# CHEP 2007

# glideinWMS
# -
# A generic pilot-based
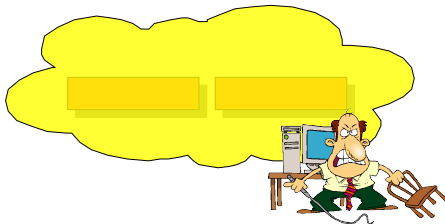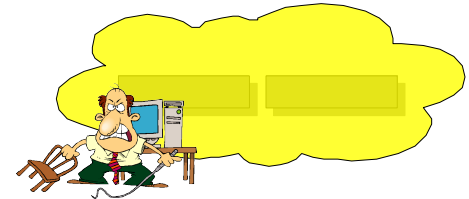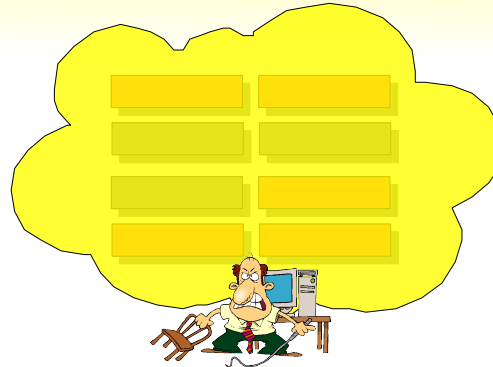# Workload Management System

by Igor Sfiligoi (FNAL)

# Outline

- What is glideinWMS?

- How does it work?

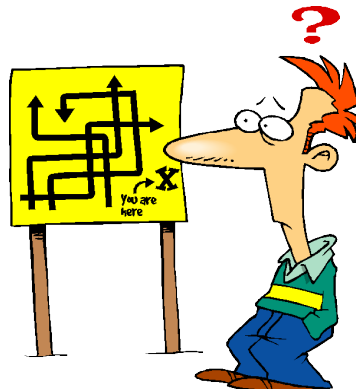- How does it perform?

- Monitoring

- Conclusions

# What is glideinWMS?
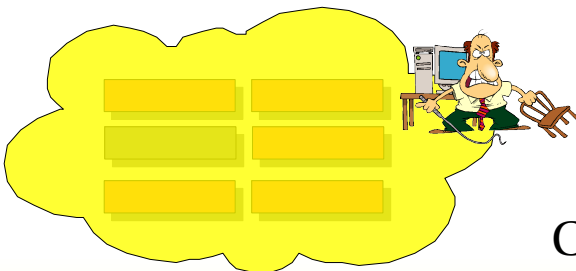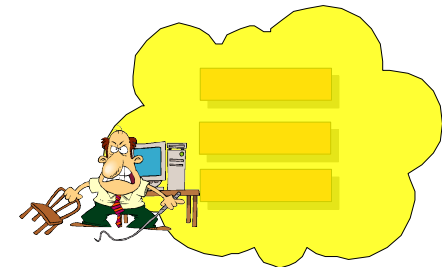
- A Condor glidein-based Workload Management System

- Developed by CMS for CMS, but generic enough to be used by other groups, too
  - A generalization of the CDF glidekeeper

- Available at: http://home.fnal.gov/~sfiligoi/glideinWMS/

# Why do we need a WMS?

"The Grid" is really a sum of hundreds of independent Grid sites.

Choosing where to try to run the jobs is not a trivial task

# What is Condor? (1)

- A widely used batch system
- Based on a fully distributed architecture



Collector

Negotiator

Have jobs,
need workers

Schedd

Schedd

Have workers,
need jobs

Starter

Starter

Starter

# What is Condor? (2)

- A widely used batch system
- Based on a fully distributed architecture

# What is a glidein? [1]

- Just a regular starter

- Submitted as a Grid job



Collector

Negotiator

Have jobs, need workers

Have worker, need job

Schedd

Schedd

Grid batch slot

Starter

The Grid

Grid batch slot

Other Grid Job

# What is a glidein? (2)

- Just a regular starter
- Submitted as a Grid job



Collector

Negotiator

Expect a job from s2

Schedd

Grid batch slot

Starter

Send job to wg

Job

Schedd

The Grid

Grid batch slot

Other Grid Job

# What is a glidein? (3)

- Just a regular starter
- Submitted as a Grid job

# What else can a glidein do?

- Make sanity checks before fetching any job

- Discover and publish batch slot characteristics:
  – OS version

  – CPU model, available RAM and disk

  – Availability of certain software

- Importing VO specific software

- Prepare the environment for the user jobs

  – Possibly putting the VO software in the path

- etc.

# Why using glideins? (1)

- ## For people already using Condor
  - ### An easy way to extend the pool
    - Or to create one from scratch
  - ### Can hide all the grid stuff from user jobs
    - Can even run standard universe jobs on the Grid!
- ## For people just wanting to use the Grid
  (even if not Condor fans)
  - ### Protect user jobs from many obvious errors
    - A dead glidein will not pull a user job
  - ### Simplifies resource selection
    - A glidein can detect what is available on the worker and user jobs get sent only to complying workers
    - No guessing involved, job sent after resource acquired

# Why using glideins? (2)

- Get all the advantages of a local batch system
  - Locally set priorities between different users
    - Including group quotas
    - Or even priorities between jobs of the same user
  - Reliable, real time monitoring
  - Reliable file transfer
    - Full file encryption supported, too
- **While still running on the Grid!**

