# How OSG Increases Goodput with Elasticsearch

Todd Tannenbaum / Jason Patton
Center for High Throughput Computing

# Usage reporting with Accounting ads

$ **condor_userprio** –usage

```
                            Wghted  Total Usage         Usage                    Last
User Name                   In Use  (wghted-hrs)     Start Time              Usage Time
------------------------    ------  ------------  ----------------    --------------------
todd@submit5.wisc.edu            0          0.04  10/10/2022 09:13    10/10/2022 09:22
fred@submit5.wisc.edu            0          1.02  10/10/2022 05:15    10/10/2022 09:29
------------------------    ------  ------------  ----------------    --------------------
Number of users: 2               0          1.06                      10/09/2022 09:29
```

# Usage reporting with Accounting ads

| Fm: 2022-05-16 | | Total | | CHTC | | OSG | | CHTC-HPC | |
|---|---|---|---|---|---|---|---|---|---|
| To: 2022-05-17 | | Hours | %Pool | Hours | %Pool | Hours | %Pool | Hours | %Pool |
| 88 | Projects | 946,364 | 100.0% | 308,057 | 32.6% | 1,511 | 0.2% | 65,066 | 6.9% |
| 1 | | 339,219 | 35.8% | 19,919 | 6.5% | 0 | 0.0% | 0 | 0.0% |
| 2 | | 117,976 | 12.5% | 17,180 | 5.6% | 605 | 40.1% | 0 | 0.0% |
| 3 | | 94,461 | 10.0% | 0 | 0.0% | 0 | 0.0% | 0 | 0.0% |
| 4 | | 44,959 | 4.8% | 44,959 | 14.6% | 0 | 0.0% | 0 | 0.0% |
| 5 | | 38,756 | 4.1% | 0 | 0.0% | 0 | 0.0% | 0 | 0.0% |
| 6 | | 32,009 | 3.4% | 32,009 | 10.4% | 0 | 0.0% | 0 | 0.0% |
| 7 | | 30,364 | 3.2% | 30,364 | 9.9% | 0 | 0.0% | 0 | 0.0% |
| 8 | | 28,576 | 3.0% | 0 | 0.0% | 0 | 0.0% | 28,576 | 43.9% |
| 9 | | 27,070 | 2.9% | 21,680 | 7.0% | 0 | 0.0% | 0 | 0.0% |
| 10 | | 23,381 | 2.5% | 23,381 | 7.6% | 0 | 0.0% | 0 | 0.0% |
| 11 | | 18,804 | 2.0% | 18,334 | 6.0% | 0 | 0.0% | 0 | 0.0% |
| 12 | | 17,412 | 1.8% | 17,412 | 5.7% | 0 | 0.0% | 0 | 0.0% |
| 13 | | 15,376 | 1.6% | 15,376 | 5.0% | 0 | 0.0% | 0 | 0.0% |
| 14 | | 12,731 | 1.3% | 0 | 0.0% | 0 | 0.0% | 12,731 | 19.6% |

Identities redacted

CENTER FOR HIGH THROUGHPUT COMPUTING

HTCondor

# Usage reporting with Accounting ads

› So, we delivered almost a million CPU core hours that day.
  - …was any of it good?

› Any other usage?
  - GPU hours?
  - Memory usage?
  - Files transferred?

› How was the user experience?
  - How often were jobs interrupted or put on hold?

CENTER FOR
HIGH THROUGHPUT
COMPUTING

HTCondor

# Storing job history in Elasticsearch

› We use the <u>condor_adstash</u> tool to periodically push job history ads from access points to Elasticsearch (ES).

