# CMS Monitoring at the LPC

*Compact Muon Solenoid*

*LHC Physics Center*

User job monitoring at the FNAL LPC,

*from the support perspective*

**Marguerite Tonjes**, LPC Computing Support

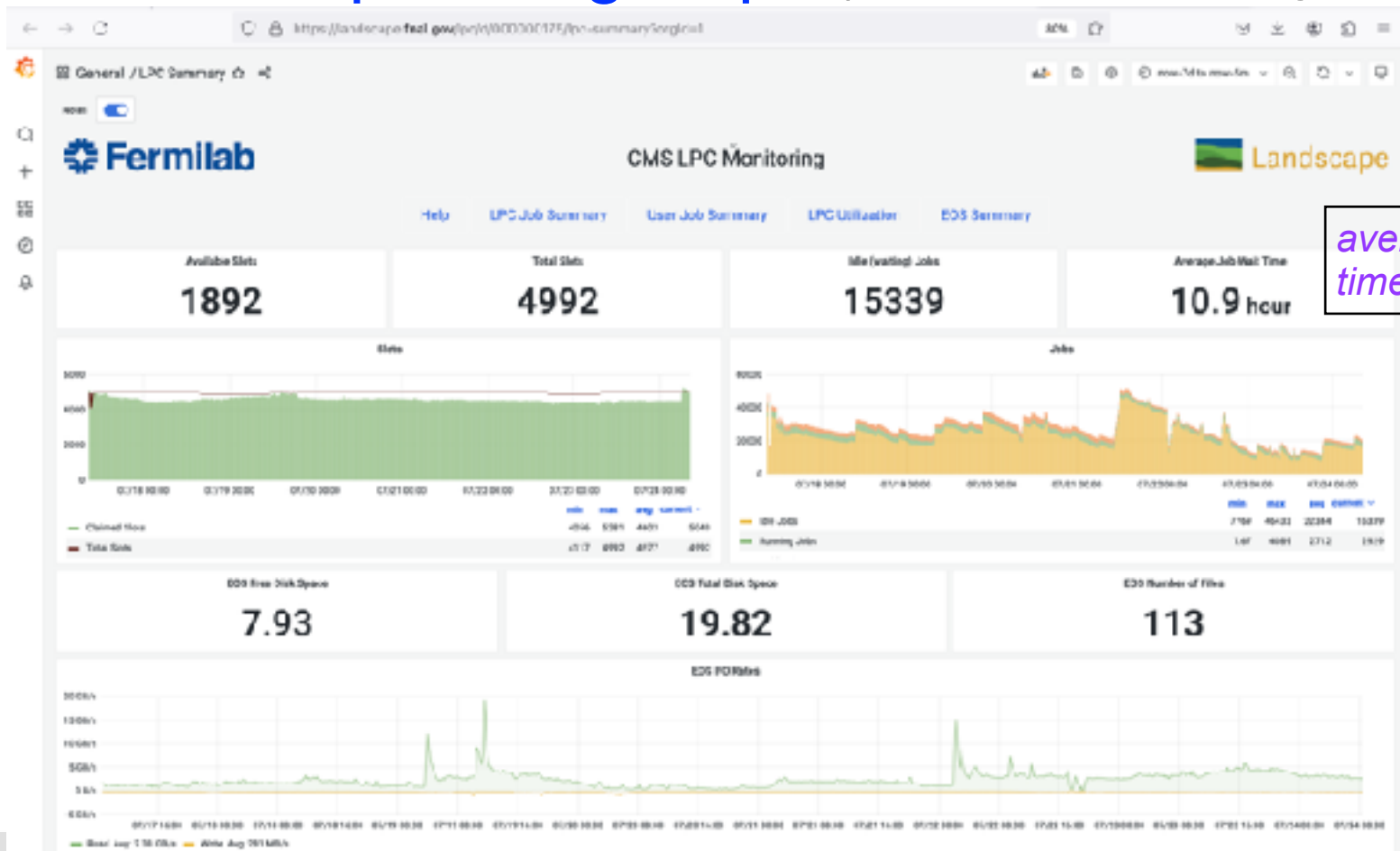*University of Illinois, Chicago*

*Senior Research Specialist*

# CMS LPC User jobs

- CMS as a whole has a fairly complex infrastructure, users often want to know:
  - ◉ Are my jobs running?
  - ◉ Why aren't they running?
  - ◉ Where and how did they fail?
- Jobs are submitted through local cmslpc interactive HTCondor, or through grid via [CRAB](#) (CMS Remote Analysis Builder) or via CMS Connect (HTCondor)
  - ◉ Note the CRAB and CMS Connect access the whole of the CMS computing grid (which includes opportunistic space on FNAL Tier1 as well as User Tier3
  - ◉ Only users with FNAL computing accounts access cmslpc (T3_US_FNALLPC - Tier3) HTCondor batch nodes from any source