Any weak points will be presented at the end

# How do I submit a glidein?

- Condor provides condor_glidein
  - Simple command line tool
  - **Useful when you have just a few jobs**
    - Will submit a single glidein per invocation
- Install a glideinWMS instance
  - Needs more resources and some initial effort to set it up
  - **Setup once, glideins will be launched as needed**
    - Will look for jobs that need resources
    - Submit glideins as needed to sites that seem to match at least an idle job
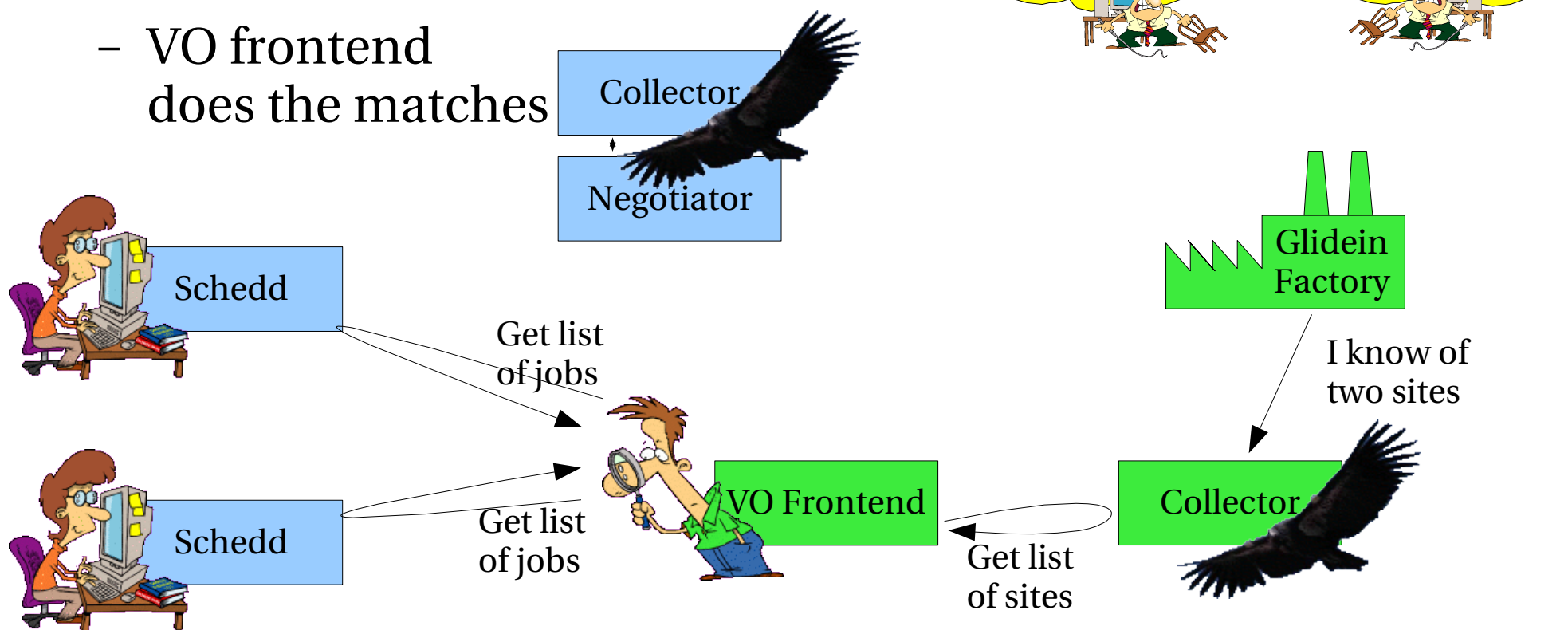
# glideinWMS

How does it work?

# glideinWMS overview [1]
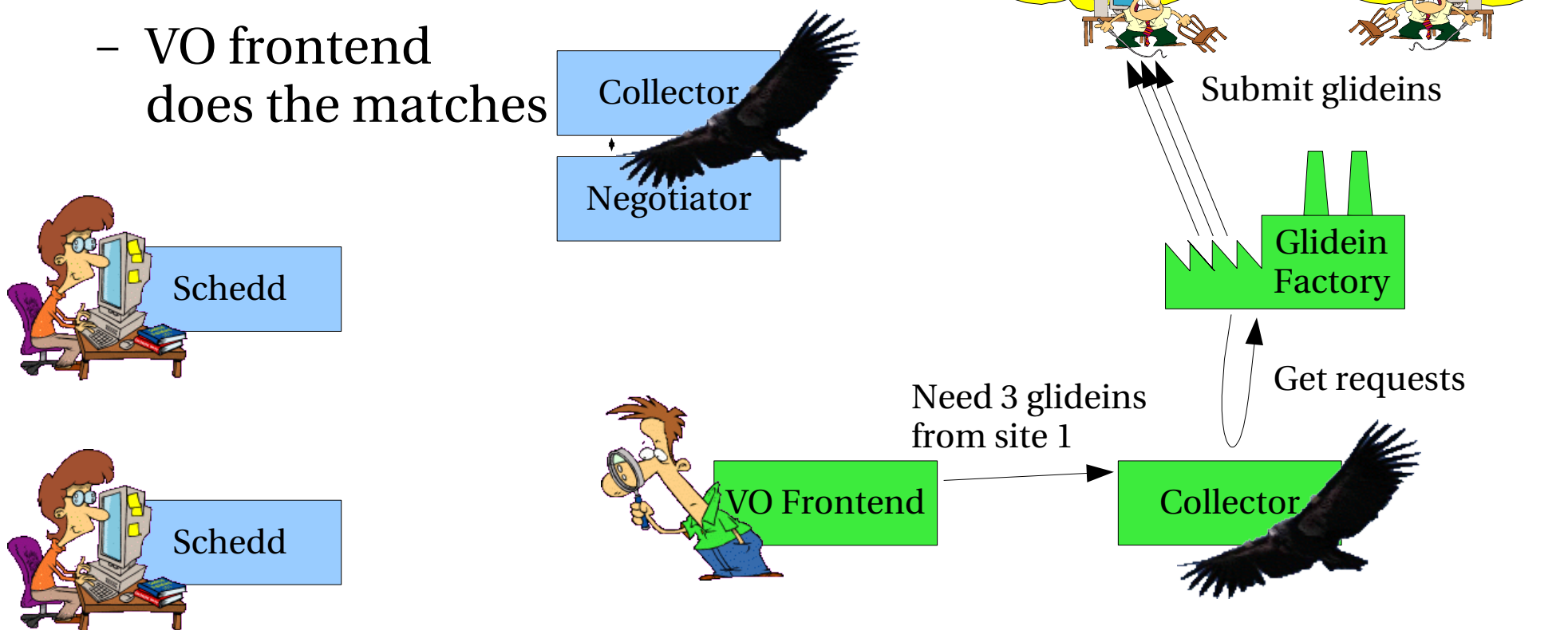
- A thin layer on top of Condor
  - VO frontend does the matches

