 2.34 s | 1.40 s | 1.37 s |
| condor_pool_jobs | - | - | - | - |
| fifebatch-pp_status | 1.18 s | 11.11 s | 1.71 s | 1.26 s |
| fifebatch2_status | 3.90 min | 5.89 min | 4.77 min | 3.93 min |
| fifebatch_status | 4.07 min | 5.72 min | 4.73 min | 4.28 min |
| fnpccm1_status | - | - | - | - |
| gpce01_status | 2.38 s | 11.15 s | 3.01 s | 2.57 s |
| gpce02_status | 3.24 s | 9.05 s | 3.79 s | 3.34 s |
| gpcollector01_status | 1.99 s | 2.04 min | 25.28 s | 2.42 s |
| gpgrid | - | - | - | - |

Hmm, we couldn't query a CE for a few minutes. I'll check the probe logs.

gpce02_status

# Case Study: Stakeholder

"Is my experiment getting the resources it needs and using them effectively?"



Hazel

**Fermilab**

# FIFE Batch Monitoring

**Landscape**

## QUICK LINKS

Help  |  About Fifemon  |  FIFE Summary  |  CMS Summary

### Experiments

NOvA  MINERvA  MINOS  DUNE

MicroBooNE  DES  Other

### For Users

User Batch Details  |  Why Isn't My Job Running?

### Grid Status

FIFE Onsite Summary  |  Fifebatch  |  GPGrid (CE)

GPGrid (Condor)

## DASHBOARDS

### Main Dashboards

About Fifemon ⭐

Experiment Overview ⭐

Fifebatch ⭐

GPGrid ⭐

Grid Utilization ⭐

Help ⭐

Jobs Exceeding Resource Request ⭐

SCD Summary - CMS ⭐

### Starred dashboards

Fifebatch Health ⭐

Fifebatch Slots ⭐

Job Cluster Summary ⭐

Probe Status ⭐

Experiment Overview

nova

Experiment Batch Details   Experiment Efficiency Details   FTS   SAM by experiment

now-6h to now-5m

BATCH

**Job Status**
15 K
10 K
5 K
0
12:00  13:00  14:00  15:00  16:00  17:00
— Running  Current: 5.18 K  — Idle  Current: 7.55 K  — Held  Current: 43
— Allocation  Current: 2.00 K

**Job Efficiency**
100%
75%
50%
25%
0%
12:00  13:00  14:00  15:00  16:00  17:00
— Overall Efficiency

**Running Jobs by User**
6 K
4 K
2 K
0
12:00  13:00  14:00  15:00  16:00  17:00
— brebel — bzamoran — crisprin — edniner — kuldeepm — lcremone
— novapro — ranjan — sedayath — siva1987 — ynitin — Allocation

SAM

**Size of active files catalogued**
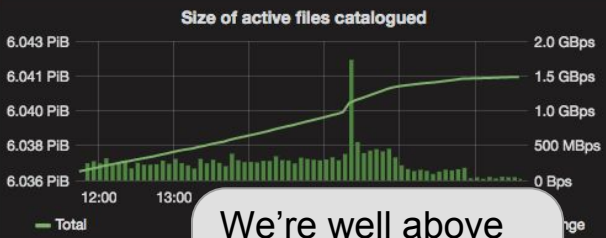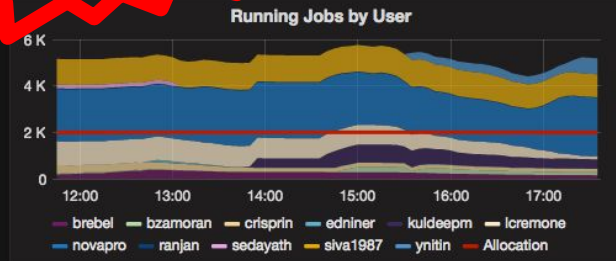6.043 PiB                    2.0 GBps
6.041 PiB                    1.5 GBps
6.040 PiB                    1.0 GBps
6.038 PiB                    500 MBps
6.036 PiB                    0 Bps
12:00  13:00
— Total