  - ◉ CMS grid certificate used for grid job authentication (needed as well for most remote file reading), must also be stored in local database
- *Scale*:
  - ◉ **723** unique users logged in interactively 2023-2024 *(.Xauthority)*
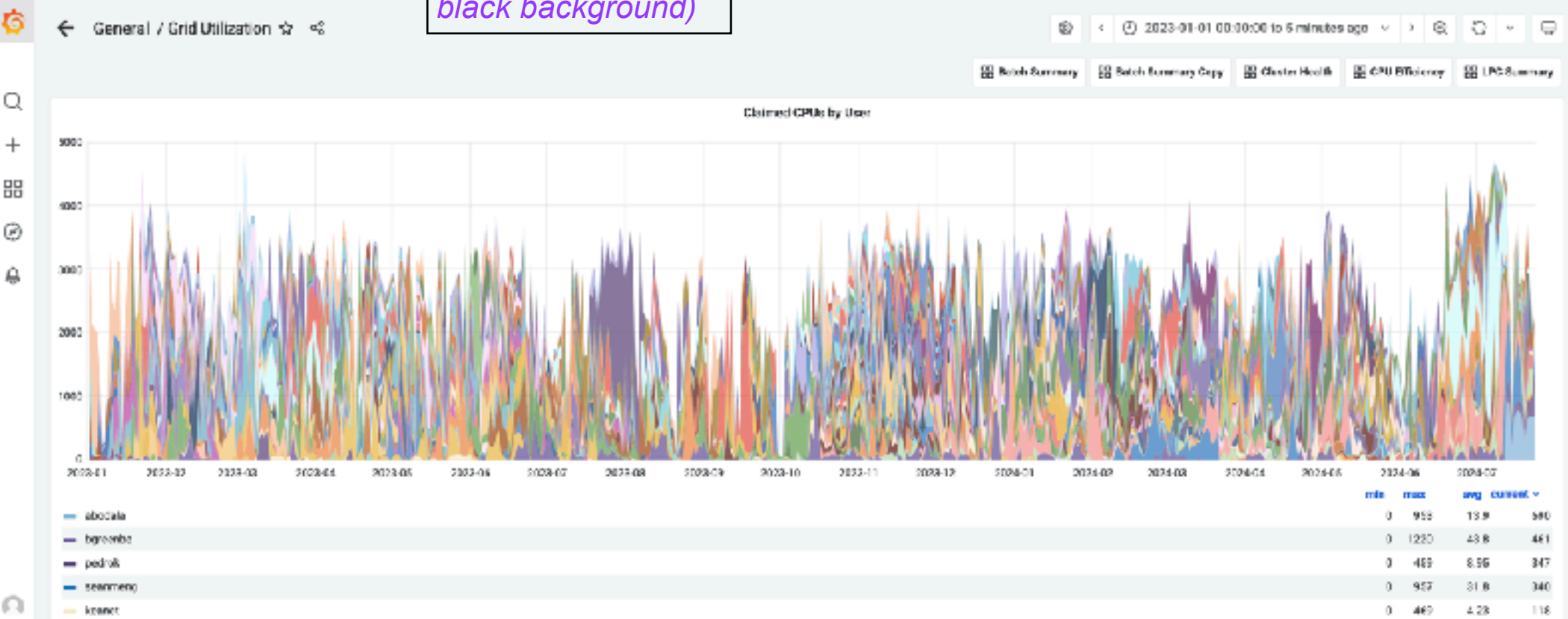  - ◉ **480** unique users running cmslpc batch 2023-2024 *(HTCondor landscape)*

# CMS LPC monitoring

- [https://uscms.org/uscms_at_work/physics/computing/status/index.shtml](https://uscms.org/uscms_at_work/physics/computing/status/index.shtml) collection of monitoring links useful to CMS LPC users at FNAL

- [https://landscape.fnal.gov/lpc](https://landscape.fnal.gov/lpc) (authenticate CMS grid certificate)



*average job wait time can be large!*

*(note - normally black background)*

- From "LPC Utilization"
  - ◉ What we show management (US CMS group leads; funding)
  - ◉ *Dip due to a bug in schedulers fixed Monday*