A Grid Site

A Grid Site

Collector

Negotiator

Schedd

Schedd

Get list of jobs

Get list of jobs

VO Frontend

Get list of sites

Glidein Factory

I know of two sites

Collector

# glideinWMS overview (2)

- A thin layer on top of Condor
  - VO frontend does the matches



A Grid Site

A Grid Site

Submit glideins

Collector

Negotiator

Schedd

Schedd

Glidein Factory

Get requests

Need 3 glideins from site 1

VO Frontend

Collector

More details at http://cd-docdb.fnal.gov/cgi-bin/ShowDocument?docid=2048

# glideinWMS overview (3)

- A thin layer on top of Condor

Job

Starter
A Grid Site

A Grid Site

Have worker, need job

Collector

Submit glideins

Negotiator

Send job to wg

Glidein Factory

Schedd

Have jobs, need workers

VO Frontend

Collector

Schedd

# glideinWMS details[1]

- Matchmaking done on two levels
  - **VO frontend** matches glideins to sites that claim to support at least one job waiting in the queue
  - The **condor negotiator** matches glidein starters to the jobs waiting in the queue
- The condor negotiator has the final word
  - If a site was lying about its capabilities, the starter will not be matched and will exit within minutes
  - The job that is sent to a starter might not be the one for which the glidein was submitted for

# glideinWMS details(2)

- The WMS logic is to keep constant pressure on the Grid sites
  - As long as there are waiting jobs that could be run on a site, it tries to keep a steady number of idle glideins in the site queues
- The VO frontend drives the WMS
  - Deciding how much pressure to put on different sites
  - The glidein factories will submit the glideins, following the orders from the VO frontend

# glideinWMS details(3)

- Communication between processes based on Condor ClassAds
  - For each site, a Glidein Factory publishes:
    - CE attributes
    - list of parameters it accepts
  - Each VO Frontend replies a ClassAd containing:
    - The target site
    - VO parameters (a subset of the above)
    - Number of idle glideins to keep in the queue
- Using a standard Condor collector

# glideinWMS details(4)

- Condor-G used for glidein submission

- The list of sites a factory serves is a configuration parameter
  - Can be set manually, fine tuning each and every site characteristics
  - Easy to script
    - Just a standard XML file
    - The installation script can use the CRONUS information system, and ReSS information system will be added soon
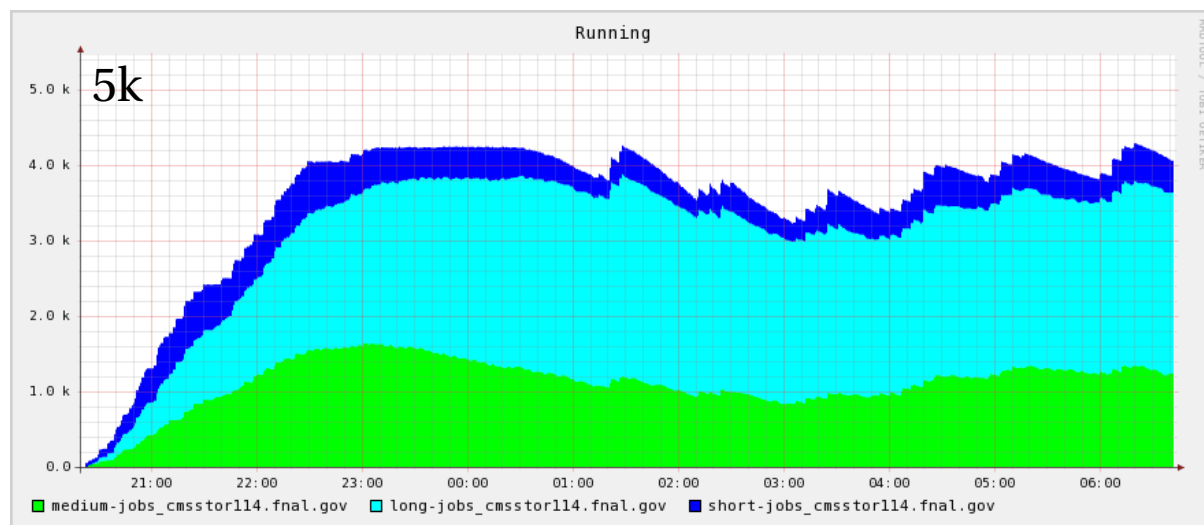  - Or can be paired with a Condor-G matchmaker, like ReSS and CRONUS

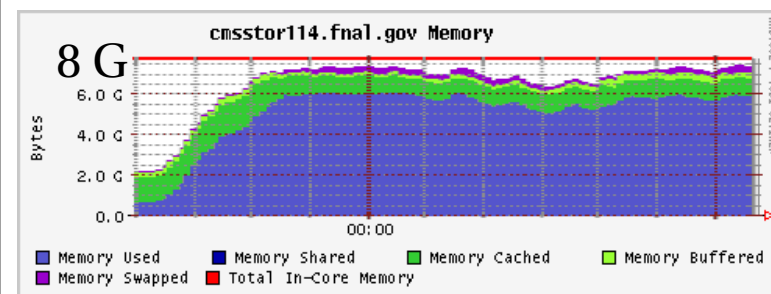# glideinWMS

## How does it perform?

# Do glideins scale? (1)

- Synthetic tests, using a single submit machine, scaled well with ~4000 running jobs
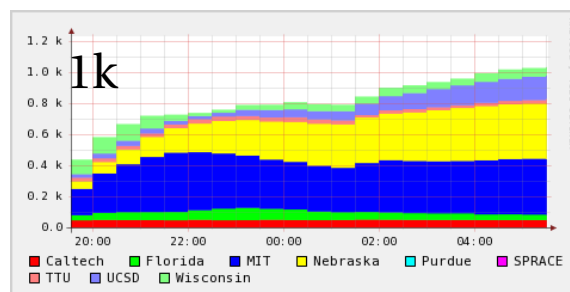
  – Memory a major limiting factor
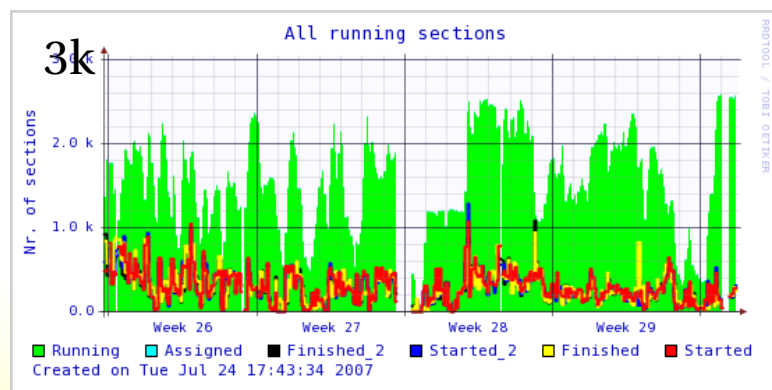
See also
Talk #216

