

This manuscript has been authored by Fermi Research Alliance, LLC under Contract No. DE-AC02-07CH11359 with the U.S. Department of Energy, Office of Science, Office of High Energy Physics.



# Transitioning Users to SciTokens and Getting them Closer to HTCondor with Jobsub\_lite

Shreyas Bhat on behalf of the Jobsub Team

July 14, 2023

Throughput Computing 2023

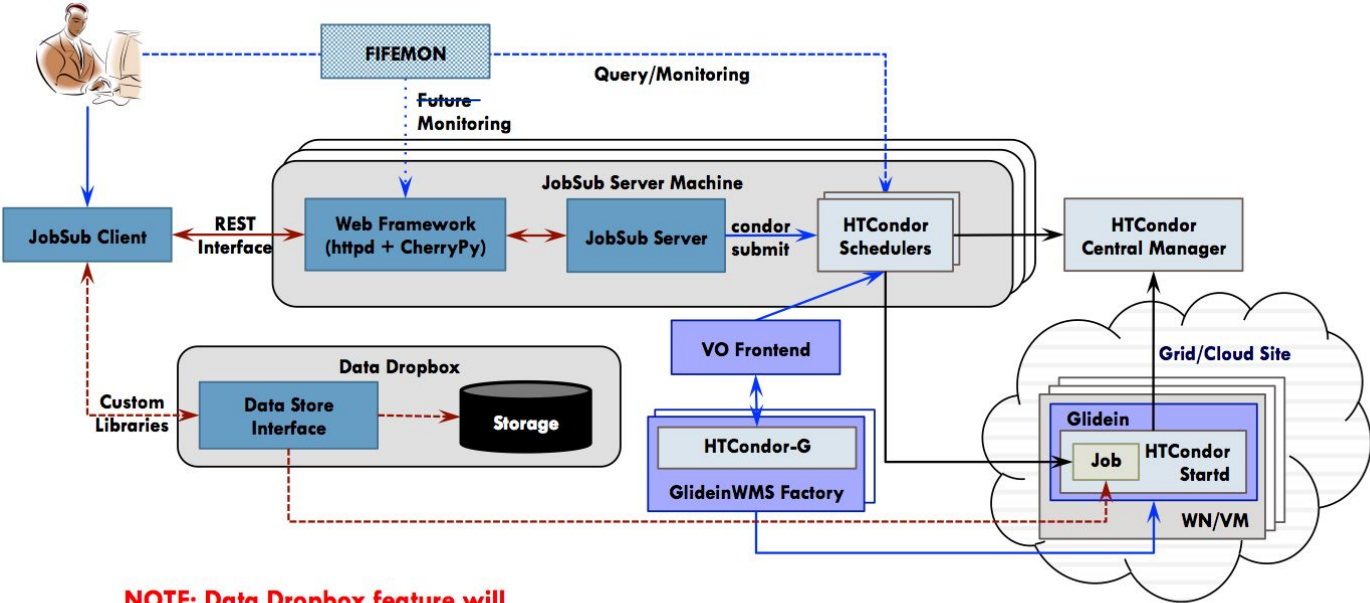
# Outline

- Background on the jobsub project and issues
- jobsub\_lite
- Tokens in jobsub\_lite
- Adoption of jobsub\_lite
- Wins and...opportunities
- Lessons learned so far

# Jobsub Project

- Many Intensity Frontier experiments at Fermilab had their own wrapper scripts written on top of HTCondor
- Fabric for Frontier Experiments (FIFE) project, Jobsub project was meant to
  - Unify wrappers/software stack, provide common job submission interface
  - Load balancing/HA, credential management, job log management among multiple schedds
- jobsub\_tools, then jobsub\_client/jobsub\_server (circa 2013)
- jobsub\_client generally installed on experiment submit nodes
- jobsub\_server, alongside HTCondor schedd run on separate machines (3 in production cluster)
- Interaction between two via REST API

# Jobsub the Old



**NOTE: Data Dropbox feature will be implemented in future releases.**

Original Image Source: [https://cdcv.s.fnal.gov/redmine/projects/fife/wiki/Introduction\\_to\\_FIFE\\_and\\_Component\\_Services#Jobsub](https://cdcv.s.fnal.gov/redmine/projects/fife/wiki/Introduction_to_FIFE_and_Component_Services#Jobsub)

# Problems with Old Jobsub

- Being too permissive with feature request acceptance (“Wouldn’t it be nice if jobsub did.....” ) led to
  - Lots of code customization for different VOs/experiments
  - → Large number of “gotchas”/accidental behavior
  - Too many ways to do the same set of operations (e.g. tarball upload)
- >21k lines of code (not including packaging scripts, tests, etc.)
- Supporting current feature set too difficult for available effort
  - Also complicates building new features and fixing bugs
- Used transition of OSG to SciToken auth and HTCondor dropping internal proxy auth to rewrite jobsub
  - Neatly avoids issues with proxies: e.g. have had instances of users accidentally deleting large swaths of data...