**Number of files catalogued**
48.54 Mil                    12.5 Hz
48.53 Mil                    10.0 Hz
48.52 Mil                    7.5 Hz
48.51 Mil                    5.0 Hz
                             2.5 Hz
48.50 Mil                    0 Hz
12:00  13:00  14:00  15:00  16:00  17:00
— Total       — Rate of change

**Active SAM Processes by User**
5.0 K
4.0 K
3.0 K
2.0 K
1.0 K
0
12:00  13:00  14:00  15:00  16:00  17:00
— amoren — arrieta1 — brebel — crisprin — novapro — pbultrag
— projas — radovic — sedayath — siva1987 — vito — ynitin

DCACHE

We're well above our quota, but efficiency could be better.

**Analysis Pool**
2.5 K
14:00  16:00
Queue — regular — restoreQueue
— WAN — XRootD — default

**Analysis Pool Queues**
1.5 K
1.0 K
500
0
12:00  14:00  16:00
— NFS — WAN — XRootD — default
— moverQueue — regular — restoreQueue

**Public Scratch Pool Active Transfers**
5.0 K
4.0 K
3.0 K
2.0 K
1.0 K
0
12:00  14:00  16:00
— NFS — WAN — XRootD — default
— moverQueue — regular — restoreQueue

**Public Scratch Pool Queues**
4.0 K
3.0 K
2.0 K
1.0 K
0
12:00  14:00  16:00
— NFS — WAN — XRootD — default
— moverQueue — regular — restoreQueue

Zoom Out  🕐 now-12h to now-5m

nova ▾

GPGrid Usage | Experiment Efficiency Details | Experiment Overview | FTS | SAM by experiment

## User Jobs

| User | I | R | C | X | H | Max Memory/Request | Max Disk/Request | Max Time/Request |
|------|---|---|---|---|---|--------------------|------------------|------------------|
| anorman | 0 | 0 | 0 | 0 | 9 | 0.78 | 0.00 | 0.00 |
| arrieta1 | 100 | 0 | 0 | 0 | 3 | 0.00 | 0.00 | 0.00 |
| bianjm | 825 | 2506 | 0 | 0 | 0 | 0.37 | 0.00 | 0.73 |
| boyd | 50 | 0 | 0 | 0 | 0 | 0.00 | 0.16 | 0.00 |
| brebel | 0 | 1 | 0 | 0 | 0 | 0.00 | 0.00 | 3.27 |
| crisprin | 0 | 3 | 0 | 0 | 0 | 0.01 | 0.00 | 8.55 |
| dmendez | 0 | 0 | 0 | 0 | 6 | 1.00 | 0.01 | 0.00 |
| kherner | 4 | 0 | 0 | 0 | 0 | 0.00 | 0.00 | 0.00 |
| kretzke | 1 | 0 | 0 | 0 | 0 | 0.00 | 0.00 | 0.00 |
| kuldeepm | 0 | 10 | 0 | 0 | 0 | 0.34 | 0.00 | 6.07 |
| lcremone | 0 | 2 | 0 | 0 | 0 | 0.29 | 0.13 | 5.54 |
| novapro | 22154 | 3464 | 0 | 0 | 14 | 1.05 | 1.01 | 12.04 |
| pavan219 | | | 0 | 0 | 11 | 0.95 | 0.11 | 0.00 |
| siva1987 | | | 0 | 0 | 0 | 0.66 | 0.00 | 3.39 |

Disk and Memory requests look good, lots of users exceeding request time though.

Memory Usage

14 TiB
11 TiB
9 TiB
7 TiB
5 TiB

Disk Usage

182 TiB
136 TiB
91 TiB

# Case Study: User

"What's the status of my jobs?"