Ignore
the colors

- Further scalability can be obtained by using multiple submission machines
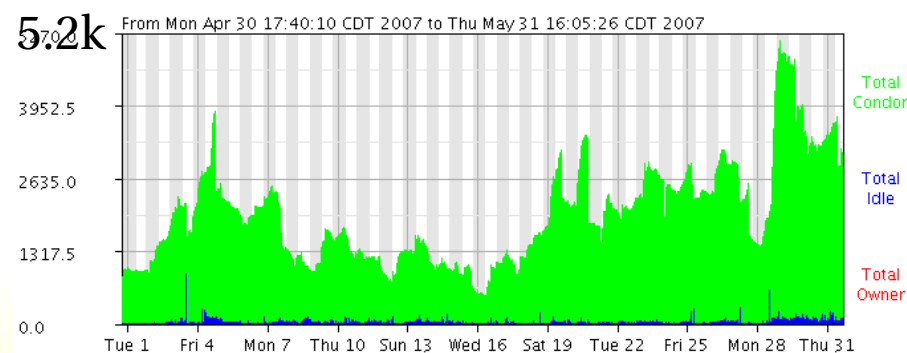
# Do glideins scale? (2)

- glideinWMS-based CMS MC production up to 1k jobs in parallel
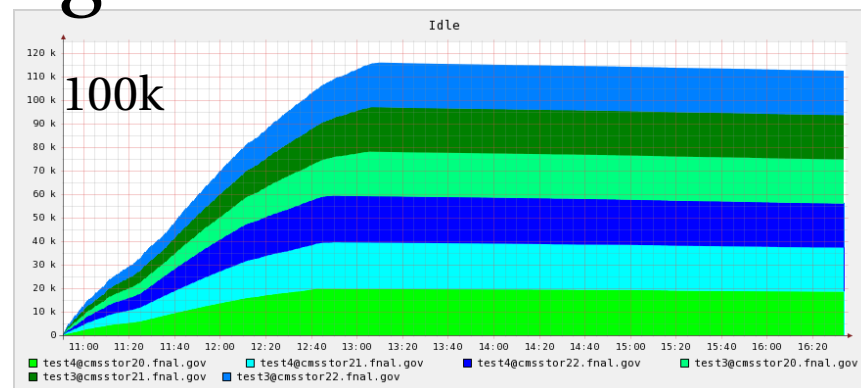


- CDF GlideCAF up to 2.5k
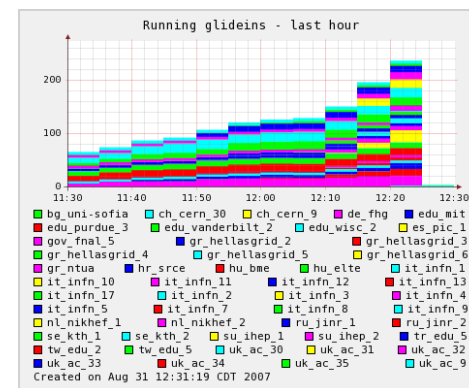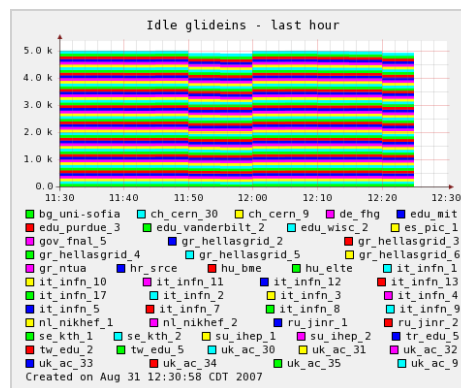


- ATLAS Cronus up to 5k

# Does glideinWMS scale?

- Synthetic tests with single VM frontend and single glidein factory

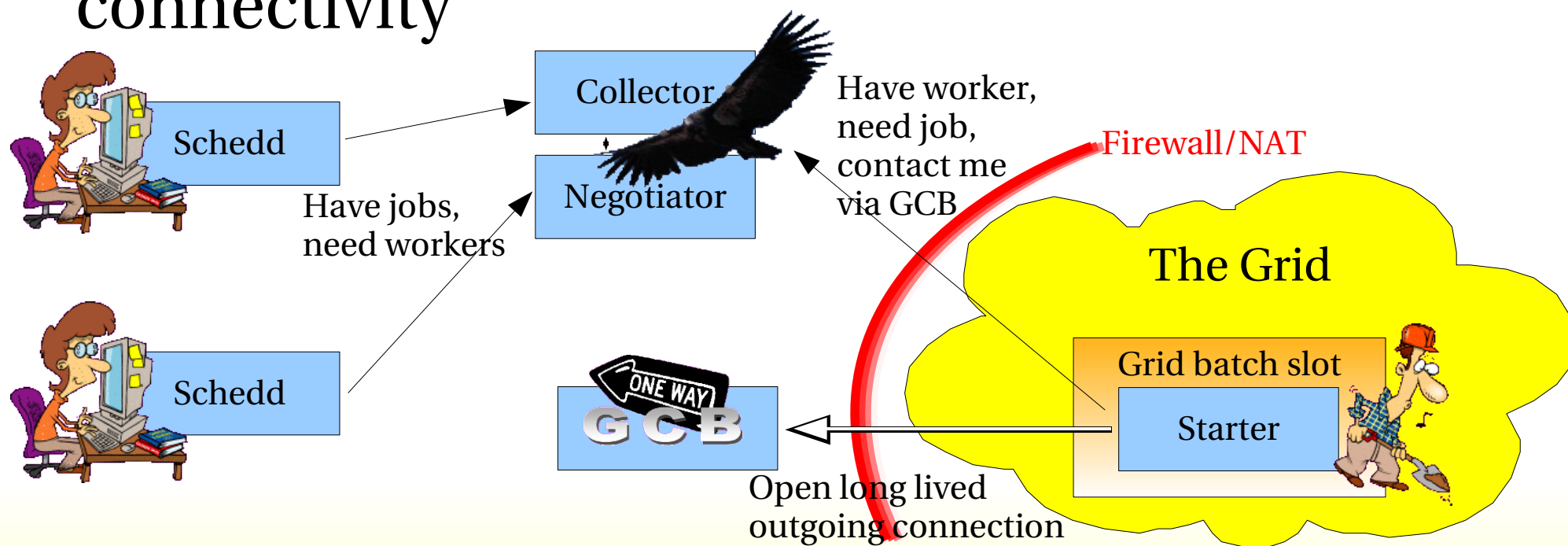  - 6 submission points with 100k queued jobs