› Wins:

+ Query-able history of all* job ads

+ New attributes do not have to be predefined before inserting ads

+ Libraries in popular languages for querying ES

+ Kibana web UI for simple queries and graphs

# Storing job history in Elasticsearch

› Gotchas

– Adstash does *remote* history queries, limited by knob setting HISTORY_HELPER_MAX_HISTORY (last 10,000 ads by default)

– *may miss ads on busy APs, especially if outages occur

– Unlike ClassAds (which may contain user-defined attrs), ES field names are case-sensitive and field values must have same type

– By default, Adstash converts unknown attr names to lowercase and types unknown fields as text

– (IMO) ES has a penchant for API breaking changes

– Example: Adstash suddenly broken with elasticsearch-py v8.0+

# Now what?

**52,918,476** hits

| GlobalJobId | QDate | RecordTime | RemoteWallClockTime |
|---|---|---|---|
| > submit3.chtc.wisc.edu#15282029.195702#1623463414 | 1,623,180,040 | 1,623,475,763 | 8,235 |
| > submit3.chtc.wisc.edu#15282029.200212#1623469354 | 1,623,180,040 | 1,623,475,760 | 2,808 |
| > submit3.chtc.wisc.edu#15282029.199736#1623468993 | 1,623,180,040 | 1,623,475,757 | 3,848 |
| > submit3.chtc.wisc.edu#15282029.200294#1623469394 | 1,623,180,040 | 1,623,475,756 | 2,670 |
| > submit3.chtc.wisc.edu#15282029.195730#1623463425 | 1,623,180,040 | 1,623,475,752 | 8,206 |
| > submit3.chtc.wisc.edu#15282029.200285#1623469392 | 1,623,180,040 | 1,623,475,752 | 2,692 |
| > submit3.chtc.wisc.edu#15282029.194963#1623462638 | 1,623,180,040 | 1,623,475,748 | 8,968 |
| > submit3.chtc.wisc.edu#15282029.199741#1623468998 | 1,623,180,040 | 1,623,475,744 | 3,829 |
| > submit3.chtc.wisc.edu#15282029.199725#1623468980 | 1,623,180,040 | 1,623,475,742 | 3,838 |
| > submit3.chtc.wisc.edu#15282029.199770#1623469047 | 1,623,180,040 | 1,623,475,737 | 3,780 |
| > submit3.chtc.wisc.edu#15282029.200232#1623469366 | 1,623,180,040 | 1,623,475,736 | 2,757 |
| > submit3.chtc.wisc.edu#15282029.200287#1623469392 | 1,623,180,040 | 1,623,475,735 | 2,674 |

CENTER FOR
HIGH THROUGHPUT
COMPUTING

HTCondor

# Usage reporting with ~~Accounting ads~~ condor_adstash

› **Was any of our usage good?**

› Any other usage?

- GPU hours?
- Memory usage?
- Files transferred?

› How was the user experience?

- How often were jobs interrupted or put on hold?

# Let's use our job history for good(put)!

**Goodput**

*noun*

1. the opposite of badput.


**Badput**

*noun*

1. claimed computing resources that did not contribute meaningfully to a requested computing task (i.e. to science).
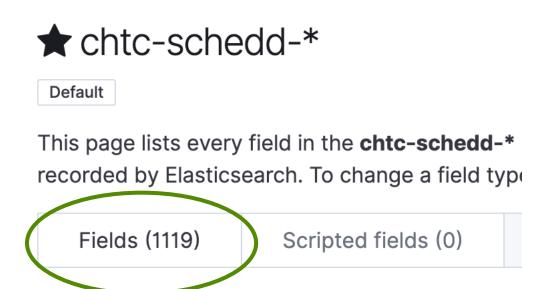
# Let's use our job history for good(put)!

› "good" CPU hours = total CPU hours – "bad" CPU hours

› What should count towards bad CPU hours?

- The time spent by any job execution that doesn't exit on its own accord or that doesn't exit due to user action.

- Evicted executions clearly lead to badput, what about held jobs and removed jobs?

- Assumption: Good CPU hours are CPU hours used in the last execution attempt (i.e. "final run") of a job.

**Can we calculate goodput from a job ad?**

# Calculating goodput from a job ad



$ condor_history -limit 1 -l | wc -l

**167**

Let's check page 490 of the HTCondor manual…

# Pop quiz!

Which pair of attributes provides the total runtime across all a job's runs and the runtime of a job's final run, respectively?

A. RemoteWallClockTime, CommittedTime

B. CommittedTime,      RemoteWallClockTime

C. RemoteWallClockTime,    LastRemoteWallClockTime

D. LastRemoteWallClockTime, RemoteWallClockTime

# Pop quiz!

Which pair of attributes provides the total runtime across all a job's runs and the runtime of a job's final run, respectively?

