# Field of Science Reporting in the OSG

## Overview

Most of the reporting in the OSG groups its data in one of two ways – by organization (VO) that submitted the jobs or by the site that ran the jobs. This is a natural approach, because it groups jobs either by producer (the VO) or the consumer (the site).

We are working on a third method of groupings: field of science. Grouping data by science has immediate strengths:

1. External parties don't recognize VOs by their names, especially smaller VOs. Classifying by science allows us to cast the data in more familiar terminology.

2. Knowing the science provides the OSG with feedback on how far they are spreading outside of the field of HEP, which has been the primary driver this far.

This document covers the current progress and results in reporting accounting data by field of science.
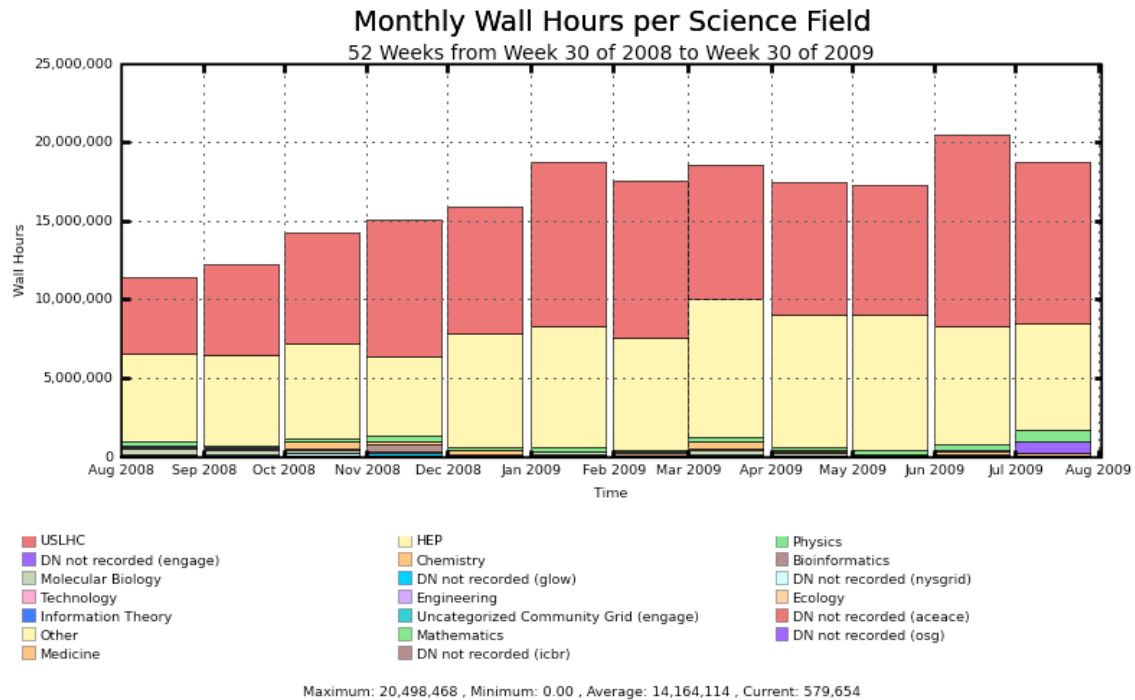
While the OSG has a thorough accounting system with Gratia, several critical pieces of data had to be developed to complete this project. The largest hurdle was handling community grid VOs. These organizations have many users with multiple science domains; unlike experiment VOs, it was not as simple as labeling the entire VO as being in one domain.