# Are my jobs running?

- cmslpc HTCondor or CMS connect: many users just use command line checks like `condor_q`
  - ◉ https://uscms.org/uscms_at_work/computing/setup/batch_troubleshoot.shtml#Troubleshooting describes useful troubleshooting techniques for cmslpc
  - ◉ Unfortunately recent HTCondor software update removed CPU Time and Memory reporting
  - ◉ Landscape User Batch Summary (example)
- CRAB command line: `crab status`

```
crab status -d crabsubmit/crab_cmsdas_minbias_test0
CRAB project directory:     /uscms_data/d3/username/cmsdas/CMSSW_13_0_13_mcgen/src/crabsubmit/
crab_cmsdas_minbias_test0
Task name:              231110_212908:username_crab_cmsdas_minbias_test0
Grid scheduler - Task Worker:   crab3@vocms0198.cern.ch - crab-prod-tw01
Status on the CRAB server:  SUBMITTED
Task URL to use for HELP:   https://cmsweb.cern.ch/crabserver/ui/task/
231110_212908%3Ausername_crab_cmsdas_minbias_test0
Dashboard monitoring URL:   https://monit-grafana.cern.ch/d/cmsTMDetail/cms-task-monitoring-task-view?
orgId=11&var-user=username&var-task=231110_212908%3Ausername_crab_cmsdas_minbias_test0&from=1699648148000&to=now
Status on the scheduler:    SUBMITTED

Jobs status:                        idle              100.0% (10/10)

No publication information available yet
Log file is /uscms_data/d3/username/cmsdas/CMSSW_13_0_13_mcgen/src/crabsubmit/crab_cmsdas_minbias_test0/crab.log
```
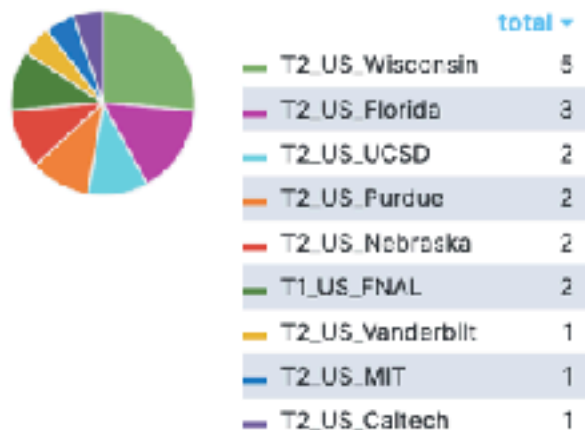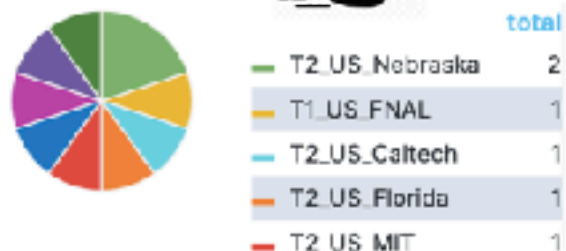
*not yet running*

# CRAB user dashboard 2

- Depending on the error code, CRAB will automatically retry jobs. This particular analysis went to US sites.



- Where did the job process?

# CRAB user dashboard list

- User can click on JobLog to look at analysis stdout/ stderr. ExitCodes have a reference *(8021: file read error)*
- PostJob is the transfer: handled centrally by FTS *(File Transfer Service)* from a temporary location on the remote processing site to the final file destination
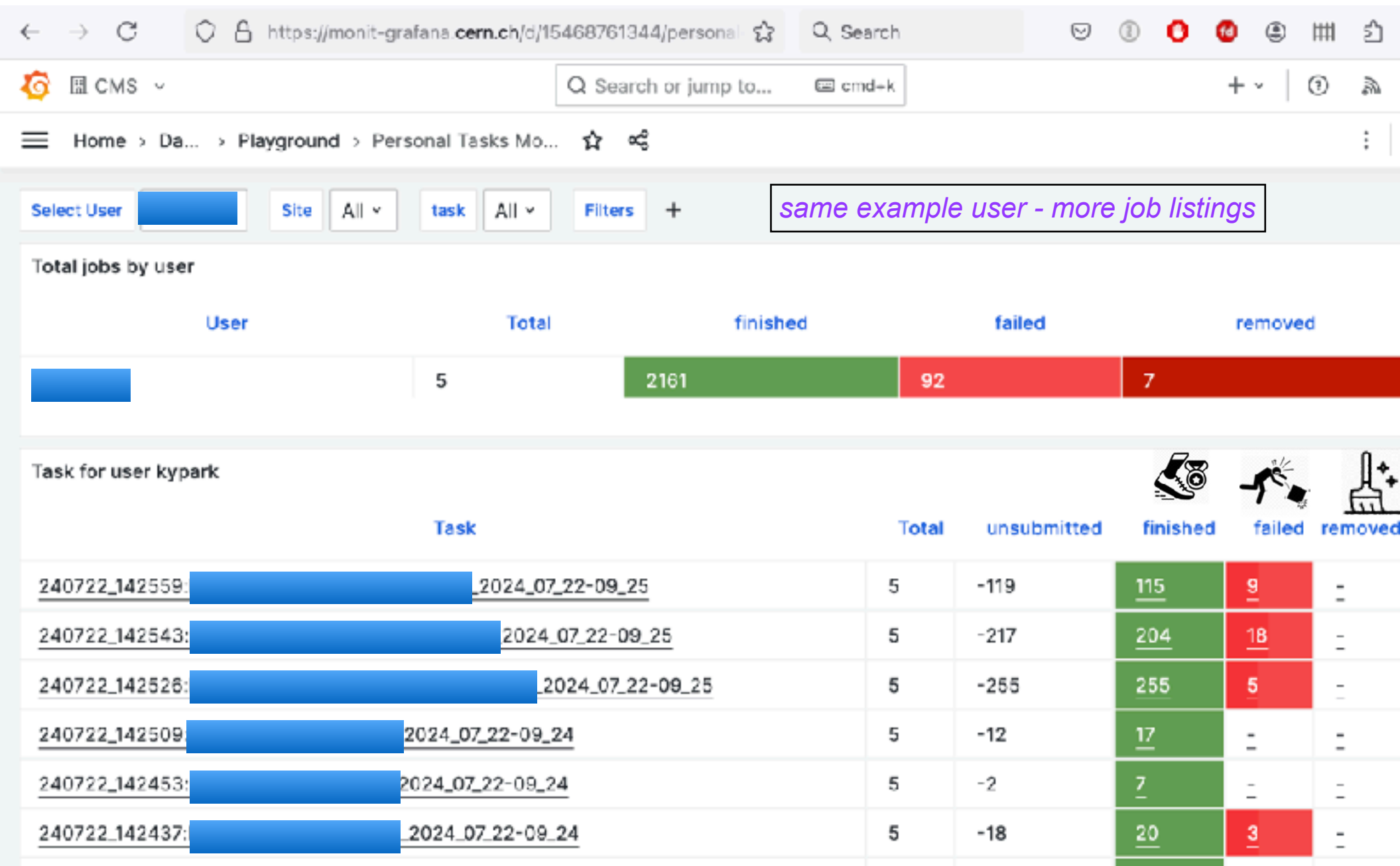
*same example user - more job listings*

