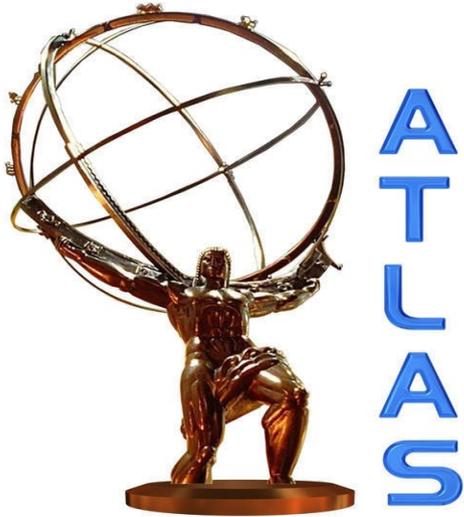


CRONUS: A Condor Glide-in Based ATLAS Production Executor



Sanjay Padhi

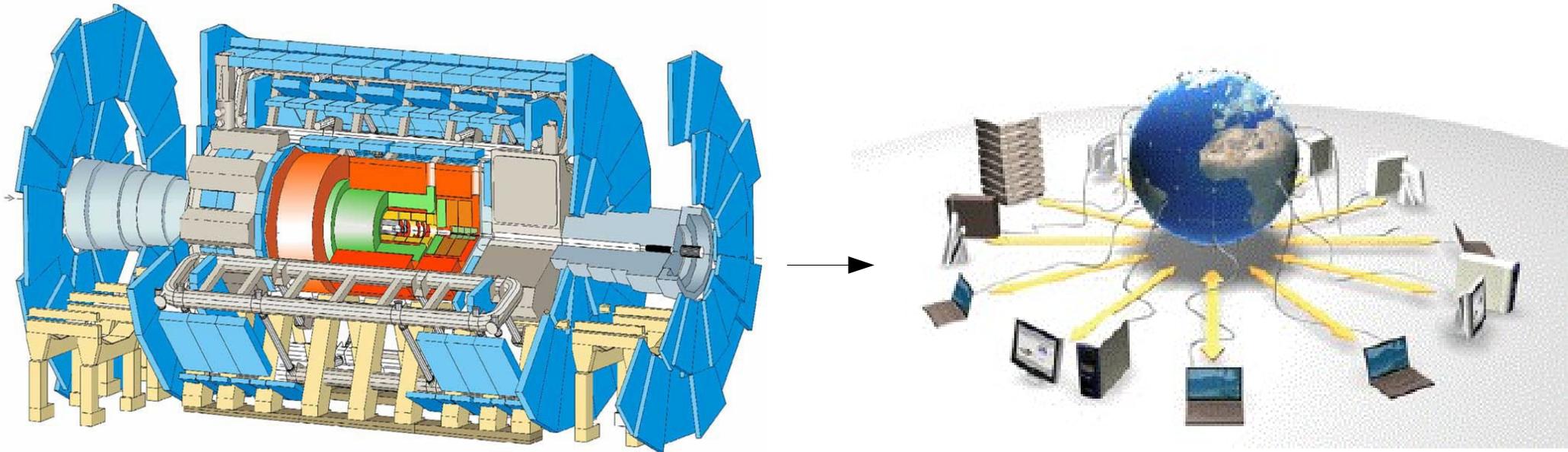
University of Wisconsin-Madison



International Conference on Computing in High Energy
and Nuclear Physics, Victoria, Canada, 02 – 07 Sept. 2007

Many thanks to R. Walker, Sau Lan Wu, Miron Livny,
the Condor Team & several members of the ATLAS Collaboration