# Heeeeere's jobsub\_lite!

- New software for job submission and monitoring, built directly on top of Condor
- Tried to keep the most-used pieces of jobsub\_client, strip out unnecessary parts
- Client-only, installed on experiment submit nodes
- jobsub\_\* counterparts to condor\_\* commands (e.g. jobsub\_submit, jobsub\_q, etc)
- Currently, submits jobs with both proxy and token, but will be phasing out proxy gradually
- Remote submission to schedd

# What happens

- `jobsub_lite` takes user command,
  - Finds schedds
  - For submit: Converts user command to Condor submission file (Job Definition File), and uses Condor commands to remotely submit the job to schedd
  - For other commands: Gets credentials, converts user command to HTCondor command and runs it

```
jobsub_q -G fermilab 12345@jobsub01.fnal.gov
```



```
_condor_CREDD_HOST=jobsub01.fnal.gov  
/usr/bin/condor_q -global -schedd-constraint  
IsJobsubLite==True -name jobsub01.fnal.gov  
<formatting args> 12345
```

# jobsub\_lite and HTCCondor

- Idea is to keep jobsub\_lite....light
- Provide lightly-wrapped Condor executables (condor\_submit, condor\_q, etc.) on submit nodes
- Provide DAG submission through jobsub\_submit\_dag (different DAG format called *dagnabbit*, which we translate to Condor DAG format)
- Most users: jobsub commands
- Advanced use-cases: Condor commands
- To help with advanced use-cases, jobsub\_submit has option to just create Condor Job Definition File to use with condor\_submit

# Lightweight condor wrappers

- Idea from CMS LHC Physics Center (LPC) deployment at Fermilab
- Parse a couple of arguments, get credentials, find correct collector/schedd, then hand the work over to Condor

```
condor_q -G fermilab 12345@jobsub01.fnal.gov
```

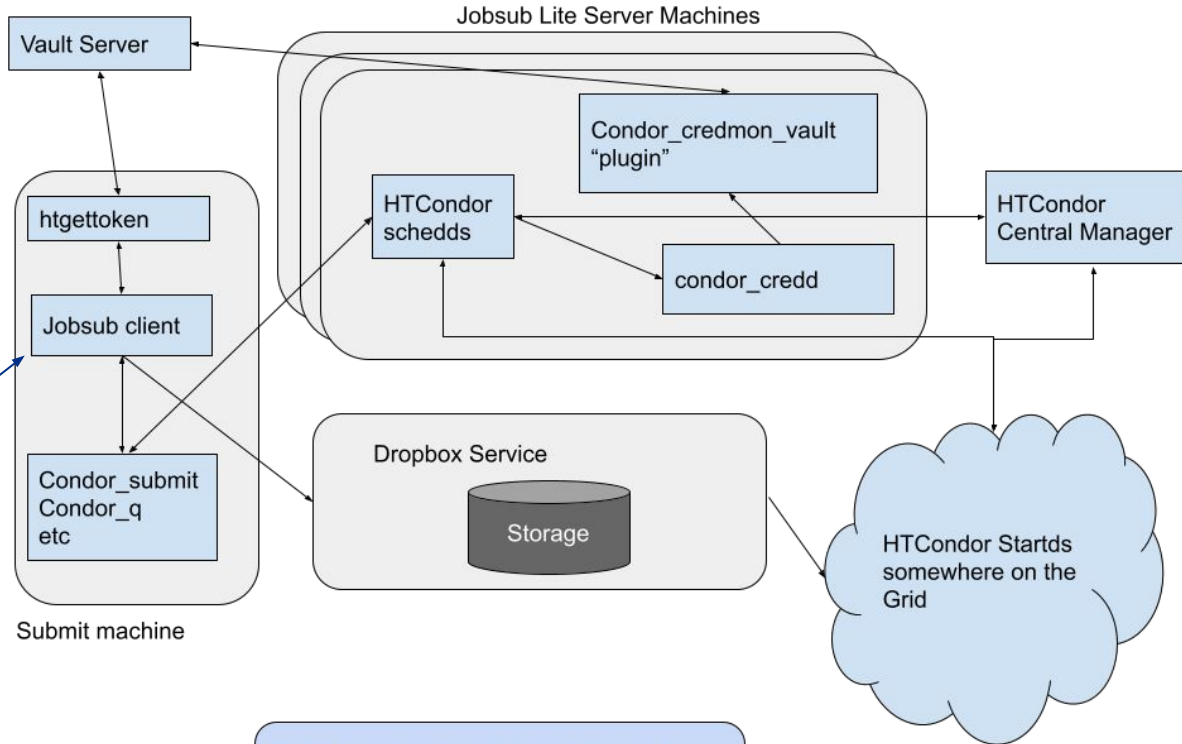


```
_condor_CREDD_HOST=jobsub01.fnal.gov  
/usr/bin/condor_q -global -schedd-constraint  
IsJobsubLite==True -name jobsub01.fnal.gov  
<formatting args> 12345
```