**Total jobs by user**

| User | Total | finished | failed | removed |
|------|-------|----------|--------|---------|
| | 5 | 2161 | 92 | 7 |

**Task for user kypark**

| Task | Total | unsubmitted | finished | failed | removed |
|------|-------|-------------|----------|--------|---------|
| 240722_142559: _____2024_07_22-09_25 | 5 | -119 | 115 | 9 | - |
| 240722_142543: _____2024_07_22-09_25 | 5 | -217 | 204 | 18 | - |
| 240722_142526: _____2024_07_22-09_25 | 5 | -255 | 255 | 5 | - |
| 240722_142509: _____2024_07_22-09_24 | 5 | -12 | 17 | - | - |
| 240722_142453: _____2024_07_22-09_24 | 5 | -2 | 7 | - | - |
| 240722_142437: _____2024_07_22-09_24 | 5 | -18 | 20 | 3 | - |

# CRAB UI (check job specifics)

- Failed to run or failed to finish properly?
  - So many possibilities:
    - Asked for too much Memory per job and few groups of slots on machines available
    - Crashed in remote dataset
    - Went over limits (48 hours; 40GB local space, etc.)
    - Failed to transfer output (they're over quota; rarer disk error at the LPC; problem with FTS/transfer from grid site)
    - Failed to read some/all files (problem at site; user(s) read too many times in parallel same file and overloaded network; parsed input wrong; file on tape and not disk, ...)
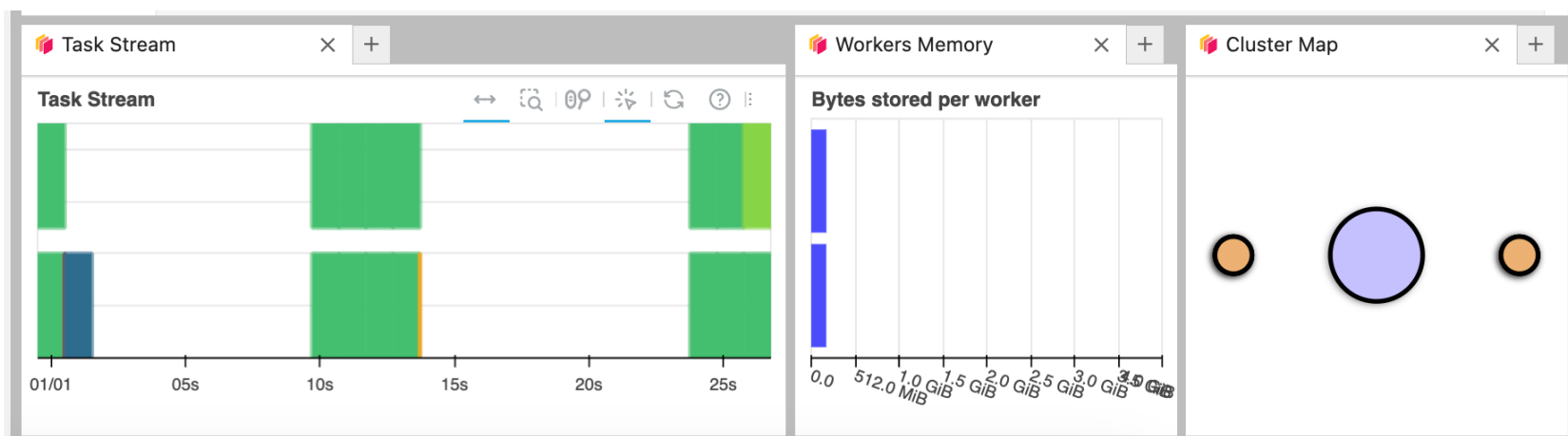    - User error
    - etc...
  - Note that CRAB produces error codes for various cases which can be used to debug many problems *(not full number range for each)*
    - 1 - 512 (Unix); 7000 - 9000 (CMSSW exit codes); 10000 - 19999 (environment setup); 50000 - 59999 (executable); 60000 - 69999 (staging out); 70000 - 79999 (WMAgent - job transfer); 80000-99999 (CRAB and other: only 6)