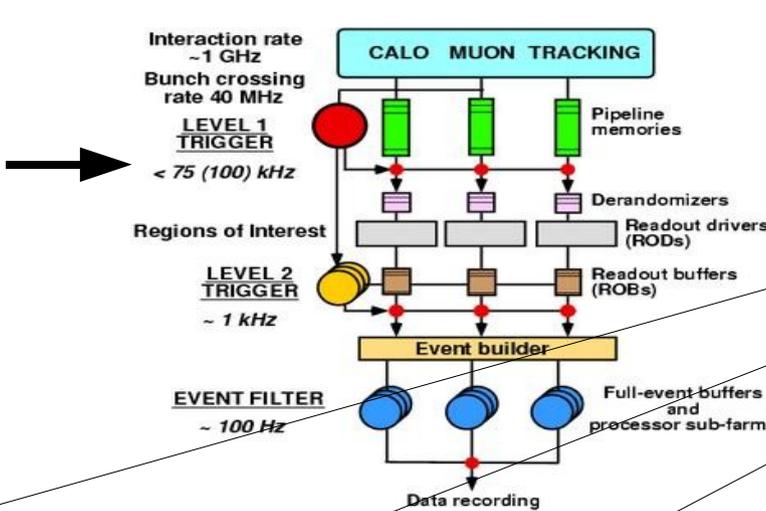
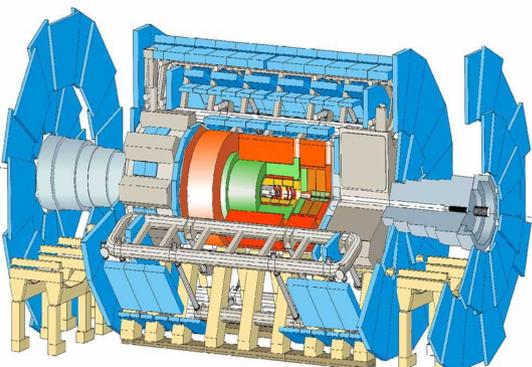
High Energy Physics and The Grid



- HEP is evolving :
 - Physics analysis is no longer at the Experimental site.
- Physicists are evolving:
 - Illiteracy in Computing is no longer possible.
- Collaborations are evolving:
 - Computing using local batch systems are just not enough.

Inter-operable Grid Computing is a necessity

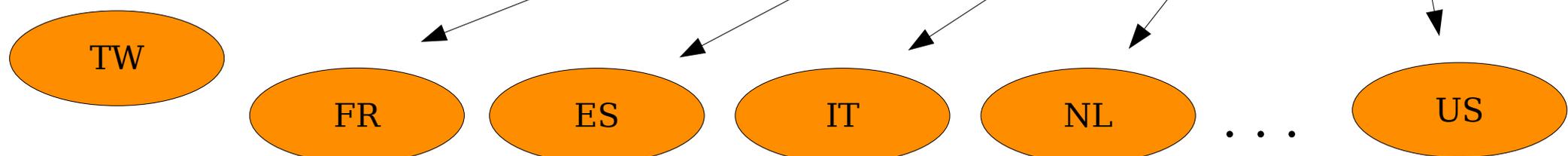
Classical LHC Tier Model



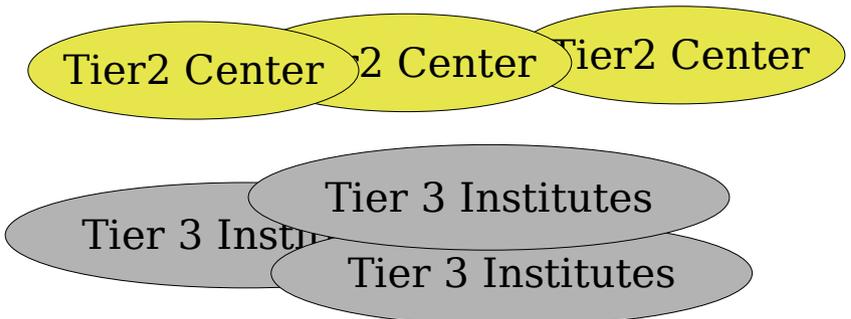
Archival and Distribution of RAW DATA



TIER 1 hosts and provide long term access & archival - subset of RAW DATA



TIER 2s - provide calibration constants, simulation & analysis



RAW Size ~ 1.6MB (1600 TB/year)
 ESD Size ~ 1.0MB (1000 TB/year)
 AOD Size ~ 0.1MB (100 TB/year)
 Number of Events (2008) - 1000M