# Infrastructure/Condor Versions

- Schedds:
  - Development: Were using Condor 9 on schedds
  - Production: Condor 10.0.3
  - Currently deploying with shared schedds, but plan to transition to one schedd per large experiment, and a couple of shared schedds (use a SupportedVOList classad attribute on schedd)
- Submit nodes:
  - Most running 9.0.17. (To be upgraded to 10 soon)

# jobsub\_lite Infrastructure with Tokens



This is all that jobsub\_lite developers have to maintain!

Active/Large experiments get their own schedd and can't be brought down by other experiments errant users.

Original Image Credit: J. Boyd

# Tokens and Authentication

# Token Authentication Flow

We use `htgettoken` to obtain vault and bearer tokens

Steps:

1. Kerberos ticket used to authenticate to Hashicorp Vault
2. Vault contacts token issuer - CILogon
3. First time, token issuer has user authenticate in browser
4. Refresh token stored in vault, Vault and Access token downloaded to user node

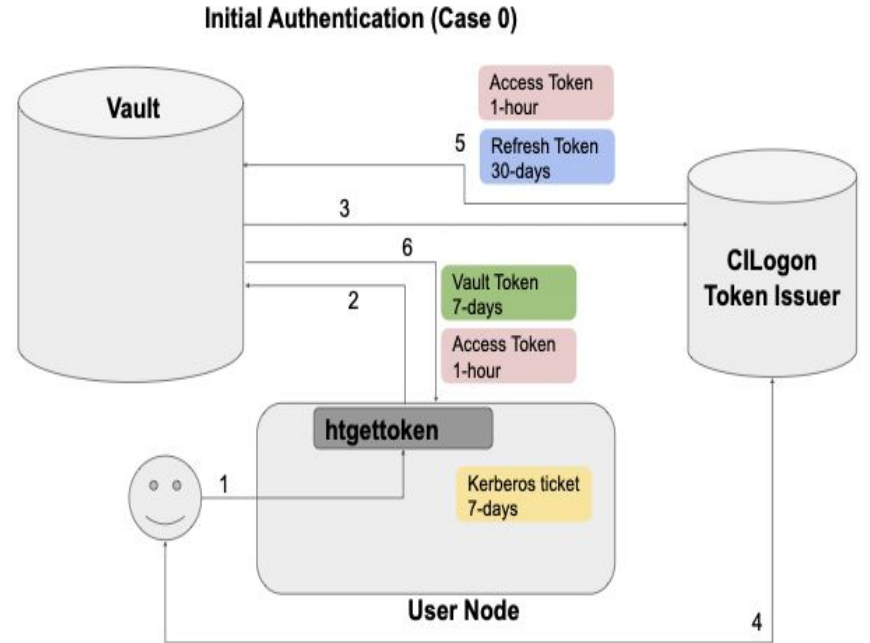


Image Credit: M. Altunay and D. Dykstra

# jobsub\_lite and Tokens

- Fine-grained access control via SciTokens!
  - DUNE, for example, has many different capability sets (different sets of pre-defined token scopes) for different sets of users
  - Outsource these decisions to VOs/experiments
- jobsub commands obtain Access token via `htgettoken` in case it's needed
- Leverage Condor to do token exchange at submission time: default tokens scopes include *storage.read*, *storage.create*, but not *storage.modify* → Users have to specifically request tokens with *storage.modify*