# HTCondor troubleshooting

- Typically we tell people to check `condor_q -better-analyze` and `.stdout`/`.stderr`

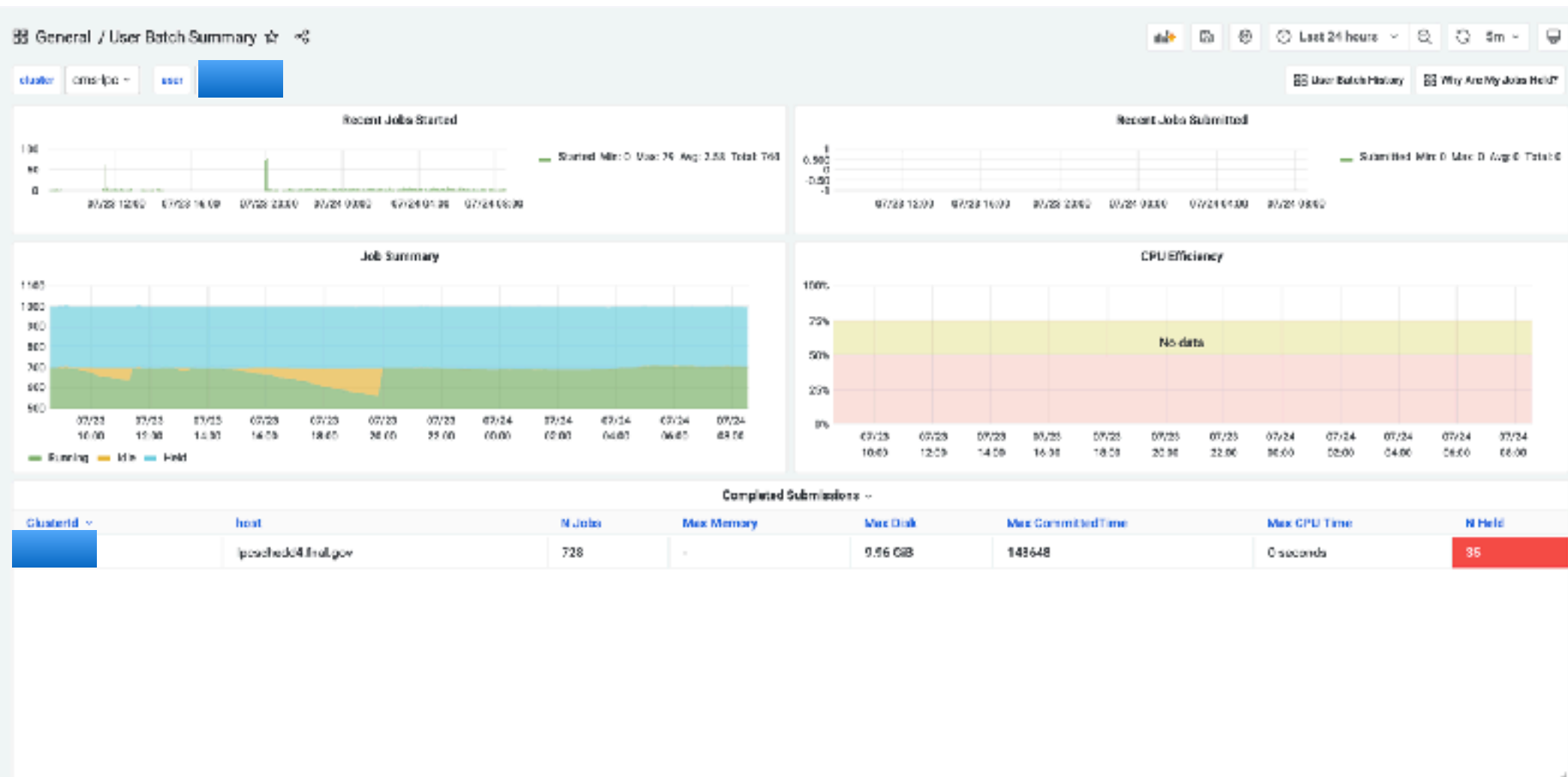  - ◉ General recommendation to put the following in .stdout:

```
echo "Starting job on " `date` #Date/time of start of job
echo "Running on: `uname -a`" #Condor job is running on this node
echo "System software: `cat /etc/redhat-release`" #Operating System job is running in
```