A. **RemoteWallClockTime, CommittedTime**

B. CommittedTime,    RemoteWallClockTime

emoteWallClockTime

**Undefined if job was removed!**

noteWallClockTime

CENTER FOR
HIGH THROUGHPUT
COMPUTING

HTCondor

# Pop quiz!

...ides the total runtime across all a...f a job's final run, respectively?

**Only exists since HTCondor 9.4.0**

A. RemoteWallClockTime, CommittedTime

B. CommittedTime, RemoteWallClockTime

C. **RemoteWallClockTime, LastRemoteWallClockTime**

D. LastRemoteWallClockTime, RemoteWallClockTime

# Calculating goodput from a job ad

› Current approach:

Total CPU Hours ~= CpusProvisioned * RemoteWallClockTime / 3600

Good CPU Hours ~= CpusProvisioned * {

                LastRemoteWallClockTime,

                CommittedTime,

                0

      } / 3600


› Finally, we can calculate goodput!

% Good CPU Hours = (Good CPU Hours/Total CPU Hours) * 100%

# Calculating goodput from a job ad

| | Project | Num Users | All CPU Hours | Num Uniq Job Ids | % Good CPU Hours | % Rm'd Jobs | % Short Jobs | % Jobs w/>1 Exec Att | % Jobs w/1+ Holds | % Jobs using S'ty | Total Files Xferd | Shadw Starts / Job Id | Exec Atts / Shadw Start | Holds / Job Id |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 39 | TOTAL | 49 | 465,485 | 429,041 | 80.8 | 5.2 | 16.7 | 5.6 | 5.6 | 48.2 | 3,528,787 | 1.18 | 0.892 | 0.13 |
| 1 | Identities redacted | 1 | 133,282 | 4,910 | 80.3 | 0.0 | 0.1 | 27.0 | 0.2 | 0.0 | 354,176 | 1.53 | 0.895 | 0.01 |
| 2 | | 1 | 71,443 | 8,902 | 89.1 | 0.0 | 0.0 | 18.6 | 0.0 | 0.0 | 98,725 | 1.24 | 0.986 | 0.00 |
| 3 | | 3 | 29,420 | 18,493 | 84.9 | 67.6 | 0.0 | 4.5 | 67.6 | 100.0 | 64,436 | 2.44 | 0.159 | 2.04 |
| 4 | | 3 | 23,586 | 36,350 | 83.5 | 3.0 | 1.5 | 1.3 | 0.9 | 98.5 | 274,211 | 1.03 | 0.974 | 0.04 |
| 5 | | 1 | 21,921 | 4,134 | 64.1 | 0.0 | 0.0 | 55.7 | 1.2 | 100.0 | 73,047 | 2.15 | 0.992 | 0.01 |
| 6 | | 2 | 19,987 | 77,567 | 90.5 | 0.0 | 0.2 | 2.5 | 0.0 | 100.0 | 551,357 | 1.04 | 0.991 | 0.00 |
| 7 | | 1 | 19,473 | 8,381 | 91.5 | 0.0 | 0.1 | 12.1 | 0.0 | 100.0 | 35,870 | 1.15 | 0.990 | 0.00 |
| 8 | | 1 | 18,405 | 5,732 | 90.1 | 0.0 | 0.0 | 15.5 | 0.0 | 0.0 | 48,050 | 1.22 | 0.980 | 0.00 |
| 9 | | 5 | 18,164 | 4,759 | 69.5 | 0.9 | 8.0 | 73.5 | 67.4 | 0.0 | 53,297 | 2.87 | 0.851 | 1.00 |
| 10 | | 1 | 17,833 | 26,623 | 80.6 | 25.2 | 0.0 | 10.0 | 25.2 | 0.0 | 182,931 | 1.23 | 0.899 | 0.25 |
| 11 | | 1 | 13,698 | 3,012 | 48.6 | 0.0 | 0.0 | 70.6 | 0.0 | 0.0 | 29,728 | 2.37 | 0.965 | 0.00 |
| 12 | | 1 | 13,606 | 1,632 | 89.5 | 0.0 | 0.0 | 20.2 | 0.0 | 0.0 | | 1.28 | 0.968 | 0.00 |
| 13 | | 1 | 13,133 | 3,505 | 81.3 | 0.0 | 0.9 | 24.3 | 0.0 | 0.0 | 30,382 | 1.42 | 0.945 | 0.00 |
| 14 | | 1 | 10,851 | 7,367 | 90.8 | 0.0 | 0.0 | 11.9 | 0.0 | 0.0 | 39,613 | 1.14 | 0.987 | 0.00 |