# Robot Tokens and the Managed Tokens Service

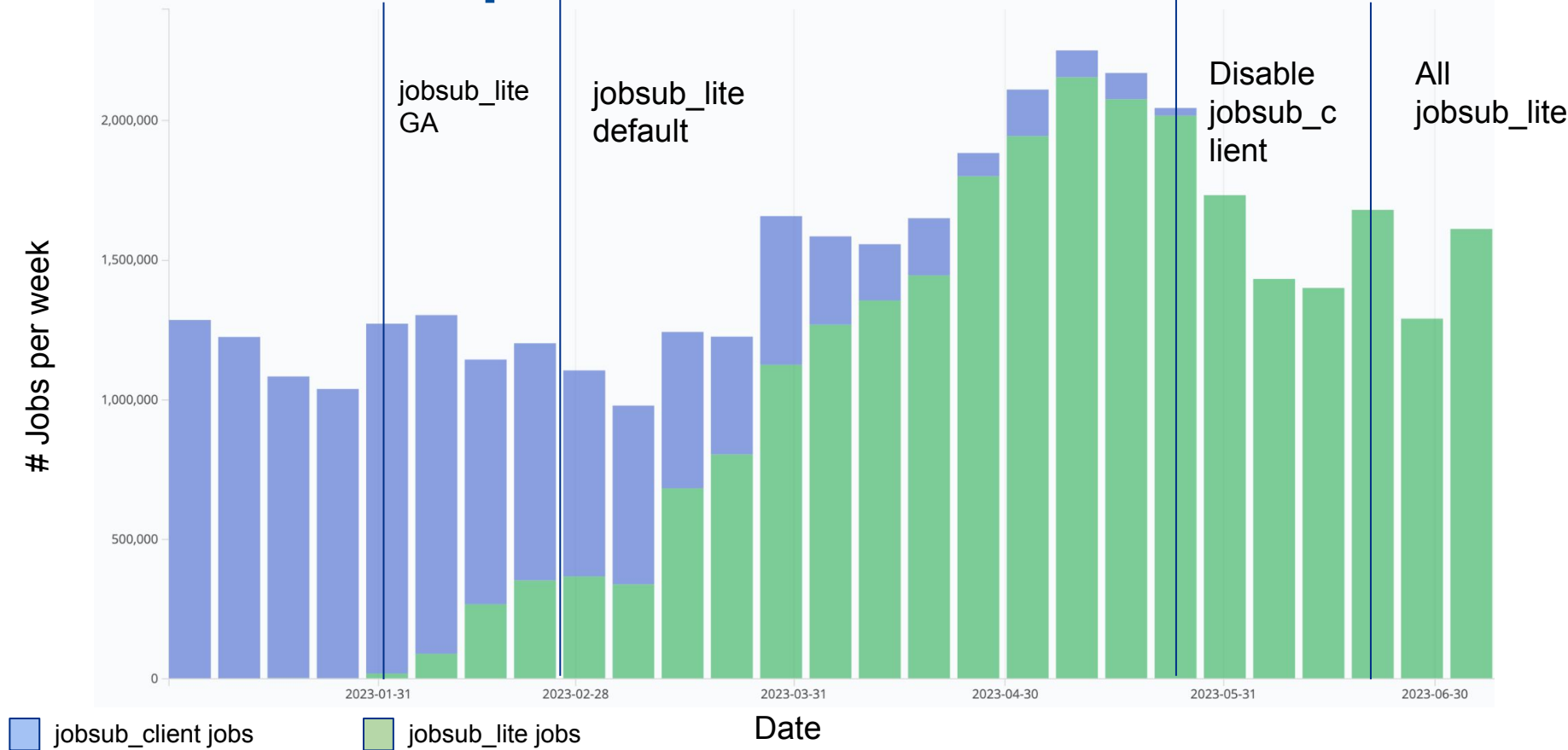
- Previously, Managed Proxies service periodically refreshed VOMS proxies on experiment interactive nodes for production activity
- Stakeholders requested same for tokens
- `htgettoken` supports use of robot kerberos creds to obtain vault tokens
- Leverage this capability for new **Managed Tokens Service** (written in Go)
- Push production vault tokens to interactive nodes, keep them refreshed:
  - Obtain kerberos credentials
  - `condor_vault_storer` for each schedd
  - `rsync` vault token to appropriate submit nodes

# Managed Tokens Service (2)

- Managed Tokens Service users should *never* have to authenticate in CILogon
  - “Onboarding” = operator running `condor_vault_storer` manually for all schedds, and authenticating (Managed Tokens Service has utility to do that)
  - The pushed vault token is used to obtain bearer token on submit node
- User steps:
  - Set `--credkey` in `HTGETTOKENOPTS` in environment
  - In `jobsub_*` command, pass `--role=production` environment to set right service for `condor_vault_storer/mappings`
- Has been running in production since November 2022 with very few issues