- Some users developed their own job management [scripts](scripts)

- The [coffea](coffea) (columnar analysis) workflow has its own job monitoring for dask ([coffea-casa](coffea-casa); [Purdue](Purdue))



*example from: Dmitry Kondratyev (Purdue University (US)) (from LPC tutorial on dask)*

# CMS LPC Landscape User job



*memory and CPU time missing: HTCondor bug :(*

# Too much memory

# Frequently: Cannot Access Data

- Note many users use xrootd: which gives **A**ny Data, **A**ny Time, **A**nywhere

    ◉ Important factor is that software needs fallback built in (CMSSW)

- We have a Data Aggregation Service (web based and command line) to locate data/Monte Carlo and understand what and where it is and if it's valid

- Can check the <u>CMS Site Status Board</u> to check outages (sites are clickable) - also tells you status of remote compute resources



*also number of site tickets*

# Disk & node health

- I find users don't show me they have checked disk/node health in troubleshooting
  - ◉ They like to open tickets and blame sites
  - ◉ Maybe I solve too much too quickly so they don't check themselves - also a lot of things are very clearly [documented](#), this is a **big help** not to be underestimated

- FNAL [EOS has a landscape link](#)

- Only if one of the 50 (Alma8) or 50 (Alma9) interactive nodes is problematic do users check [SSI metrics](#) (requires FNAL VPN or onsite fgz)

- Worker nodes have been so stable and monitored well I haven't checked their status to support tickets in a couple years! Same with NFS *(users don't have access to those checks)*

- Quota: command line (`quota -s`; `eosquota`)

# Other monitoring

- We have a community mailing list as well as community chat ([LPC gethelp web page](#)) and people will frequently just ask "this thing broke with this error, is FNAL EOS disk working?"
  - ◉ No seriously, a high number of "yes I also have this error" shows a problem, just like a high number of tickets shows it
  - ◉ This has led to tickets (strongly encouraged)
  - ◉ This has led to community instruction on how to access data, fixing problems of misunderstanding old analysis recipes, etc.

- CMS LPC emails a user list for outages and I also post updates on the chat (*a lot of emails were @fnal.gov for users that dropped off when forwarding stopped*)
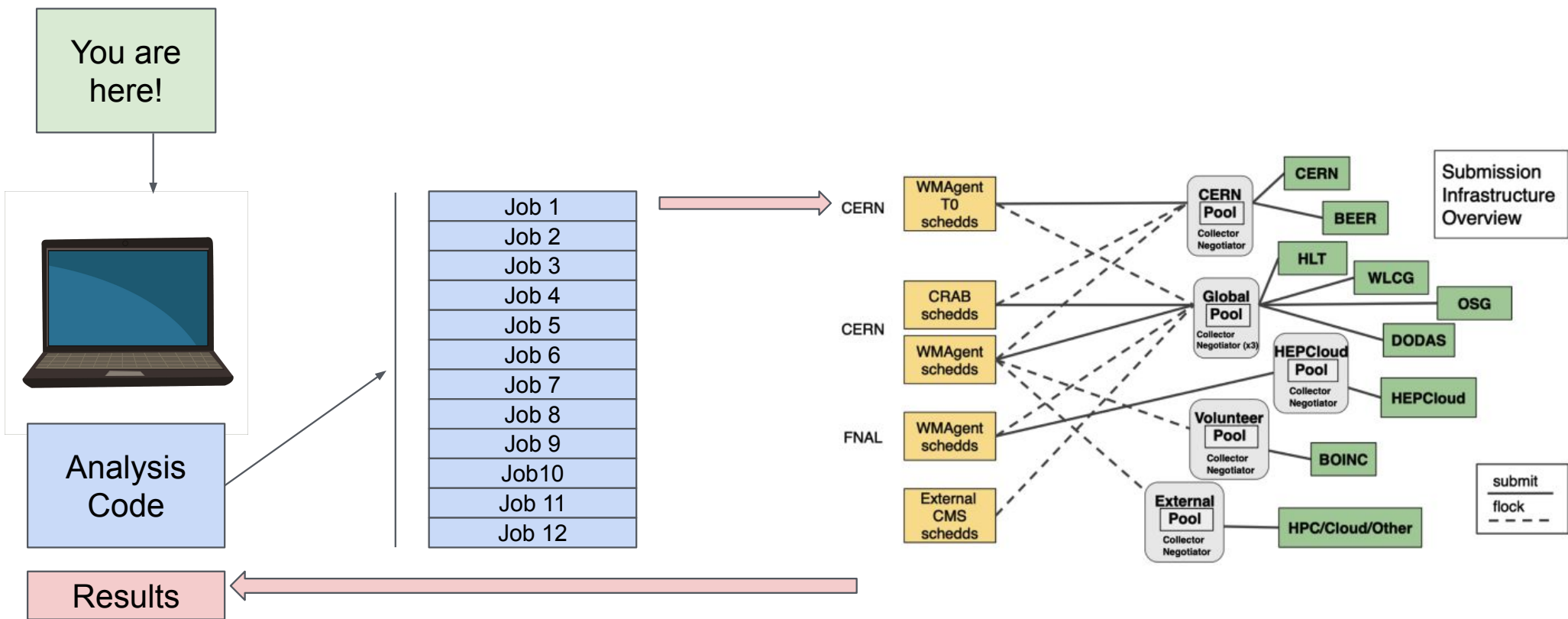
# Conclusion/Questions?