CENTER FOR HIGH THROUGHPUT COMPUTING

HTCondor

# Usage reporting with ~~Accounting ads~~ condor_adstash

› **Was any of our usage good?** ✔

› Any other usage?
- GPU hours?
- Memory usage?
- Files transferred?

› How was the user experience?
- How often were jobs interrupted or put on hold?

# Usage reporting with ~~Accounting ads~~ condor_adstash

› Was any of our usage good? ✔

› **Any other usage?** ✔

- GPU hours?

- Memory usage?

- Files transferred?

› How was the user experience?

- How often were jobs interrupted or

| | User | All CPU Hours | All GPU Hours | Num Uniq Job Ids | % Good CPU Hours | % Good GPU Hours | % Ckpt Able | % Rm'd Jobs |
|---|---|---|---|---|---|---|---|---|
| 13 | TOTAL | 6,837 | 2,564 | 1,728 | 97.9 | 94.5 | 0.0 | 1.6 |
| 1 | Identities redacted | 3,230 | 100 | 1 | 100.0 | 100.0 | 0.0 | 0.0 |
| 2 | | 2,054 | 2,054 | 1,575 | 93.9 | 93.9 | 0.0 | 1.4 |
| 3 | | 866 | 219 | 12 | 99.9 | 99.9 | 0.0 | 25.0 |
| 4 | | 259 | 25 | 3 | 100.0 | 100.0 | 0.0 | 0.0 |

CENTER FOR HIGH THROUGHPUT COMPUTING

HTCondor

# Usage reporting with ~~Accounting ads~~ condor_adstash

› Was any of it good? ✔

› Any other usage? ✔

- GPU hours?
- Memory usage?
- Files transferred?

› **How was the user experience?**

- How often were jobs interrupted or put on hold?