# Adoption/Lessons Learned

# Jobsub\_lite adoption



jobsub\_client jobs

jobsub\_lite jobs

Date



# The Good News (For Users)

- Users have full access to Condor commands, Condor JDFs
  - No more passing through constraints through jobsub
  - Don't have to wrap Condor DAG commands if you don't want
- With more focused interface, easier to get new users started on jobsub\_lite
- Horizontal scaling (adding more schedds) much easier
  - Started with one schedd in production, have now scaled to four
  - Can do rolling upgrades to various components of system

# The Good News (For Everyone)

- Fine-grained access control via SciTokens!
- Less code (6906 lines INCLUDING tests and templates) = less to maintain, easier to add features/fix bugs...but
- Have had to be strict about feature requests
  - Decreases our support load
  - Users should be using Condor (for anything beyond basics)!

# The Bad News - Infrastructure Issues

- If credd has issues talking to vault, sometimes not enough info in the logs
  - If token expired in credmon, in certain cases, jobs would just fail to start, with no user notification until they went held for SHADOW exceptions (Fixed in 10.0.3)
- Duty Cycle Issues on Schedds:
  - Above token issues → Tons of shadow starts (Fixed in 10.0.3)
  - With no jobsub\_server throttling user submissions, had duty cycle issues when under heavy load →
    - Tweak MAX\_JOBS\_PER\_SUBMISSION to limit cluster size
    - Tweak CURB\_MATCHMAKING to throttle job matchmaking
- HTCCondor team SUPER helpful in assisting with/fixing these issues

# The Bad News - Everything Else

- Monitoring: FIFEMon didn't fully support jobsub\_lite until phase 2, which led to resistance to adoption
  - Remote submission with -spool → jobs didn't leave queue for users to get logs PLUS our monitoring looks at condor\_history
- Resistance to adopting tokens - “But why do I have to do this when the old way just works?!”
- Have had to be strict about feature requests, which didn't make some happy → Users should be using HTCondor!

# Lessons Learned

- Inertia is real...
- Better to delay a go-live of this magnitude if monitoring, fetching of logs, etc., not ready
- Robust monitoring of old system allowed us to choose which features to implement in `jobsub_lite`
  - We looked at 6 months of command-line options used with `jobsub_client` to pick `jobsub_lite` flags
- Running both systems in parallel for a time was helpful in transition in case there were issues
  - Classad mechanism for schedds allowed mixed cluster of `jobsub_client/jobsub_lite` submit nodes and `jobsub_server/condor` schedds (control which submit nodes submit to which schedds)
  - Users could fall back to old system with a single flag, giving us breathing room to fix bugs
  - Could migrate schedds one or two at a time

## Lessons Learned (2)

- Tokens  $\neq$  Proxies  $\rightarrow$  Training users is key
- Heterogeneous environments  $\rightarrow$  Need stakeholders from ALL parties to test, not just the willing...
  - But if you only have one test schedd, don't ask everyone to test at once
- Those who were willing to test early had a much easier transition, so that is KEY.

# Future work

- Phasing out use of X509 proxies in job submission:
  - Next minor release: Create opt-out flag for obtaining proxy
  - When we receive approval from experiments, convert to opt-in flag (probably in a couple of years)
- Bugfixes
- Usability features
- Test with EL9 (Currently running on SL7 machines)
- Shift to maintenance mode (try not to add any major features)

# References/Links

- Jobsub\_lite git repository/documentation:  
[https://github.com/fermitools/jobsub\\_lite](https://github.com/fermitools/jobsub_lite)
- Managed Tokens Service git repository:  
<https://github.com/shreyb/managed-tokens>
- Original jobsub paper: Dennis Box 2014 J. Phys.: Conf. Ser. 513 032010  
DOI 10.1088/1742-6596/513/3/032010  
<https://iopscience.iop.org/article/10.1088/1742-6596/513/3/032010>

# Thank you!

The jobsub project team:

Shreyas Bhat, Joe Boyd, Vito Di Benedetto, Lisa Goodenough, Marc Mengel, Nick Peregonow, Kevin Retzke

This work was produced by Fermi Research Alliance, LLC under Contract No. DE-AC02-07CH11359 with the U.S. Department of Energy. Publisher acknowledges the U.S. Government license to provide public access under the DOE Public Access Plan: [DOE Public Access Plan](#)