## Results

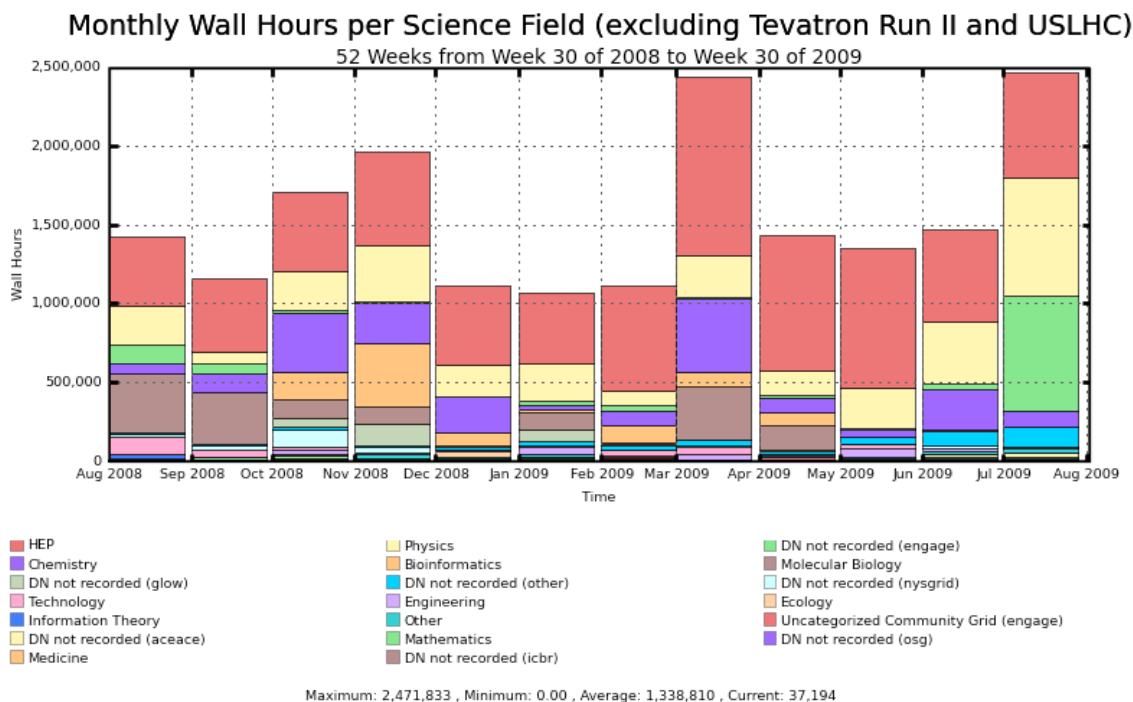The graph below shows the number of wall hours spent on jobs by field of science[1].

---

[1] http://t2.unl.edu/gratia/bar_graphs/monthly_field?starttime=2008-08-01%2000:00:00&endtime=2009-08-01%2023:59:59

**Monthly Wall Hours per Science Field**
52 Weeks from Week 30 of 2008 to Week 30 of 2009

Legend:
- USLHC
- HEP
- Physics
- DN not recorded (engage)
- Chemistry
- Bioinformatics
- Molecular Biology
- DN not recorded (glow)
- DN not recorded (nysgrid)
- Technology
- Engineering
- Ecology
- Information Theory
- Uncategorized Community Grid (engage)
- DN not recorded (aceace)
- Other
- Mathematics
- DN not recorded (osg)
- Medicine
- DN not recorded (icbr)

Maximum: 20,498,468 , Minimum: 0.00 , Average: 14,164,114 , Current: 579,654

Because USLHC (CMS and ATLAS) and Tevatron Run II (D0 and CDF) experiments are so dominant, it helps to remove those four VOs[2]:
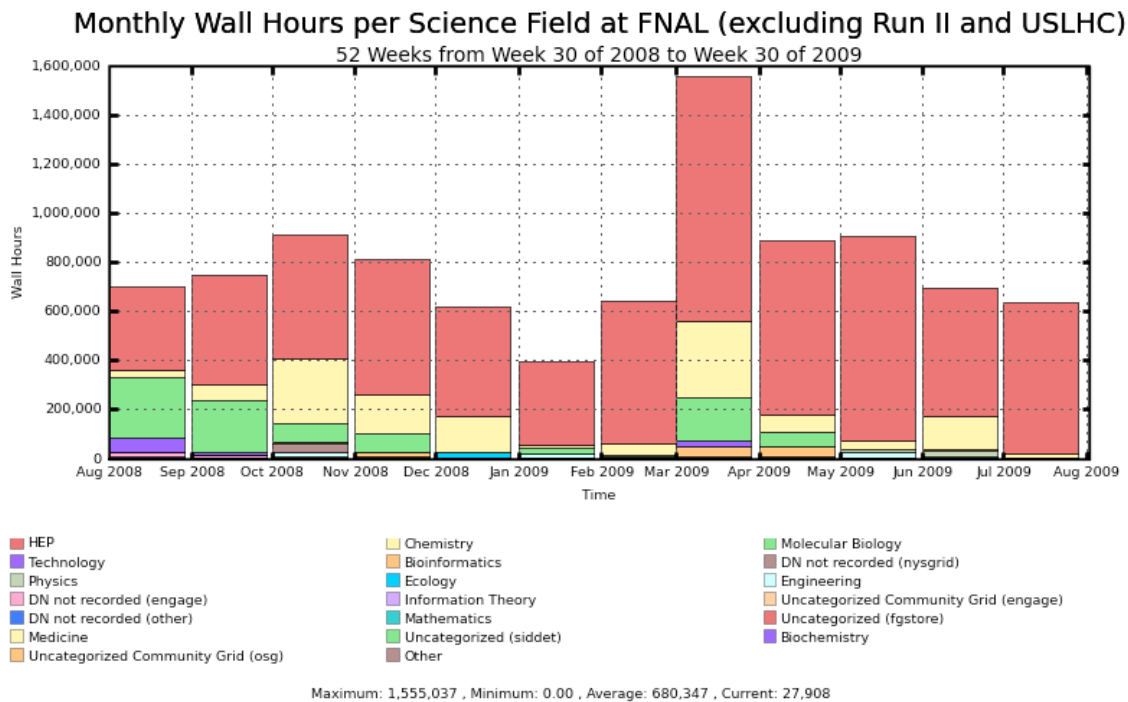
---