# Finding users who are having a bad time

| | Project | Num Users | All CPU Hours | Num Uniq Job Ids | % Good CPU Hours | % Rm'd Jobs | % Short Jobs | % Jobs w/>1 Exec Att | % Jobs w/1+ Holds | % Jobs using S'ty | Total Files Xferd | Shadw Starts / Job Id | Exec Atts / Shadw Start | Holds / Job Id |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 39 | TOTAL | 49 | 465,485 | 429,041 | 80.8 | 5.2 | 16.7 | 5.6 | 5.6 | 48.2 | 3,528,787 | 1.18 | 0.892 | 0.13 |
| 1 | Identities redacted | 1 | 133,282 | 4,910 | 80.3 | 0.0 | 0.1 | 27.0 | 0.2 | 0.0 | 354,176 | 1.53 | 0.895 | 0.01 |
| 2 | | 1 | 71,443 | 8,902 | 89.1 | 0.0 | 0.0 | 18.6 | 0.0 | 0.0 | 98,725 | 1.24 | 0.986 | 0.00 |
| 3 | | 3 | 29,420 | 18,498 | 84.9 | 67.6 | 0.0 | 4.5 | 67.6 | 100.0 | 64,436 | 2.44 | 0.159 | 2.04 |
| 4 | | 3 | 23,586 | 36,350 | 83.5 | 3.0 | 1.5 | 1.3 | 0.9 | 98.5 | 274,211 | 1.03 | 0.974 | 0.04 |
| 5 | | 1 | 21,921 | 4,134 | 64.1 | 0.0 | 0.0 | 55.7 | 1.2 | 100.0 | 73,047 | 2.15 | 0.992 | 0.01 |
| 6 | | 2 | 19,987 | 77,567 | 99.5 | 0.0 | 0.2 | 2.5 | 0.0 | 100.0 | 551,357 | 1.04 | 0.991 | 0.00 |
| 7 | | 1 | 19,473 | 8,381 | 91.5 | 0.0 | 0.1 | 12.1 | 0.0 | 100.0 | 35,870 | 1.15 | 0.990 | 0.00 |
| 8 | | 1 | 18,405 | 5,732 | 90.1 | 0.0 | 0.0 | 15.5 | 0.0 | 0.0 | 48,050 | 1.22 | 0.980 | 0.00 |
| 9 | | 5 | 18,164 | 4,759 | 69.5 | 0.9 | 8.0 | 73.5 | 67.4 | 0.0 | 53,297 | 2.87 | 0.851 | 1.00 |
| 10 | | 1 | 17,833 | 26,623 | 80.6 | 25.2 | 0.0 | 10.0 | 25.2 | 0.0 | 182,931 | 1.23 | 0.899 | 0.25 |
| 11 | | 1 | 13,698 | 3,012 | 48.6 | 0.0 | 0.0 | 70.6 | 0.0 | 0.0 | 29,728 | 2.37 | 0.965 | 0.00 |
| 12 | | 1 | 13,606 | 1,632 | 89.5 | 0.0 | 0.0 | 20.2 | 0.0 | 0.0 | | 1.28 | 0.968 | 0.00 |
| 13 | | 1 | 13,133 | 3,505 | 81.3 | 0.0 | 0.9 | 24.3 | 0.0 | 0.0 | 30,382 | 1.42 | 0.945 | 0.00 |
| 14 | | 1 | 10,851 | 7,367 | 90.8 | 0.0 | 0.0 | 11.9 | 0.0 | 0.0 | 39,613 | 1.14 | 0.987 | 0.00 |

CENTER FOR HIGH THROUGHPUT COMPUTING

HTCondor

# Finding users who are having a bad time

| | Project | Num Users | All CPU Hours | Num Uniq Job Ids | % Good CPU Hours | % Rm'd Jobs | % Short Jobs | % Jobs w/>1 Exec Att | % Jobs w/1+ Holds | % Jobs using S'ty | Total Files Xferd | Shadw Starts / Job Id | Exec Atts / Shadw Start | Holds / Job Id |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 39 | TOTAL | 49 | 465,485 | 429,041 | 80.8 | 5.2 | 16.7 | 5.6 | 5.6 | 48.2 | 3,528,787 | 1.18 | 0.892 | 0.13 |
| 1 | Identities redacted | 1 | 133,282 | 4,910 | 80.3 | 0.0 | 0.1 | 27.0 | 0.2 | 0.0 | 354,176 | 1.53 | 0.895 | 0.01 |
| 2 | | 1 | 71,443 | 8,902 | 89.1 | 0.0 | 0.0 | 18.6 | 0.0 | 0.0 | 98,725 | 1.24 | 0.986 | 0.00 |
| 3 | | 3 | 29,420 | 18,498 | 84.9 | 67.6 | 0.0 | 5 | 67.6 | 100.0 | 64,436 | 2.44 | 0.159 | 2.04 |
| 4 | | 3 | 23,586 | 36,350 | 83.5 | 3.0 | 15 | 13 | 0.9 | 98.5 | 274,211 | 1.03 | 0.974 | 0.04 |
| 5 | | 1 | 21,921 | 4,134 | 64.1 | 0.0 | 0.0 | 55.7 | 1.2 | 100.0 | 73,047 | 2.15 | 0.992 | 0.01 |
| 6 | | 2 | 19,987 | 77,567 | 90.5 | 0.0 | 0.2 | 2.5 | 0.0 | 100.0 | 551,357 | 1.04 | 0.991 | 0.00 |
| 7 | | 1 | 19,473 | 8,381 | 91.5 | 0.0 | 0.1 | 12.1 | 0.0 | 100.0 | 35,870 | 1.15 | 0.990 | 0.00 |
| 8 | | 1 | 18,405 | 5,732 | 90.1 | 0.0 | 0.0 | 15.5 | 0.0 | 0.0 | 48,050 | 1.22 | 0.980 | 0.00 |
| 9 | | 5 | 18,164 | 4,759 | 69.5 | 0.9 | 8.0 | 73.5 | 67.4 | 0.0 | 53,297 | 2.87 | 0.851 | 1.00 |
| 10 | | 1 | 17,833 | 26,623 | 80.6 | 25.2 | 0.0 | 10.0 | 25.2 | 0.0 | 182,931 | 1.23 | 0.899 | 0.25 |
| 11 | | 1 | 13,698 | 3,012 | 48.6 | 0.0 | 0.0 | 70.6 | 0.0 | 0.0 | 29,728 | 2.37 | 0.965 | 0.00 |
| 12 | | 1 | 13,606 | 1,632 | 89.5 | 0.0 | 0.0 | 20.2 | 0.0 | 0.0 | | 1.28 | 0.968 | 0.00 |
| 13 | | 1 | 13,133 | 3,505 | 81.3 | 0.0 | 0.9 | 24.3 | 0.0 | 0.0 | 30,382 | 1.42 | 0.945 | 0.00 |
| 14 | | 1 | 10,851 | 7,367 | 90.8 | 0.0 | 0.0 | 11.9 | 0.0 | 0.0 | 39,613 | 1.14 | 0.987 | 0.00 |