Cocoa

# FIFE Batch Monitoring

Fermilab       Landscape

## QUICK LINKS

Help    About Fifemon    FIFE Summary    CMS Summary

### Experiments

NOvA   Mu2e   MINERvA   MINOS   DUNE

MicroBooNE   DES   Other

### For Users

User Batch Details    Why Isn't My Job Running?
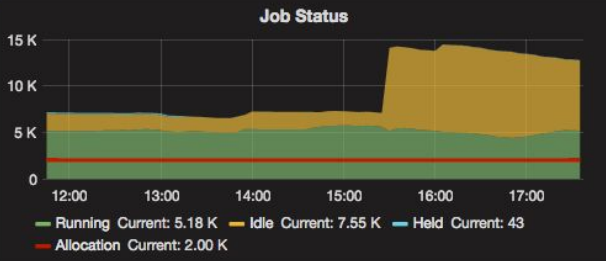
### Grid Status

FIFE Onsite Summary   Fifebatch   GPGrid (CE)

GPGrid (Condor)

## DASHBOARDS

### Main Dashboards

| | |
|---|---|
| About Fifemon | ☆ |
| Experiment Overview | ☆ |
| Fifebatch | ☆ |
| GPGrid | ☆ |
| Grid Utilization | ☆ |
| Help | ☆ |
| Jobs Exceeding Resource Request | ☆ |
| SCD Summary - CMS | ☆ |
| SCD Summary - FIFE | ☆ |

### Starred dashboards

| | |
|---|---|
| Fifebatch Health | ⭐ |
| Fifebatch Slots | ⭐ |
| Job Cluster Summary | ⭐ |
| Probe Status | ⭐ |

Top timeline axis: 05:00 06:00 07:00 08:00 09:00 10:00

Legend (left): — Fermigrid  — Fermigridosg1  — FNAL  — GPGrid          — Idle

Top-right timeline axis: 05:00 06:00 07:00 08:00 09:00 10:00

Legend (right):
— BNL  — Caltech  — Clemson  — Cornell  — FNAL_HEPCLOUD  — FZU  — Hyak_CE  — MIT  — MWT2  — Michigan
— Nebraska  — NotreDame  — OSC  — Omaha  — SMU  — SMU_HPC  — SU-OG  — TTU  — UCSD  — UChicago
— USCMS-FNAL-WC1  — Wisconsin  — unknown          — Idle