- I only focused on user monitoring and did not cover a lot more of CMS central production

- There are many more tools available to understand all of these systems (CRAB servers; database servers; xrootd file servers: FTS servers: etc...) not covered here



*most icons from IconScout.com*

taken from Andrew Melo's presentation to USCMS PURSUE interns 2024

# CRAB user dashboard 3



**Jobs by Execution Site - all retries**

**All Completed Jobs**

| | total ▾ |
|---|---|
| T2_US_Wisconsin | 5 |
| T2_US_Florida | 3 |
| T2_US_UCSD | 2 |
| T2_US_Purdue | 2 |
| T2_US_Nebraska | 2 |
| T1_US_FNAL | 2 |
| T2_US_Vanderbilt | 1 |
| T2_US_MIT | 1 |
| T2_US_Caltech | 1 |

**Failed and Retried Jobs**

No data points

*bug? There should be data here...*

**Jobs in PostProcessing**

| | total |
|---|---|
| T2_US_Purdue | 1 |
| T2_US_Florida | 1 |
| T1_US_FNAL | 1 |

*CRAB is reading memory even though our HTCondor doesn't right now*

**Average resource use by site (All Completed jobs)**

**Average memory (GB) by site**

| | |
|---|---|
| T2_US_Wisconsin | 1.90 |
| T2_US_Florida | 1.76 |
| T2_US_UCSD | 1.55 |
| T2_US_Purdue | 1.73 |
| T2_US_Nebraska | 1.76 |
| T1_US_FNAL | 1.72 |
| T2_US_Vanderbilt | 1.29 |

**Average wall time (hours) by site**

| | |
|---|---|
| T2_US_Wisconsin | 1.56 |
| T2_US_Florida | 0.97 |
| T2_US_UCSD | 0.25 |
| T2_US_Purdue | 1.23 |
| T2_US_Nebraska | 0.79 |
| T1_US_FNAL | 1.44 |
| T2_US_Vanderbilt | 0.91 |

**Average CPU time (Hours) by site**

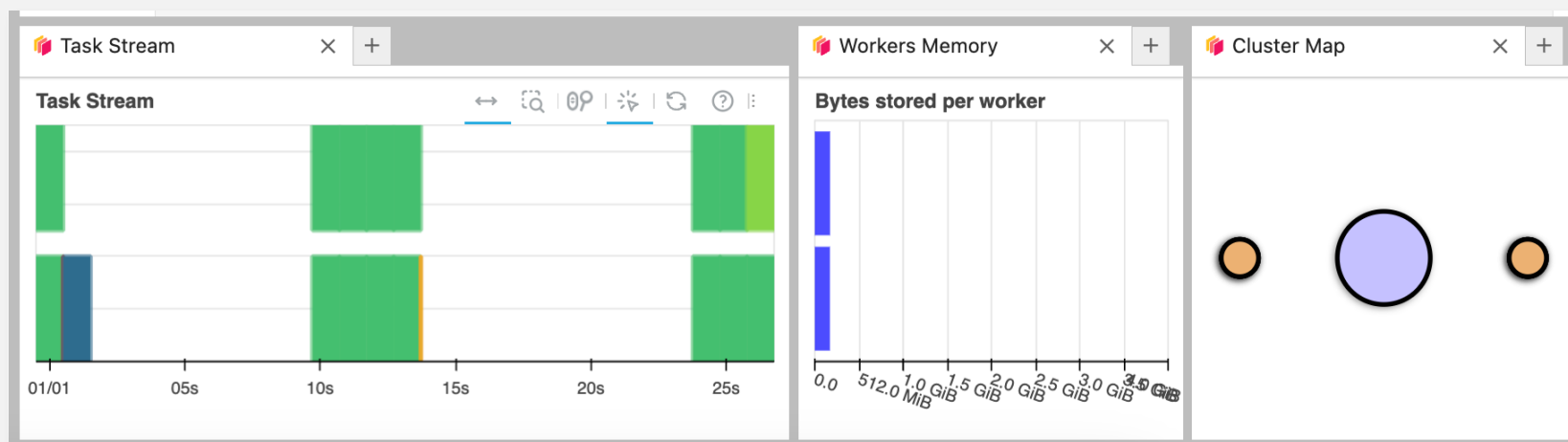| | |
|---|---|
| T2_US_Wisconsin | 1.24 |
| T2_US_Florida | 0.87 |
| T2_US_UCSD | .11 |
| T2_US_Purdue | 1.00 |
| T2_US_Nebraska | 0.67 |
| T1_US_FNAL | 1.31 |
| T2_US_Vanderbilt | 0.17 |