For all this to work:

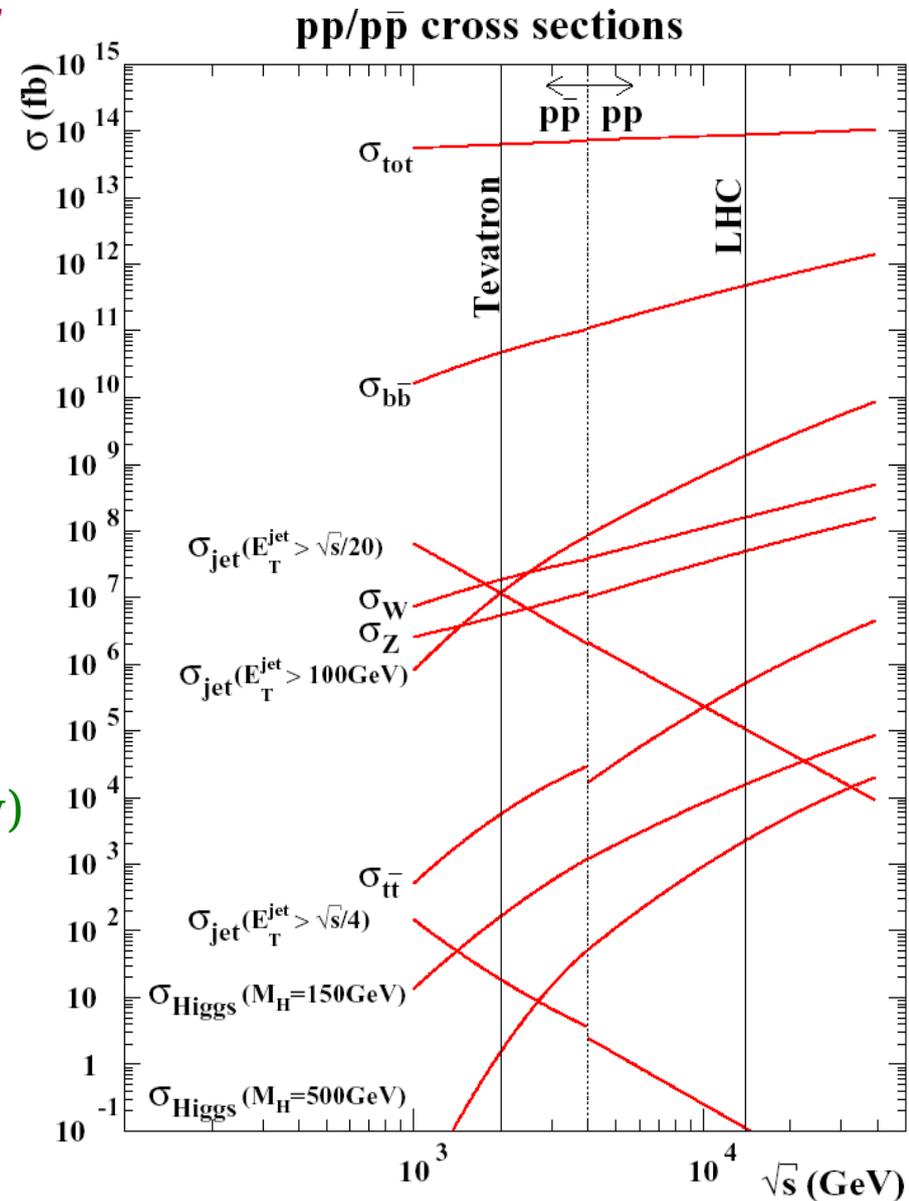
Need a strong relationship between scientists and "the Grid Community"

New Physics (SUSY) is expected to have “small” cross sections

Standard Physics processes have relatively large uncertainties but with huge cross section
 [10 pb⁻¹ ~ 1 day at 1/100 of the design lumi.]