## Current Jobs

| Cluster | I | R | H | Submit Time/Command | Memory (MB) | Disk (MB) | Time (hr) | Max Eff. | Starts |
|---|---|---|---|---|---|---|---|---|---|
| 7989120 | 7 | 0 | 0 | 2016-03-08T02:22:51.000Z  qhuang-qhuang-reco-keepup-Offsite-3000-S16-03-04-neardet-BNB-25_days_ago-20160308_0222.sh_20160308_022251_3281177_0_1_wrap.sh | 0 / 3000 | 0 / 10240 | 0 / 11 | ---- | 0 |
| 7989126 | 7 | 0 | 0 | 2016-03-08T02:23:06.000Z  qhuang-qhuang-reco-keepup-Offsite-3000-S16-03-04-neardet-BNB-27_days_ago-20160308_0222.sh_20160308_022305_3282245_0_1_wrap.sh | 0 / 3000 | 0 / 10240 | 0 / 11 | ---- | 0 |
| 7989131 | 7 | 0 | 0 | 2016-03-08T02:23:18.000Z  qhuang-qhuang-reco-keepup-Offsite-3000-S16-03-04-neardet-BNB-29_days_ago-20160308_0223.sh_20160308_022318_3283095_0_1_wrap.sh | 0 / 3000 | 0 / 10240 | 0 / 11 | ---- | 0 |
| 7989137 | 7 | 0 | 0 | 2016-03-08T02:23:30.000Z  qhuang-qhuang-reco-keepup-Offsite-3000-S16-03-04-neardet-BNB-31_days_ago-20160308_0223.sh_20160308_022329_3284093_0_1_wrap.sh | 0 / 3000 | 0 / 10240 | 0 / 11 | ---- | 0 |
| 7991809 | 0 | 244 | 0 | 2016-03-08T04:58:24.000Z  bzamoran-prod_full_chain_R16-03-03-prod2reco.a_ND_numi_epoch3c-20160308_0458.sh_20160308_045824_3734826_0_1_wrap.sh | 1952 / 2000 | 451 / 34180 | 6 / 6 | 62.3% | 2 |
| 7993210 | 7998 | 1719 | 0 | 2016-03-08T08:50:35.000Z  tghosh-tghosh_prod_daq_R16-02-11-prod2genie.b_fd_genie_nonswap_fhc_nova_v08_full_batch1_v1_birksmodB-20160308_0850.sh_20160308_085035_4016786_0_1_wrap.sh | 1353 / 2000 | 2178 / 4000 | 2 / 3 | 36.8% | 1 |
| 4937278 | | | | 2016-03-08T09:21:32.000Z  ...chain_R16-03-03-prod2reco.a_ND_numi_period1-20160308_0921.sh_20160308_092132_3138891_0_1_wrap.sh | 1925 / 2000 | 114 / 34180 | 1 / 6 | 57.5% | 1 |
| 4937313 | | | | ...4:26.000Z  ...chain_R16-03-03-prod2reco.a_ND_numi_epoch3b-20160308_0924.sh_20160308_092426_3149550_0_1_wrap.sh | 1929 / 2000 | 79 / 34180 | 1 / 6 | 55.2% | 1 |
| | | | | ...7:01.000Z  bzamoran-prod_full_chain_R16-03-03-prod2reco.a_ND_numi_period2-20160308_0936.sh_20160308_093701_3191659_0_1_wrap.sh | 1920 / 2000 | 85 / 34180 | 1 / 6 | 53.2% | 1 |

Speech bubble: This cluster has poor efficiency, let's take a look at it.

**COMPLETED JOBS**

**RESOUCE GRAPHS**

Zoom Out  ⏱ Last 24 hours  Refresh every 5m

cluster: 7714932 ▾

◄ PAGE HELP

**JOB INFORMATION**

| | | | |
|---|---|---|---|
| **Job ID:** | 7714932.0@fifebatch2.fnal.gov | **Resources Requested** | |
| **Submit Date:** | 2016-02-26T18:09:46 | **CPU:** | 1 |
| **Experiment:** | mu2e | **Memory:** | 3994 MB |
| **User:** | mu2epro (mu2epro/cron/mu2egpvm01.fnal.gov@FNAL.GOV) | **Disk:** | 9216 MB |
| **Usage Model:** | OFFSITE | **Runtime:** | 9 hr |
| **Sites Requested:** | BNL,Caltech,FERMIGRID,FNAL,MIT,Michigan,Nebraska,Omaha,SU-OG,Wisconsin,UCSD,NotreDame,MWT2 | | |

View sandbox files          View available slots

**PROCESS STATUS**

| Total Processes | Idle Processes | Running Processes | Held Processes |
|---|---|---|---|
| 9175 | 6065 | 2898 | 4 |

A few failed processes, and a bunch are disconnected.

| | Failed Processes (nonzero exit code) | Disconnected Processes |
|---|---|---|
| | 26 | 408 |

**RESOURCES USED**

| Max Memory Usage | Max Disk Usage | Max Walltime |
|---|---|---|

| Competed Processes (exit code 0) | Failed Processes (nonzero exit code) | Disconnected Processes |
|:---:|:---:|:---:|
| **1011** | **26** | **408** |