  

  - 50 grid sites

  

- Further scalability by using multiple VO frontends and multiple glidein factories
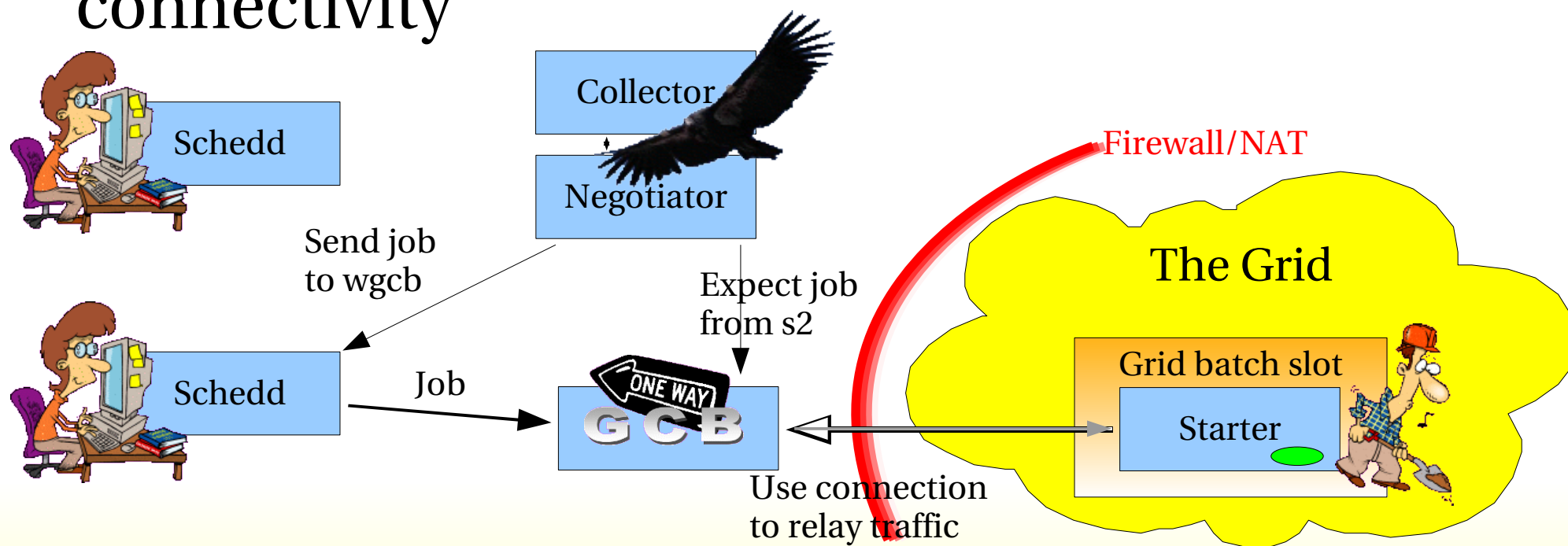
# Do glideins really work over WAN? [1]

- Yes, but it needs GCB to work
  - A Condor proxy server
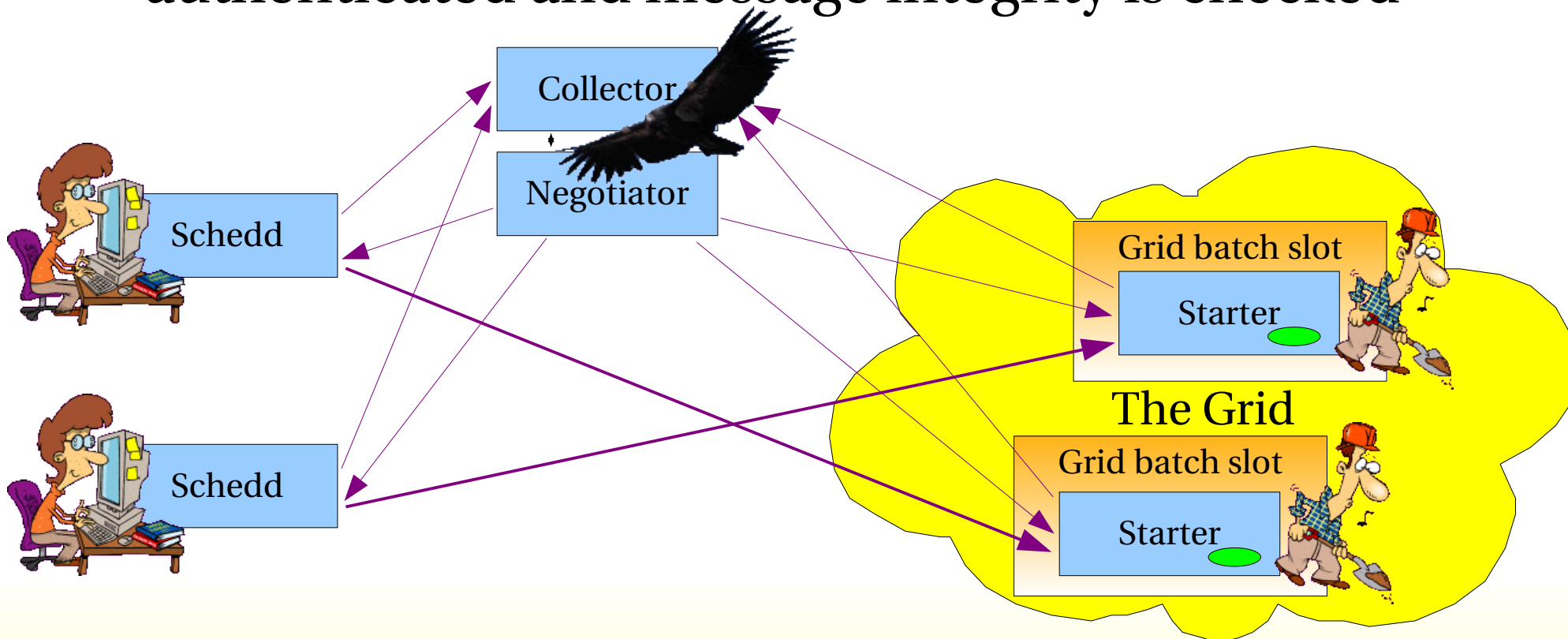- The only requirement is that there is outgoing connectivity