# Backup Slides

# Token Authentication Flow

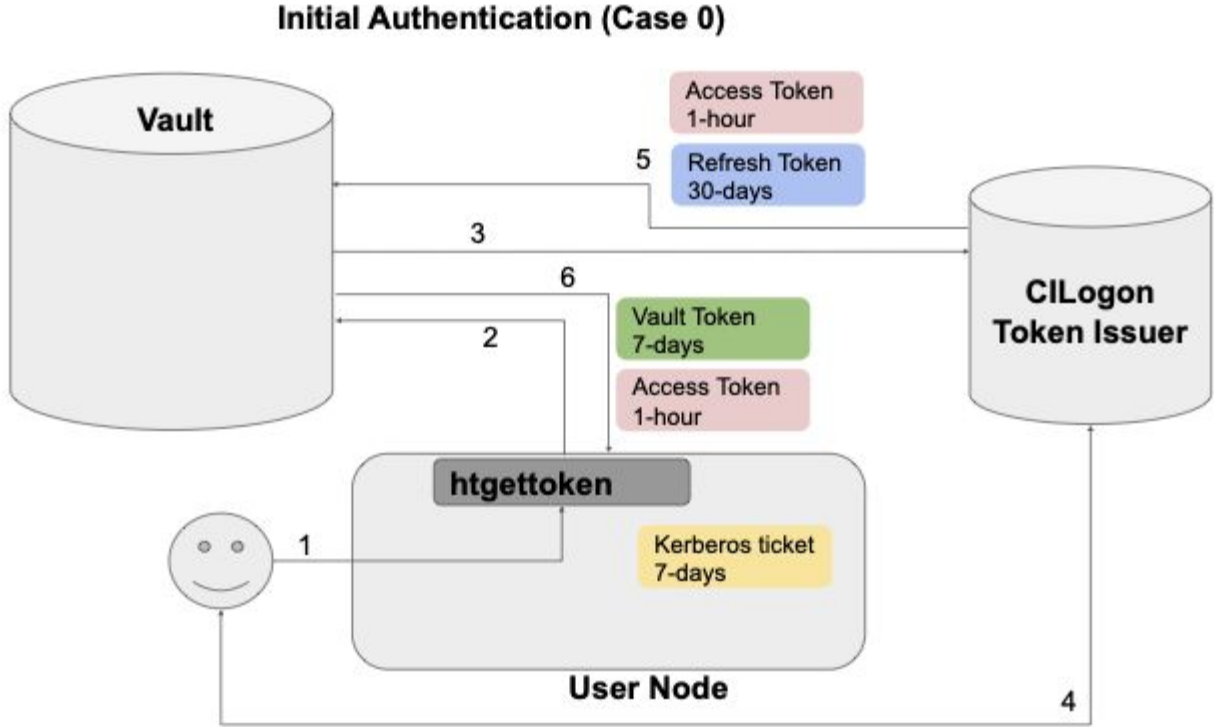


Image Credit: M. Altunay and D. Dykstra

# First-time authentication

- Authentication happens for most grid operations - now X509 Proxy, soon tokens (jobsub\_lite, ifdhc commands)
- Absence of vault or refresh token → Authenticate with CILogon

Attempting OIDC authentication with `https://htvaultprod.fnal.gov:8200`

Complete the authentication at:

`https://cilogon.org/device/?user_code=_redacted_user_code`

No web open command defined, please copy/paste the above to any web browser

Waiting for response in web browser

- Will need to copy/paste that link into browser

# First-time authentication, continued



Consent to Attribute Release ▼

[htvaultst-fermilab-vault](#) requests access to the following information. If you do not approve this request, do not proceed.

- User Code: Z4H-CHD-WPJ
- Your CILogon user identifier
- Your name
- Your email address
- Your username and affiliation from your identity provider

Select an Identity Provider

Fermi National Accelerator Laboratory ? ←

Remember this selection ?

←

By selecting "Log On", you agree to the [privacy policy](#).

- Select "Fermi National Accelerator Laboratory"
- Then click "Log On"
- Log in with Services Credentials

For questions about this site, please see the FAQs or send email to [help@cilogon.org](mailto:help@cilogon.org).  
Know your responsibilities for using the CILogon Service.  
See acknowledgements of support for this site.