## RESOURCES USED

| Max Memory Usage | Max Disk Usage | Max Walltime |
|:---:|:---:|:---:|
| **1.934 GiB** | **7.91 GiB** | **11.11 hour** |

| Memory Usage | | | Disk Usage | | | Walltime | | |
|---|---|---|---|---|---|---|---|---|
| Min ▾ | Max | Average | Min ▾ | Max | Average | Min ▾ | Max | Average |
| 10.02 MiB | 1.93 GiB | 1.31 GiB | 1.75 GiB | 7.91 GiB | 5.34 GiB | 33.70 min | 11.11 hour | 5.36 hour |

◀ PROCESS LIST

## CONDOR EVENTS

### All Events



Legend: ecuteEvent, JobDisconnectedEvent, JobHeldEvent, JobImageSizeEvent, JobReconnectFailedEvent, ReconnectedEvent, JobReleaseEvent, JobTerminatedEvent, ShadowExceptionEvent

### Abnormal Events



Legend: JobDisconnectedEvent, JobHeldEvent, JobReconnectFailedEvent, JobReconnectedEvent, JobReleaseEvent, ShadowExceptionEvent

STATS BY SITE

JOBSUB

# Case Study: Upper Management

*"What does the computing division do again?"*

Sage (and minions)

⚡ Notes ✔                                    ⚏ SCD Summary - CMS    ⚏ SCD Summary - FIFE 2

🟩 **Landscape**          **Fermilab Scientific Computing Summary**          ✦ **Fermilab**
                                         **FIFE**
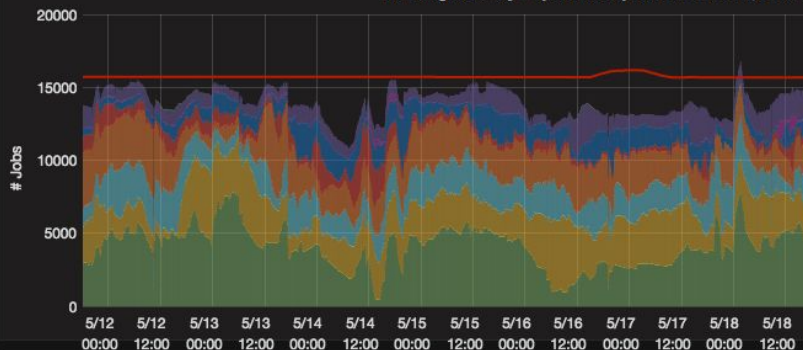
**Average Number of Jobs Running Concurrently**                **Total Jobs Run**

# 13899                                          # 818008

**Running Jobs by Experiment (includes Onsite, OSG & Cloud)**

| | min | max | avg |
|---|---|---|---|
| NOvA | 432 | 7863 | 4091 |
| Mu2e | 52 | 5249 | 2643 |
| MINERvA | 602 | 3748 | 1931 |
| MINOS | 285 | 5464 | 2281 |
| DUNE | 10 | 2804 | 715 |
| MicroBooNE | 117 | 3761 | 1121 |
| DES | 0 | 1383 | 48 |
| Other Experiments | 265 | 2891 | 1067 |
| Projects | 0 | 39 | 2 |
| Onsite Slots (GPGrid) | 15664 | 16175 | 15735 |

Y-axis: # Jobs (0, 5000, 10000, 15000, 20000)
X-axis: 5/12 00:00, 5/12 12:00, 5/13 00:00, 5/13 12:00, 5/14 00:00, 5/14 12:00, 5/15 00:00, 5/15 12:00, 5/16 00:00, 5/16 12:00, 5/17 00:00, 5/17 12:00, 5/18 00:00, 5/18 12:00

| **Percent Jobs Run Onsite** | **Percent Jobs Run on HEP Cloud** | **Percent Jobs Run on OSG** |
|---|---|---|
| # 92.4% | # 0.0% | # 7.6% |