- Total inelastic: ~ 0.1 barn
- Inclusive QCD multijets ~ 10⁵ nb (10⁹/day)
- Inclusive bbar: ~ 10³ nb (10⁷/day)
- Inclusive $\gamma/Z/W$: ~ 0.1/1.5/10 nb (1- 10k /day)
- Inclusive top ~ 0.89 nb (8.9 k /day)
- Z(vv) + Jets ~ 0.3 nb (3 k /day)

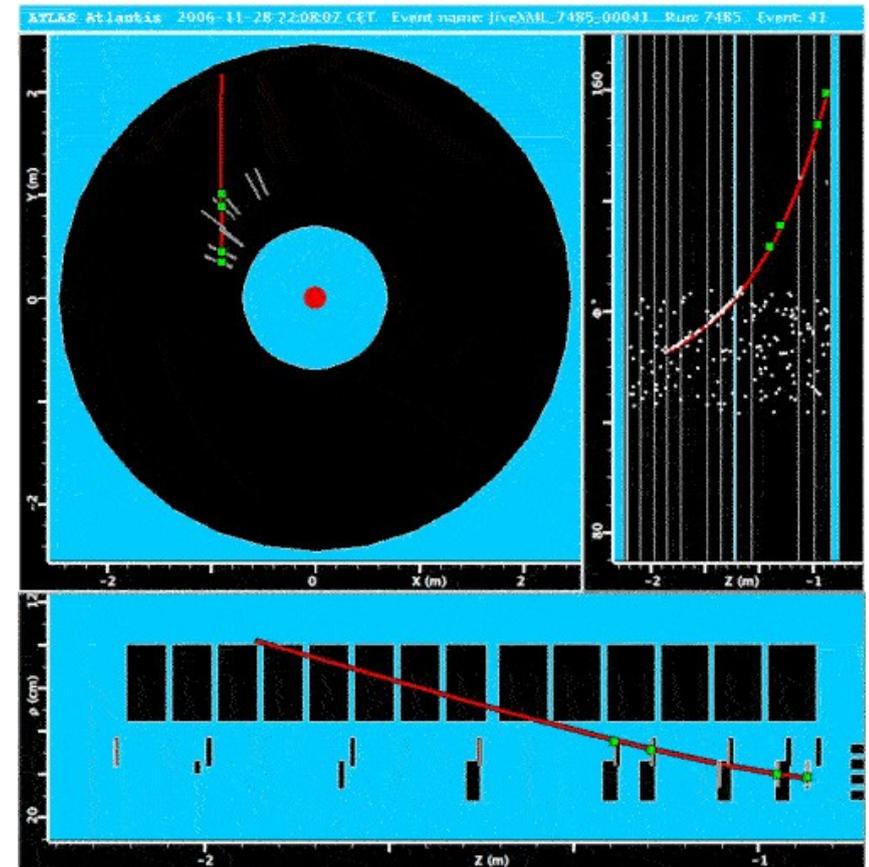
Compare this to 1 TeV SUSY ~ 200 pb



Need to have an accurate understanding of the background – **to claim a discovery**

In order to fulfill ATLAS Timeline ...

- Running continuously throughout the year (increasing rates):
 - Simulation production
 - Cosmic ray data-taking (detector commissioning)
- January to June:
 - Data streaming tests
- February and May:
 - Intensive Tier-0 tests
- From February onwards:
 - Data Distribution tests
- From March onwards:
 - Distributed Analysis (intensive tests)
- May to July:
 - Calibration Data Challenge
- June to October:
 - Full Dress Rehearsal
- November:
 - **GO!**



Dario Barberis: ATLAS Computing Plans for 2007

CRONUS was introduced !!! [Motivated by CDF-CAF Igor, Frank et. al.]
– Based on Condor Glide-in Technology