[2] http://t2.unl.edu/gratia/bar_graphs/monthly_field?starttime=2008-08-01%2000:00:00&endtime=2009-08-01%2000:00:00&exclude-vo=cms|atlas|cdf|dzero|unknown

**Monthly Wall Hours per Science Field (excluding Tevatron Run II and USLHC)**
52 Weeks from Week 30 of 2008 to Week 30 of 2009

Legend:
- HEP
- Chemistry
- DN not recorded (glow)
- Technology
- Information Theory
- DN not recorded (aceace)
- Medicine
- Physics
- Bioinformatics
- DN not recorded (other)
- Engineering
- Other
- Mathematics
- DN not recorded (icbr)
- DN not recorded (engage)
- Molecular Biology
- DN not recorded (nysgrid)
- Ecology
- Uncategorized Community Grid (engage)
- DN not recorded (osg)

Maximum: 2,471,833 , Minimum: 0.00 , Average: 1,338,810 , Current: 37,194

Several entries in the legend are worth explaining, as they point out the limitations of our data collection method:

- **Other**: Any entry whose usage is too small to display on the graph gets aggregated into this catch-all entry.

- **DN not recorded (VO name)**: Not all sites report the certificate's Distinguished Name (DN) associated with a job, due to misconfiguration, complex site issues, or historical data that was reported prior to Gratia having this capability.  This entry corresponds to jobs without reported DNs whose VO does not have a single science classification.  These entries are due to missing information in the job record, and it is not possible to retroactively correct past data; any site has the ability to report this correctly.

- **Uncategorized Community Grid (VO name)**: This is usage in a community grid which has a DN associated it, but the VO has not yet determined what that user's science community is.  When the VO has made that determination, we will be able to label this historical data.

- **Uncategorized (VO name)**: This is usage in a VO that has not been classified at all.  Usually, this signals a VO that runs at a site but has not registered with the OSG.

One requested use of this data is to classify the science being performed at a single site (i.e., so the executives know what areas they are contributing to).  Below, we have an example of the science being done at the FNAL facility[3].



Monthly Wall Hours per Science Field at FNAL (excluding Run II and USLHC)
52 Weeks from Week 30 of 2008 to Week 30 of 2009

Maximum: 1,555,037 , Minimum: 0.00 , Average: 680,347 , Current: 27,908

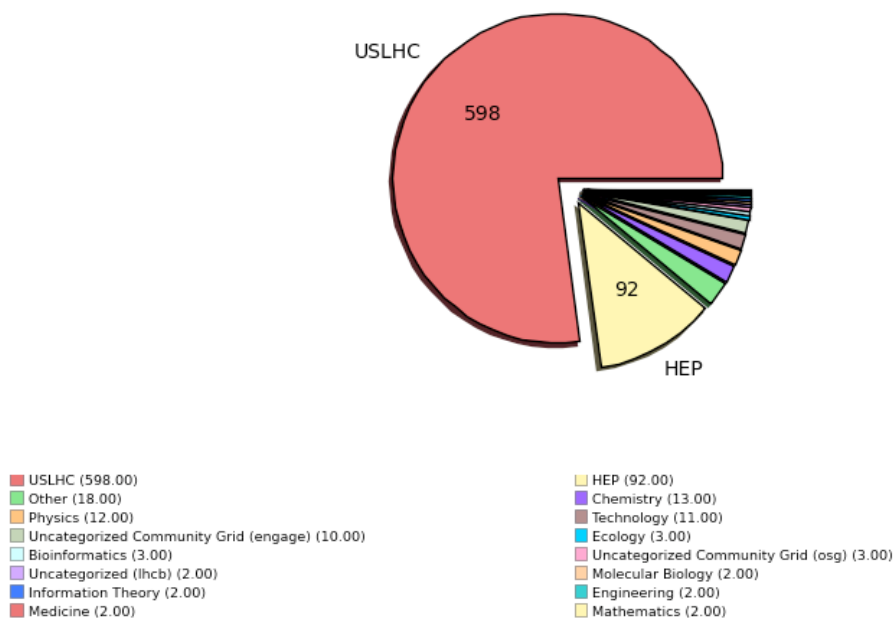Below is the number of "significant" users recorded per science field[4].  Here, we define a user to be "significant" if they've used more than 1000 wall hours.