# Further Notes about Authentication

- After initial authentication, as long as you use token-enabled grid tools for the same experiment/role at least every 30 days, you should *not* have to reauthenticate
- This is because refresh token (kept in vault) expires after 30 days of inactivity
- Tokens downloaded to user machine:
  - Vault Token: Used to authenticate to vault
  - Access (or Bearer) Token: SciToken (JWT) that is actually used for grid operations
- More information on SciTokens: <https://scitokens.org/>

# Submit and Manage Simple Job

# jobsub\_submit

- Much easier than before. Just login, and jobsub\_submit

```
$ jobsub_submit -G fermilab file:///usr/bin/printenv
Attempting to get token from https://fermicloud543.fnal.gov:8200 ... failed
Attempting kerberos auth with https://fermicloud543.fnal.gov:8200 ... succeeded
Attempting to get token from https://fermicloud543.fnal.gov:8200 ... failed
Attempting OIDC authentication with https://fermicloud543.fnal.gov:8200
```

Complete the authentication at:

```
https://cilogon.org/device/?user_code=<code>
```

No web open command defined, please copy/paste the above to any web browser

Waiting for response in web browser

```
Storing vault token in /tmp/vt_u10610
```

```
Storing bearer token in /tmp/bt_token_fermilab_Analysis_10610
```

```
Submitting job(s).
```

```
1 job(s) submitted to cluster 57106734.
```

```
Use job id 57106734.0@jobsub01.fnal.gov to retrieve output
```

# jobsub\_submit, continued

- Like before, -G/--group is required to submit job (and run all jobsub executables)
- Group dictates which token issuer is used to get a bearer token
- FOR NOW, jobsub will obtain a bearer token and VOMS-proxy (valid for ~one week) and send these to the job
  - Future - no VOMS proxy

# Manage jobs

- `jobsub_q`, `jobsub_hold`, `jobsub_release`, `jobsub_rm`, etc.  
written as lightweight wrappers around `condor_*` commands
- Tried to keep backward-compatibility
- Examples on following slides

# jobsub\_q

```
$ jobsub_q -G fermilab
```

JOBSUBJOBID COMMAND	OWNER	SUBMITTED	RUNTIME	ST	PRIO	SIZE
57106962.0@jobsub01.fnal.gov simple.sh	sbhat	11/30 14:59	0+06:00:10	C	0	1953.1
57106973.0@jobsub01.fnal.gov simple.sh	sbhat	12/01 13:48	0+06:00:27	R	0	0.0

```
$ jobsub_q -G fermilab 57106973.0@jobsub01.fnal.gov
```

JOBSUBJOBID COMMAND	OWNER	SUBMITTED	RUNTIME	ST	PRIO	SIZE
57106973.0@jobsub01.fnal.gov simple.sh	sbhat	12/01 13:48	0+06:00:27	R	0	0.0

# jobsub\_hold

```
$ jobsub_hold -G fermilab 57106973.0@jobsub01.fnal.gov
```

```
Job 57106973.0 held
```

```
$ jobsub_q -G fermilab 57106973.0@jobsub01.fnal.gov
```

JOBSUBJOBID	OWNER	SUBMITTED	RUNTIME	ST	PRIO
57106973.0@jobsub01.fnal.gov	sbhat	12/01 13:48	0+06:00:27		H
0 0.0 simple.sh					

# jobsub\_release

```
$ jobsub_release -G fermilab 57106973.0@jobsub01.fnal.gov
```

```
Job 57106973.0 released
```

```
$ jobsub_q -G fermilab 57106973.0@jobsub01.fnal.gov
```

JOBSUBJOBID	OWNER	SUBMITTED	RUNTIME	ST	PRIO
57106973.0@jobsub01.fnal.gov	sbhat	12/01 13:48	0+06:00:27	I	
0 0.0 simple.sh					

# jobsub\_rm

```
$ jobsub_rm -G fermilab 57106973.0@jobsub01.fnal.gov
```

```
Job 57106973.0 marked for removal
```

```
$ jobsub_q -G fermilab 57106973.0@jobsub01.fnal.gov
```

JOBSUBJOBID	OWNER	SUBMITTED	RUNTIME	ST	PRIO
SIZE	COMMAND				

# Singularity/Apptainer

- By default, jobs run in `fnal-wn-sl7:latest` singularity image
- Opt out by either:
  - Specifying singularity image:  
“`--singularity-image=/path/to/singularity/image`”
  - Passing “`--no-singularity`”: Site-dependent. To truly get outside a singularity container, pass `--no-singularity` and request a site that you know does not run singularity containers
- Have `--apptainer-image` and `--no-apptainer` flags

# DAGs

# Submit DAGs

- jobsub\_lite supports dagnabbit syntax to describe DAGs
- Example file mywork.dagnabbit:

```
<serial>
jobsub_submit  --mail_on_error $SUBMIT_FLAGS file://jobA.sh
jobsub_submit  --mail_on_error $SUBMIT_FLAGS file://jobB.sh
</serial>
<parallel>
jobsub_submit  --mail_on_error $SUBMIT_FLAGS file://jobC.sh
jobsub_submit  --mail_on_error $SUBMIT_FLAGS file://jobD.sh
</parallel>
<serial>
jobsub_submit  --mail_on_error $SUBMIT_FLAGS file://jobE.sh
</serial>
```