CENTER FOR HIGH THROUGHPUT COMPUTING

HTCondor

# Finding users who are having a bad time

| | Project | Num Users | All CPU Hours | Num Uniq Job Ids | % Good CPU Hours | % Rm'd Jobs | % Short Jobs | % Jobs w/>1 Exec Att | % Jobs w/1+ Holds | % Jobs using S'ty | Total Files Xferd | Shadw Starts / Job Id | Exec Atts / Shadw Start | Holds / Job Id |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 39 | TOTAL | 49 | 465,485 | 429,041 | 80.8 | 5.2 | 16.7 | 5.6 | 5.6 | 48.2 | 3,528,787 | 1.18 | 0.892 | 0.13 |
| 1 | Identities redacted | 1 | 133,282 | 4,910 | 80.3 | 0.0 | 0.1 | 27.0 | 0.2 | 0.0 | 354,176 | 1.53 | 0.895 | 0.01 |
| 2 | | 1 | 71,443 | 8,902 | 89.1 | 0.0 | 0.0 | 18.6 | 0.0 | 0.0 | 98,725 | 1.24 | 0.986 | 0.00 |
| 3 | | 3 | 29,420 | 18,498 | 84.9 | 67.6 | 0.0 | 4.5 | 67.6 | 100.0 | 64,436 | 2.44 | 0.159 | 2.04 |
| 4 | | 3 | 23,586 | 36,350 | 83.5 | 3.0 | 1.5 | 1.3 | 0.9 | 98.5 | 274,211 | 1.03 | 0.974 | 0.04 |
| 5 | | 1 | 21,921 | 4,134 | 64.1 | 0.0 | 0.0 | 55.7 | 1.2 | 100.0 | 73,047 | 2.15 | 0.992 | 0.01 |
| 6 | | 2 | 19,987 | 77,567 | 90.5 | 0.0 | 0.2 | 2.5 | 0.0 | 100.0 | 551,357 | 1.04 | 0.991 | 0.00 |
| 7 | | 1 | 19,473 | 8,381 | 91.5 | 0.0 | 0.1 | 12.1 | 0.0 | 100.0 | 35,870 | 1.15 | 0.990 | 0.00 |
| 8 | | 1 | 18,405 | 5,732 | 90.1 | 0.0 | 0.0 | 15.5 | 0.0 | 0.0 | 48,050 | 1.22 | 0.980 | 0.00 |
| 9 | | 5 | 18,164 | 4,759 | 69.5 | 0.9 | 8.0 | 73.5 | 67.4 | 0.0 | 3,297 | 2.87 | 0.851 | 1.00 |
| 10 | | 1 | 17,833 | 26,623 | 80.6 | 25.2 | 0.0 | 10.1 | 25.2 | 0.0 | 182,931 | 1.23 | 0.899 | 0.25 |
| 11 | | 1 | 13,698 | 3,012 | 48.6 | 0.0 | 0.0 | 70.6 | 0.0 | 0.0 | 29,728 | 2.37 | 0.965 | 0.00 |
| 12 | | 1 | 13,606 | 1,632 | 89.5 | 0.0 | 0.0 | 20.2 | 0.0 | 0.0 | | 1.28 | 0.968 | 0.00 |
| 13 | | 1 | 13,133 | 3,505 | 81.3 | 0.0 | 0.9 | 24.3 | 0.0 | 0.0 | 30,382 | 1.42 | 0.945 | 0.00 |
| 14 | | 1 | 10,851 | 7,367 | 90.8 | 0.0 | 0.0 | 11.9 | 0.0 | 0.0 | 39,613 | 1.14 | 0.987 | 0.00 |