Grid Federation - EGEE

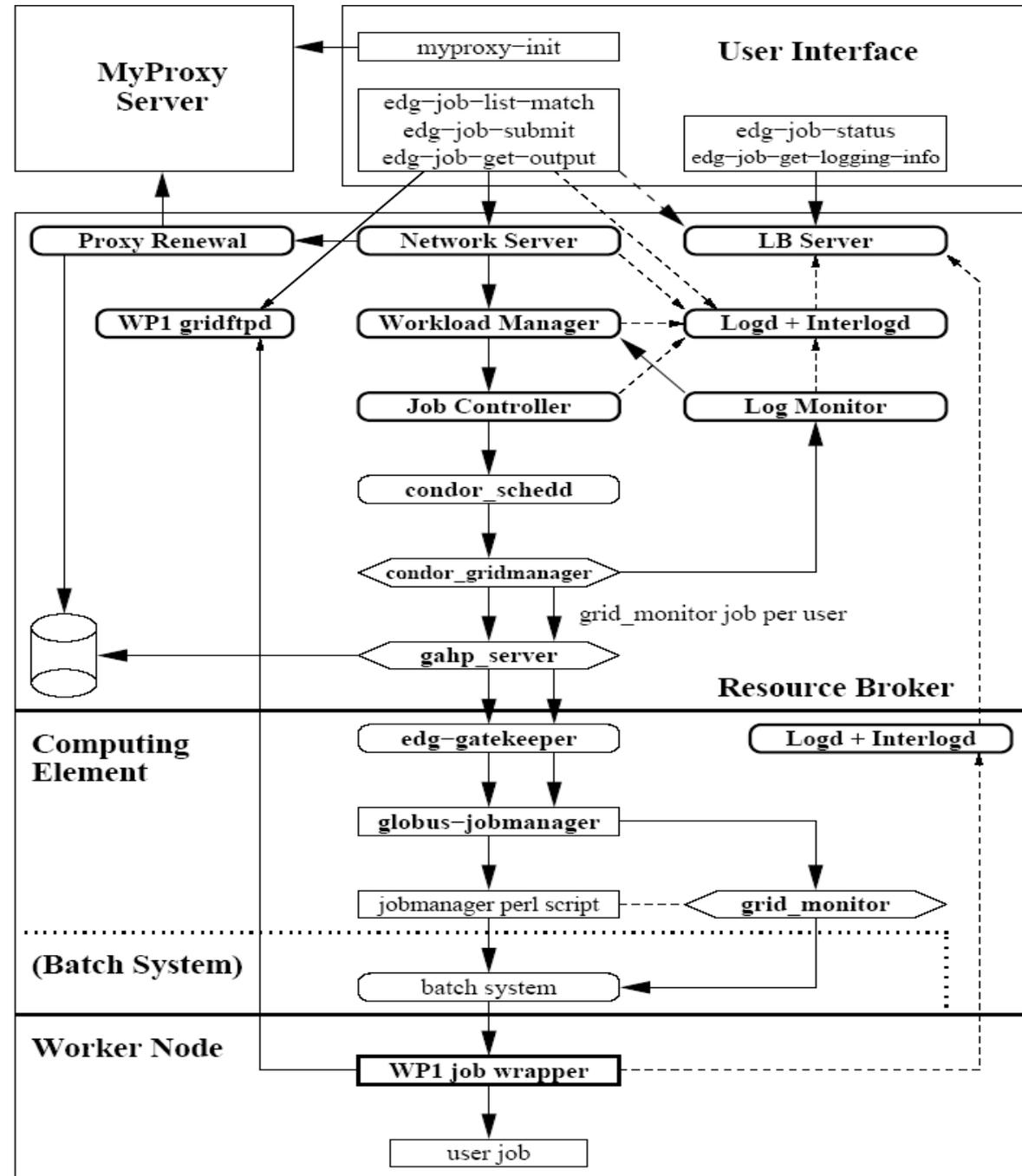
Well designed components:

- A Workload Management System
- A Data Management System
- An Information System
- An Authorization and Authentication
- An Accounting System
- Various monitoring services