# Submit DAGs, continued

- Submit DAG:

```
export SUBMIT_FLAGS="-G fermilab"  
jobsub_submit $SUBMIT_FLAGS --dag file://mywork.dagnabbit
```

# Tarfiles

# -f and --tar-file-name

- All use Rapid Code Distribution Service (RCDS) via CVMFS by default
- *--tar-file-name*: specify TAR\_FILE or DIRECTORY to be transferred to worker node
  - TAR\_FILE will be accessible to the user job on the worker node via the environment variable \$INPUT\_TAR\_FILE
  - The unpacked contents will be in the same directory as \$INPUT\_TAR\_FILE
  - Successive --tar\_file\_name options will be in \$INPUT\_TAR\_FILE\_1, \$INPUT\_TAR\_FILE\_2, etc.
  - Use with `dropbox://` for pre-made tarfile, `tardir://` to specify directory to be tarred up

## -f and --tar-file-name (2)

- *-f*: Copy INPUT\_FILE file at runtime
- INPUT\_FILE copied to directory \$CONDOR\_DIR\_INPUT on the execution node.
- Example :  
-f /grid/data/minerva/my/input/file.xxx  
copied to \$CONDOR\_DIR\_INPUT/file.xxx
- Specify as many -f INPUT\_FILE\_1 -f INPUT\_FILE\_2 args as you need.
- To copy file at submission time use -f dropbox://INPUT\_FILE to copy the file

# Condor commands

# Using Condor commands

- One major change with `jobsub_lite` is that users have access to condor commands
- We recommend users use the `jobsub_lite`-wrapped condor commands, as they handle authentication, but using HTCondor-provided condor commands is an option

# Production Jobs and Managed Tokens

# Roles in token-world

- No VOMS-server signing proxies in the token-world
- Role = entry in “wlcg.groups” entry of token
  - This entry is mapped to “capability set” in LDAP/FERRY, which defines your “scopes” entry
  - “scopes” controls authorization
- Production tokens usually have access to read/write to ALL of an experiment’s dCache area, but this is configurable

# Managed Tokens

- New service to push production vault tokens to interactive nodes, keep them refreshed
- Production users should *never* have to authenticate in CILogon
  - The pushed vault token is used to obtain bearer token
- Set in environment:

```
export
```

```
HTGETTOKENOPTS="--credkey=<account>/managedtokens/fifeutilgpvm01.fnal.gov"
```

```
export X509_USER_PROXY=/path/to/production/proxy
```

- Then, in `jobsub_*` command, pass `--role=production` (note lower-case "p")

# Deployment Overview

- November 2022: Iron out deployment details with mu2e
- December 2022:
  - Deploy to experiment “test” interactive nodes, get feedback
  - Announce to general users the go-live date
- January 2023: Run two more demos of jobsub\_lite
- **February 1, 2023:**
  - **Go-live of jobsub\_lite (see next slide)**
  - Plan changed from before
- June 21, 2023: Turn off jobsub servers → jobsub\_client will no longer work

# Moving and Deleting Files from Grid Jobs

- For a bearer token to authorize a user to move or remove a file, it must have the **storage.modify** scope on the path containing the file
- Due to security concerns, this is *not* granted to users by default; it must be requested on the jobsub\_submit command line
- This is done with the “--need-storage-modify <path>” flag
- jobsub\_lite will evaluate whether the storage.modify request is valid
- Similarly, “--need-scope <scope>” will request a scope be added to token

# Moving and Deleting Files from Grid Jobs (2)

## Examples:

- Request storage.modify on /pnfs/mu2e/scratch/users/username

```
$ jobsub_submit -G mu2e --need-storage-modify /mu2e/scratch/users/username file:///bin/true
```

# Go live!

- Go live on February 1, 2023
- Phased go-live:
  - Phase 1: Make jobsub\_lite available
  - Phase 2: Make jobsub\_lite default job-submission tool
  - Phase 3: Turn off job submission from old jobsub
  - Phase 4: Turn off old jobsub infrastructure
- Phase 4 ended June 21, 2023

# User Training/Support Efforts

- 4 training sessions for jobsub\_lite
  - 1 for power users
  - 3 for anyone
  - First one well-attended
  - Second and third, not so much
  - Fourth was better
- During Phase 1-2, “If there’s a problem either with tokens or job submission, we’ll help users directly”
- Phase 3-4: Problems need to be brought up with VO power-users first