---

[3] http://t2.unl.edu/gratia/bar_graphs/monthly_field?starttime=2008-08-01%2000:00:00&endtime=2009-08-01%2000:00:00&probe=fnal.gov&exclude-vo=cdf|dzero|atlas|cms|unknown

[4] http://t2.unl.edu/gratia/pie_graphs/science_field_user_count

## Number of Users by Science Field (Sum: 775.00 People)
### 52 Weeks from Week 31 of 2008 to Week 31 of 2009



| | |
|---|---|
| ■ USLHC (598.00) | □ HEP (92.00) |
| ■ Other (18.00) | ■ Chemistry (13.00) |
| ■ Physics (12.00) | ■ Technology (11.00) |
| ■ Uncategorized Community Grid (engage) (10.00) | ■ Ecology (3.00) |
| ■ Bioinformatics (3.00) | ■ Uncategorized Community Grid (osg) (3.00) |
| ■ Uncategorized (lhcb) (2.00) | ■ Molecular Biology (2.00) |
| ■ Information Theory (2.00) | ■ Engineering (2.00) |
| ■ Medicine (2.00) | □ Mathematics (2.00) |

## Method of Data Collection

All OSG VOs are classified according to their science field or "community grid" if their users can be classified under multiple fields. For these community grids, we have asked that they self-classify their users according to science field. We have developed a simple CSV format (so it can be edited in Microsoft Excel) that allows the community grids to classify a certificates Distinguished Name (DN), corresponding to a single user, or a FQAN, corresponding to a group of users. They then host this CSV file on their own servers and update it as users come or go[5]. Ideally, this would allow Metrics and Measurements to not have to maintain any data centrally. For this report, we have used a separate CSV to classify users for VOs that are not yet up to speed.

Each job record has an associated VO and may also have a DN and FQAN. The information we have received from VOs and OIM is applied to each record; it is applied first to the DN, then FQAN, then VO. If information is missing, we may apply one of three labels: DN not recorded, Uncategorized, or Uncategorized Community Grid. These labels are explained above.

The two largest (in terms of wall hours) community grids Engage and NYSGrid are currently participating; we expect the OSG VO to participate in the near future.

We provide feedback to the VOs to help them know which users have not yet been classified, sorted by wall hours. The pie graphs below show this information for the
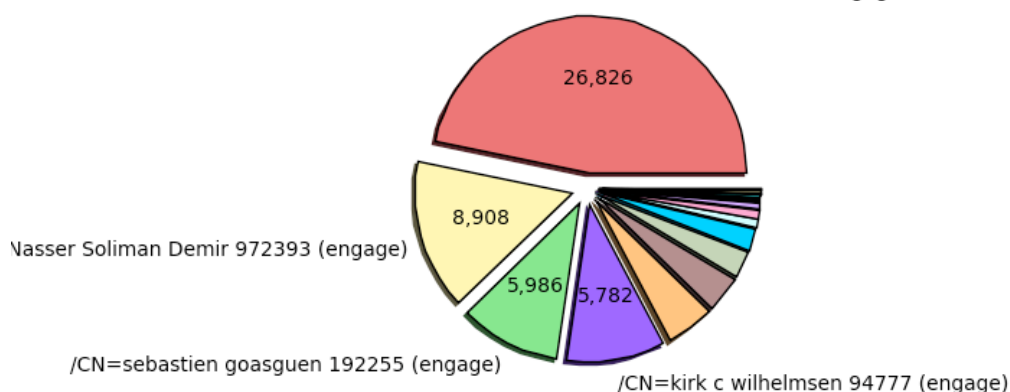
---

[5] For an example of a VO's CSV, see the engage one: http://engage-central.renci.org/engage-sciences.csv

whole OSG[6], then just for Engage[7].  Currently, all the uncategorized usage adds up to about 10 CPU years; once we remove the top 5 users, uncategorized usage will be about 3 CPU years.