**NOTES**

⚡ Notes ✔️     ▦ SCD Summary - CMS    ▦ SCD Summary - FIFE

Landscape

# Fermilab Scientific Computing Summary
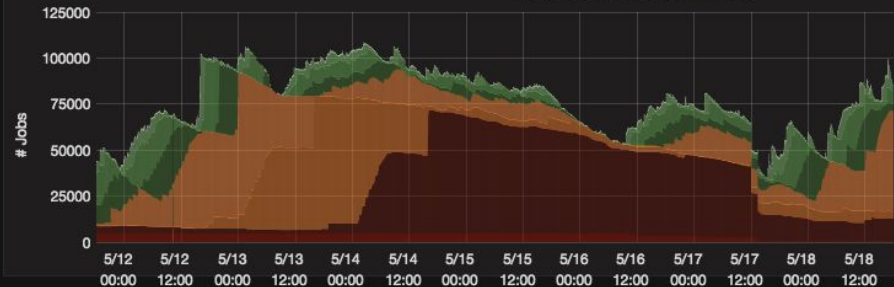## FIFE (continued)

❖ Fermilab

### Average Number of Jobs Waiting in Queue

# 77565

### Average Time Spent Waiting in Queue

# 13.06 hour

### Time Spent Waiting in Queue

| | min | max | avg |
|---|---|---|---|
| > 7 days | 14 | 4849 | 3454 |
| 2-7 days | 219 | 67478 | 27966 |
| 24-48 hours | 13 | 72728 | 12665 |
| 8-24 hours | 40 | 76989 | 16414 |
| 4-8 hours | 0 | 37261 | 6939 |
| 1-4 hours | 0 | 39132 | 6631 |
| < 1 hour | 0 | 39315 | 3091 |
| new | 1 | 12008 | 417 |

### New Data Cataloged

# 314.9 TB

**NOTES**

Zoom Out   Last 7 days

⚡ Notes ☑

Help

Select Experiment: ANNIE  CDF  CDMS  D0  DUNE  LArIAT  MINERvA  MINOS  MicroBooNE  Mu2e  NOvA  SBND  SeaQuest  g-2

Landscape

# NOvA Computing Summary
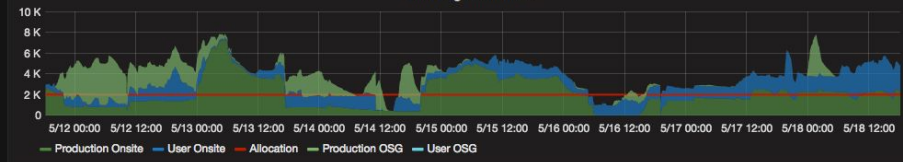
Fermilab

### Average Jobs Running Concurrently
# 4093

### Total Jobs Run
# 152018

### Average Time Spent Waiting in Queue (Production)
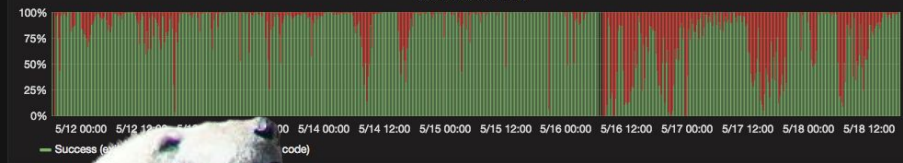# 3.440 hour

**Running Batch Jobs**



— Production Onsite  — User Onsite  — Allocation  — Production OSG  — User OSG

**Queued Production Jobs by Wait Time**



— > 7 days  — 2-7 days  — 24-48 hours  — 8-24 hours  — 4-8 hours  — 1-4 hours  — < 1 hour  — new

### Total Jobs Failed (nonzero exit code)
# 15132

### Average CPU Efficiency
# 43.6%

**Job Success Rate**



— Success (code)

**CPU Efficiency**



— Overall  — Onsite

### New Data Cataloged
# 116.4 TB

### Total Data Cataloged
# 6.8 PB

Comprehensive grid monitoring with Fifemon has improved resource utilization, job throughput, and computing visibility at Fermilab.

Probes, dashboards, and docs at:
https://github.com/fifemon