# Example of [coffea](#) dask monitoring

- At Purdue Tier2, for batch jobs: Dmitry Kondratyev (Purdue University (US)) (from LPC tutorial on dask)

- You can drag-and-drop panels to place them side by side with other tabs.

# LPC [Coffea](#) Dask tutorial

- At Purdue Tier2, for batch jobs: Dmitry Kondratyev (Purdue University (US)) (from LPC tutorial on dask)
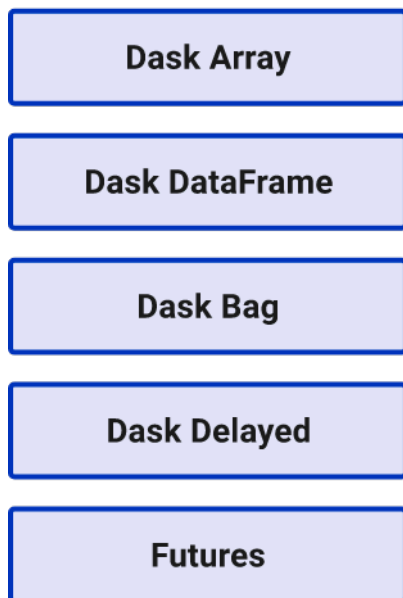
- Monitor the execution in the dashboard panels:



- **These tasks now run in a distributed mode on remote machines!**

# So what does that set of words really mean? [Coffea](#) **Dask**
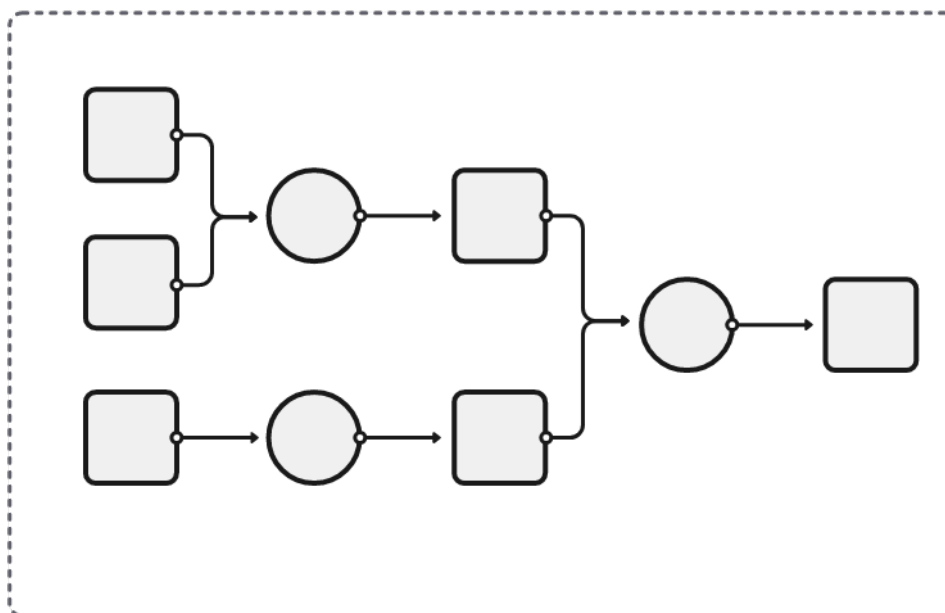
**Collections** → **Task Graph** → **Schedulers**
(create task graphs) (execute task graphs)



- You use collections to write straightforward python

- That code generates an abstract, declarative, description of your analysis
  - It can then be executed by anything that implements the collection's array interface!
  - This makes analysis code extremely portable for tradeoff in underlying complexity

- I hope to dig into this complexity enough so you can reason about task graphs

20 June 2024   L. Gray | Intro. to Dask and Dask-Awkward

🔶 **Fermilab**

# Columnar Analysis

In a "traditional" analysis, each event is processed one-at-a-time, meaning:

- An event is loaded and the appropriate values are extracted
- Some computation happens over this single event
- The temporary space is cleaned and the next event is loaded

In a "columnar" analysis, whole *batches* of events are processed at once, meaning

- 100s or even 1000s of events are loaded at once
- A bulk computation is done over the whole batch
- The temporary space is cleaned and the next event is loaded

The lower overhead for this method is particularly attractive for Data Scientists and is being embraced by CMS

Marguerite B. Tonjes          *CMS User Monitoring at the LPC*          *July 24, 2024*