Schedd → Collector

Have jobs, need workers

Negotiator

Have worker, need job, contact me via GCB

Firewall/NAT

The Grid

Schedd

ONE WAY GCB

Grid batch slot

Starter

Open long lived outgoing connection

# Do glideins really work over WAN? (2)

- ## Yes, but it needs GCB to work
  - ### A Condor proxy server
- ## The only requirement is that there is outgoing connectivity

Schedd

Collector

Negotiator

Send job to wgcb

Schedd

Job

Expect job from s2

Firewall/NAT

The Grid

ONE WAY

G C B

Use connection to relay traffic
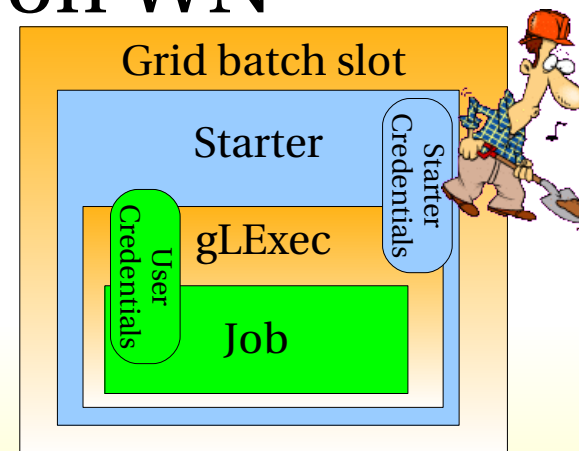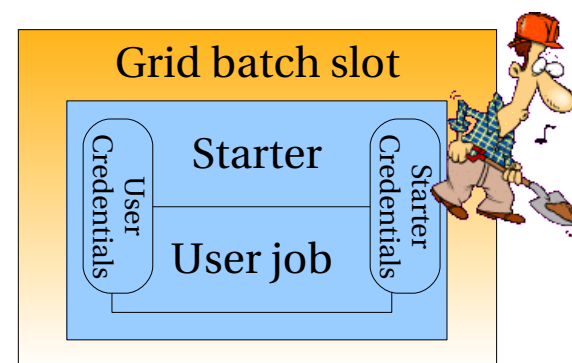
Grid batch slot

Starter

# What about security? [1]

- ## glideinWMS glideins use GSI for authentication
  - ### All daemon to daemon communication is fully authenticated and message integrity is checked

# What about security? (2)

- **However, starter does not run as root!**
  - Without help the user and the starter have to run under the same account!
    - The malicious user job can use starter privileges
    - In this scenario the glidein should only run jobs from the factory user

- **OSG is starting to deploy gLExec on WN**
  - Allows starter to start user job under appropriate UID
  - See gLExec talks #43 and #94

# Any other drawbacks? (1)

- Condor uses a lot of resources
  - Be prepared to budget 1.5Mb of RAM per running process on the submit node
  - Possibly distribute job submission over multiple nodes
- GCB is still in active development phase
  - Production version stable only to ~600 running jobs per GCB/schedd pair
    - Need to deploy many of them to scale
    - Development team is promising much higher scalability

See also Talk #216

# Any other drawbacks? (2)

- Glidein factory load scales with number of Grid sites
  - Budget one node every few dozen sites for current version
    - or pair it with Condor-G matchmakers, like ReSS
  - Work in progress to reduce the load
- If anything goes wrong with the setup, the debugging can be challenging
  - Glidein log files are returned only when the job finishes
    - May not get them back, if it never ends
    - This is a fundamental Grid limitation, not much that can be done about it