## Wall Hours by Uncategorized User (Sum: 85,473 Hours)
### 52 Weeks from Week 31 of 2008 to Week 31 of 2009



/CN=zhi sun 343390 (engage)

26,826

belov/CN=633969/CN=Sergey Belov (glow)

11,860

Iasser Soliman Demir 972393 (engage)

8,908

- /CN=zhi sun 343390 (engage) (26,826)
- /CN=Nasser Soliman Demir 972393 (engage) (8,908)
- /CN=sebastien goasguen 192255 (engage) (5,987)
- /CN=Xi Li 905786 (osg) (4,758)
- /CN=osg-vo-test (osg) (2,796)
- /CN=Tai Boon Tan 288566 (engage) (1,590)
- /CN=Poornima 105779 (engage) (1,291)
- /CN=John McGee 941820 (engage) (498.00)
- /CN=osgmm (engage) (405.00)
- /CN=tino vazquez (osg) (264.00)

- /CN=sbelov/CN=633969/CN=Sergey Belov (glow) (11,860)
- /CN=Zhengxiong Hou 510271 (osg) (6,544)
- /CN=kirk c wilhelmsen 94777 (engage) (5,782)
- /CN=Terry Farrah 12349 (engage) (2,872)
- /CN=Erin Hodgess 974801 (engage) (2,090)
- Other (1,402)
- /CN=OSG Education student 02 885004 (osgedu) (561.00)
- /CN=Kieran Nolan 602232 (engage) (478.00)
- /CN=Anand Padmanabhan 75166 (osg) (303.00)
- /CN=Scot Kronenfeld (RSV) 622930 (osgedu) (259.00)

[6] http://t2.unl.edu/gratia/pie_graphs/unclassified_users
[7] http://t2.unl.edu/gratia/pie_graphs/unclassified_users?vo=engage

Wall Hours by Uncategorized User (Sum: 57,441 Hours)
52 Weeks from Week 31 of 2008 to Week 31 of 2009

/CN=zhi sun 343390 (engage)
26,826

8,908
Nasser Soliman Demir 972393 (engage)

5,986 5,782

/CN=sebastien goasguen 192255 (engage)

/CN=kirk c wilhelmsen 94777 (engage)

| | |
|---|---|
| /CN=zhi sun 343390 (engage) (26,826) | /CN=Nasser Soliman Demir 972393 (engage) (8,908) |
| /CN=sebastien goasguen 192255 (engage) (5,987) | /CN=kirk c wilhelmsen 94777 (engage) (5,782) |
| /CN=Terry Farrah 12349 (engage) (2,872) | /CN=Erin Hodgess 974801 (engage) (2,090) |
| /CN=Tai Boon Tan 288566 (engage) (1,590) | /CN=Poornima 105779 (engage) (1,291) |
| /CN=John McGee 941820 (engage) (498.00) | /CN=Kieran Nolan 602232 (engage) (478.00) |
| /CN=osgmm (engage) (405.00) | /CN=Michael Fenn 687928 (engage) (147.00) |
| /CN=LaToya Green 652696 (engage) (138.00) | /CN=vikas patel 555007 (engage) (133.00) |
| /CN=Vincent Bloch (engage) (101.00) | /CN=Michael Stealey 848926 (engage) (91.00) |
| /CN=Melinda Chin 47783 (engage) (62.00) | /CN=Shantanu Sharma 1118 (engage) (35.00) |
| /CN=Mary Kurz 660673 (engage) (7.00) | Other (0.00) |

VOs were given the following initial list of science fields to use to classify their users (but were able to add fields as needed):

- Astronomy

- Biochemistry

- Bioinformatics

- Biophysics

- Botany

- Cellular Biology

- Chemistry

- Community Grid

- Computer Science

- Earth Sciences

- Ecology

- Engineering

- HEP

- Information Theory

- Logic

- Mathematics

- Medicine

- Microbiology

- Molecular Biology

- Physics

- Physiology

- Statistics

- Technology

- USLHC

- Zoology

Some VOs did not fill in this information at registration time (many VOs registered before it was possible to fill in field of science). For these VOs, and for inactive VOs, we have put in a field of science to the best of our ability. VOs that are still active will be contacted and asked to review the data we filled in for them. The active VOs whose data we altered are:

- CIGI

- CompBioGrid

- Engagement

- Fermilab/Astro

- Fermilab

- Fermilab/Minerva

- Fermilab/Mu2e

- Fermilab/Nova

- Fermilab/Numi

- Fermilab/Patriot

- Fermilab/Test

- GLOW

- GPN

- GROW
- i2u2
- JDEM
- mariachi
- MIS
- nanoHUB
- Ops
- OSG
- OSGEDU