# Finding users who are having a bad time

› Additional reports have shown to be helpful, such as reporting on all jobs that had at least one hold event.

| | Project | Num Users | All CPU Hours | Num Uniq Job Ids | Most Common Hold Reason | % Holds Most Comm Reas | Holds / Job Id | % Good CPU Hours | % Rm'd Jobs | |
|---|---|---|---|---|---|---|---|---|---|---|
| 14 | TOTAL | 16 | 27,201 | 24,199 | Download FileError | 54.6 | 2.30 | 39.8 | 86.4 | |
| 1 | Identities redacted | 2 | 14,031 | 3,207 | JobPolicy | 99.9 | 1.48 | 67.8 | 1.1 | |
| 2 | | 1 | 6,335 | 1,152 | StartdHeld Job | 75.5 | 3.90 | 0.0 | 100.0 | |
| 3 | | 1 | 2,049 | 48 | StartdHeld Job | 100.0 | 3.12 | 38.6 | 0.0 | |
| 4 | | 1 | 1,785 | 9 | Job Execute Exceeded | 100.0 | 3.00 | 20.7 | 0.0 | |
| 5 | | 1 | 1,255 | 6,715 | JobPolicy | 100.0 | 1.00 | 0.0 | 100.0 | |
| 6 | | 2 | 1,039 | 12,499 | Download FileError | 80.5 | 3.02 | 0.4 | 100.0 | |
| 7 | | 1 | 317 | 50 | UploadFile Error | 100.0 | 1.00 | 44.2 | 0.0 | |
| 8 | | 1 | 302 | 329 | System Policy | 99.9 | 4.75 | 0.0 | 100.0 | |

CENTER FOR HIGH THROUGHPUT COMPUTING

HTCondor

# Usage reporting with ~~Accounting ads~~ condor_adstash

› Was any of it good? ✔

› Any other usage? ✔

- GPU hours?

- Memory usage?

- Files transferred?

› **How was the user experience?** ✔

- How often were jobs interrupted or put on hold?

# Improving HTCondor

› This project has prompted many additions to the job ad:

LastRemoteWallClockTime = 3764

NumHoldsByReason = [ UserRequest = 2;      JobPolicy = 10; UnableToOpenInput = 1 ]

TransferInputStats = [ CedarFilesCountTotal =        5; CedarFilesCountLastRun = 5 ]

# Remaining challenges

› How to find strangely behaved or broken "sites"?

- Example: Job runs 3 times at Site A, failing to transfer output each time, before running and completing successfully at Site B.

- Job ads lack information about intermediate job runs, must infer from cumulative and last run stats.
  - Startd History file on the EP
  - Upcoming: Instance History file on the AP

› How to determine if jobs are checkpointing *correctly*?

- Are intermediate runs contributing to goodput or not?

# Want to try it ?

› Setup Elasticsearch

- Create an index named "`htcondor-000001`"

› Install HTCSS

- Add to HTCSS config:

  ```
  use feature: adstash
  ```

- Details: See Manual / Admin Manual / Monitoring / Elasticsearch

› Checkout and/or customize our report generation scripts

- https://github.com/CHTC/JobAccounting

# Thank You!

## Follow us on Twitter!
## https://twitter.com/HTCondor

**PATh** PARTNERSHIP to ADVANCE THROUGHPUT COMPUTING

HT CENTER FOR HIGH THROUGHPUT COMPUTING

HTCondor