# glideinWMS

# Monitoring

# Condor collector monitoring

```
[sfiligoi@cmssrv13 tools]$ condor_status -any

MyType              TargetType        Name

Scheduler           None              cmssrv13.fnal.gov
DaemonMaster        None              cmssrv13.fnal.gov
Negotiator          None              cmssrv13.fnal.gov
glidefactory        None              gov_fnal_1@v41a@cmssrv13
glidefactoryclient  None              gov_fnal_1@v41a@cmssrv13@fnal_
glideclient         None              gov_fnal_1@v41a@cmssrv13@fnal_
glidefactory        None              gov_fnal_2@v41a@cmssrv13
glidefactoryclient  None              gov_fnal_2@v41a@cmssrv13@fnal_
glideclient         None              gov_fnal_2@v41a@cmssrv13@fnal_
Database            None              quill@cmssrv13.fnal.gov
Database            None              quill_glideins1@cmssrv13.fnal.
Database            None              quill_glideins2@cmssrv13.fnal.
Database            None              quill_glideins3@cmssrv13.fnal.
Database            None              quill_glideins4@cmssrv13.fnal.
Scheduler           None              schedd_glideins1@cmssrv13.fnal
DaemonMaster        None              schedd_glideins1@cmssrv13.fnal
Scheduler           None
DaemonMaster        None
Scheduler           None
DaemonMaster        None
Scheduler           None
DaemonMaster        None
Submitter           None
Submitter           None
```
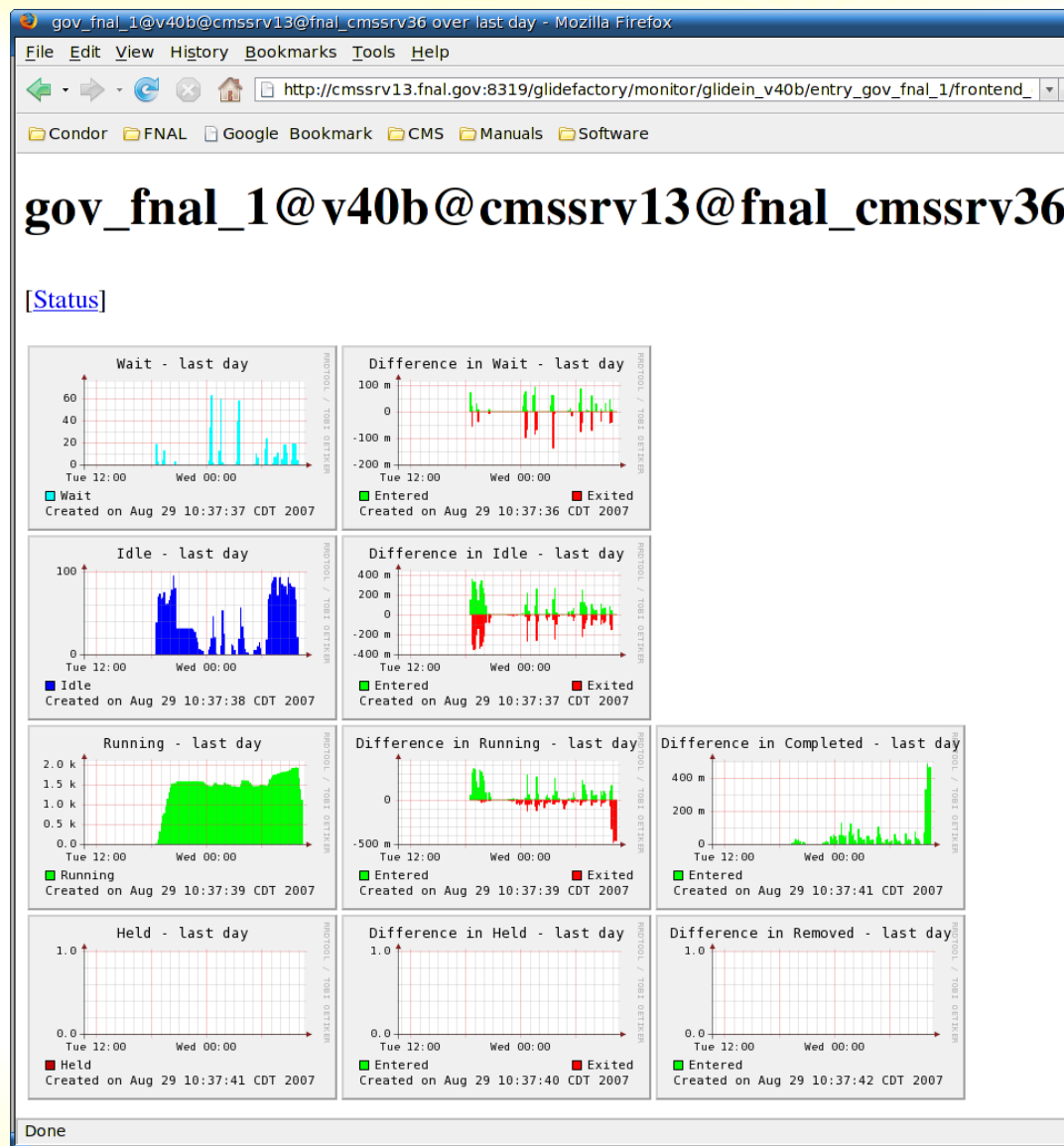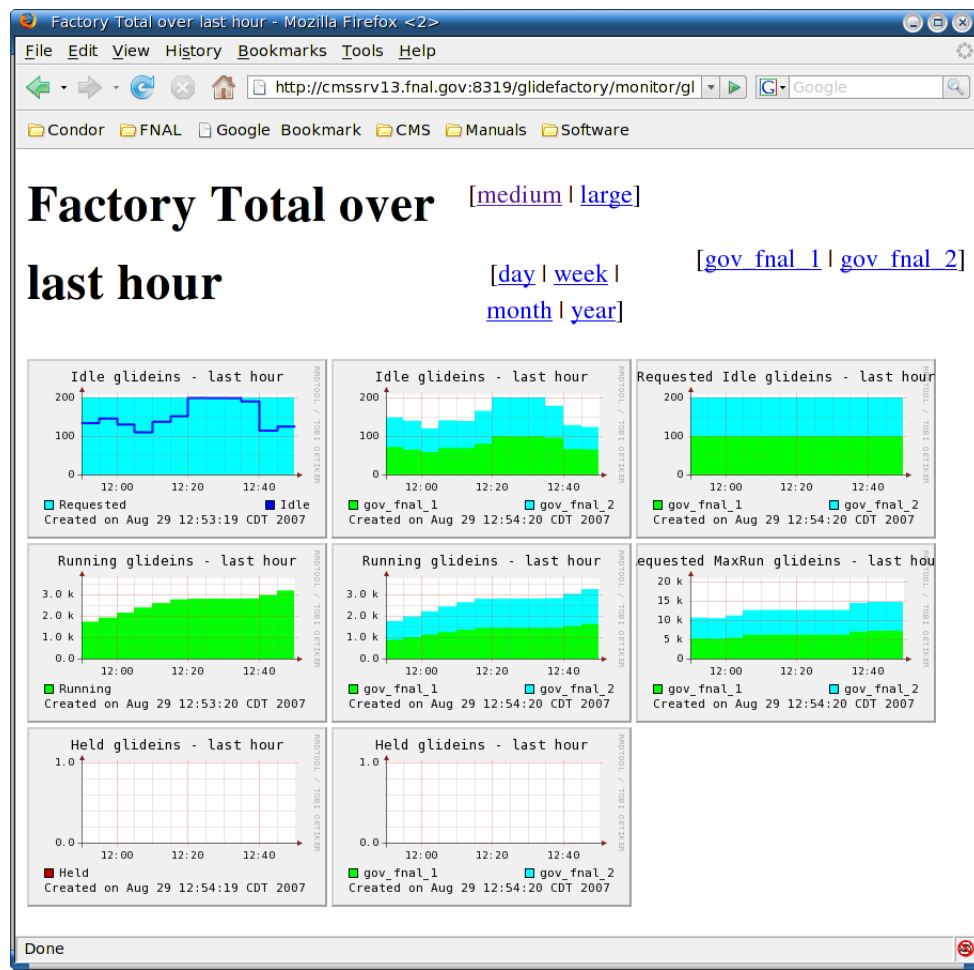
```
MyType = "glidefactoryclient"
TargetType = ""
GlideinMyType = "glidefactoryclient"
Name = "gov_fnal_1@v41a@cmssrv13@fnal_cmssrv36"
ReqGlidein = "gov_fnal_1@v41a@cmssrv13"
ReqFactoryName = "cmssrv13"
ReqGlideinName = "v41a"
ReqEntryName = "gov_fnal_1"
ReqClientName = "fnal_cmssrv36"
ReqClientReqName = "gov_fnal_1@v41a@cmssrv13"
GLIDEIN_Site = "gov_fnal"
GlideinParamGLIDEIN_Collector = "cmssrv37.fnal.gov"
GlideinMonitorRequestedIdle = 28
GlideinMonitorStatusIdle = 100
GlideinMonitorStatusRunning = 1982
GlideinMonitorStatusHeld = 0
GlideinMonitorRequestedMaxRun = 8228
MyAddress = "<131.225.205.230:0>"
LastHeardFrom = 1188412173
UpdatesTotal = 65
UpdatesSequenced = 0
UpdatesLost = 0
UpdatesHistory = "0x00000000000000000000000000000000"
```

```
[sfiligoi@cmssrv13 tools]$ python wmsXMLView.py
<glideinWMS>
    <factory name="gov_fnal_2@v41a@cmssrv13">
        <default_params GLIDEIN_Collector="Fake"/>
        <monitor TotalRequestedMaxRun="8300" TotalRequestedIdle="28" Tota
lStatusRunning="1985" TotalStatusIdle="100" TotalStatusHeld="0"/>
        <attrs GLIDEIN_Site="gov_fnal"/>
        <clients>
            <client name="gov_fnal_2@v41a@cmssrv13@fnal_cmssrv36">
                <client_monitor Idle="3934" Running="3977"/>
                <factory_monitor StatusHeld="0" RequestedMaxRun="8300" Requ
estedIdle="28" StatusIdle="100" StatusRunning="1985"/>
                <params GLIDEIN_Collector="cmssrv37.fnal.gov"/>
                <requests MaxRunningGlideins="8228" IdleGlideins="28"/>
            </client>
        </clients>
    </factory>
    <factory name="gov_fnal_1@v41a@cmssrv13">
        <default_params GLIDEIN_Collector="Fake"/>
        <monitor TotalRequestedMaxRun="8300" TotalRequestedIdle="28" Tota
lStatusRunning="1982" TotalStatusIdle="100" TotalStatusHeld="0"/>
        <attrs GLIDEIN_Site="gov_fnal"/>
        <clients>
            <client name="gov_fnal_1@v41a@cmssrv13@fnal_cmssrv36">
                <client_monitor Idle="3934" Running="3977"/>
                <factory_monitor StatusHeld="0" RequestedMaxRun="8300" Requ
estedIdle="28" StatusIdle="100" StatusRunning="1982"/>
                <params GLIDEIN_Collector="cmssrv37.fnal.gov"/>
                <requests MaxRunningGlideins="8228" IdleGlideins="28"/>
            </client>
        </clients>
    </factory>
</glideinWMS>
```
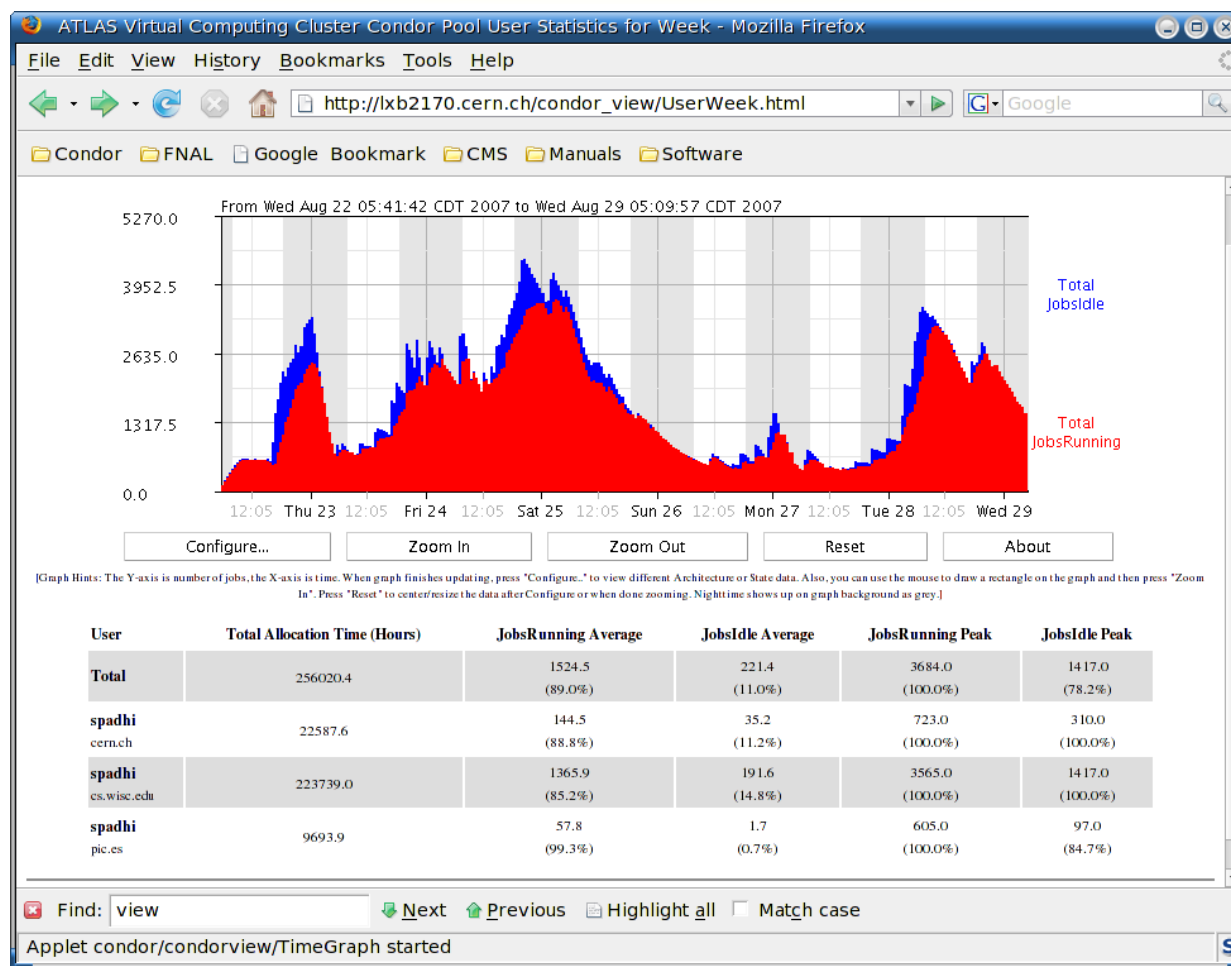
# Status Web monitoring

# Status XML monitoring

# CondorView Monitoring



**Standard Condor tool, not glideinWMS specific**

This one is actually from the CRONUS site

# Conclusions

## Condor Glideins

- Can shield user jobs from the Grid

- Give you total control over your jobs

- Allow you to have more control over the jobs scheduling

## GlideinWMS

- An automatic way to create glidein pools on the fly

- Needs some initial effort, but then it operates on its own

# glideinWMS home page

http://home.fnal.gov/~sfiligoi/glideinWMS/

sfiligoi@fnal.gov