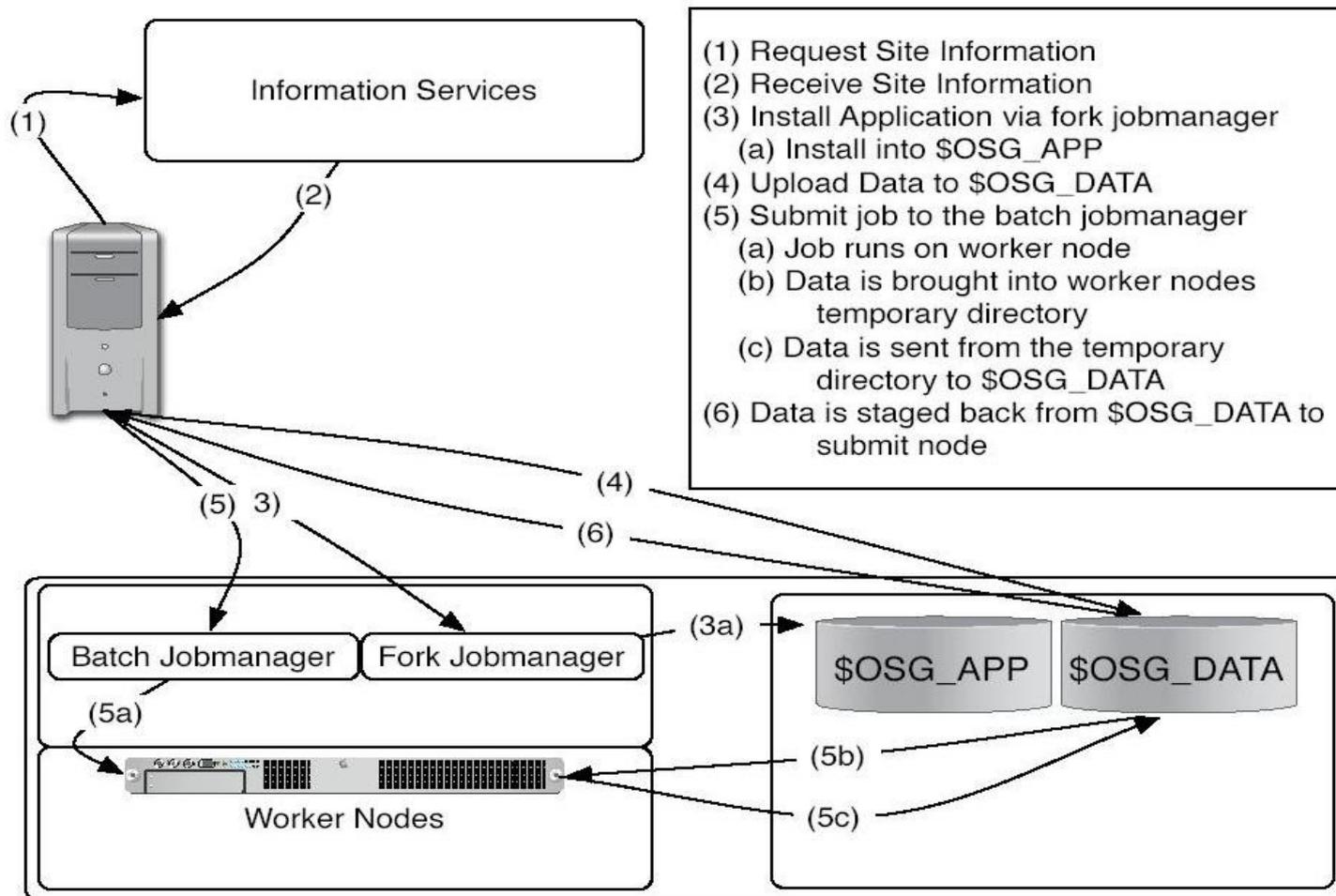
CE or the gatekeeper

- Accepts jobs from Condor-G
- Creates a **job manager (JM)** per job
- Generic interface to the batch system
- The JM *only* submits or cancel a job
- The **grid monitor** queries the status

Self consistent early binding model



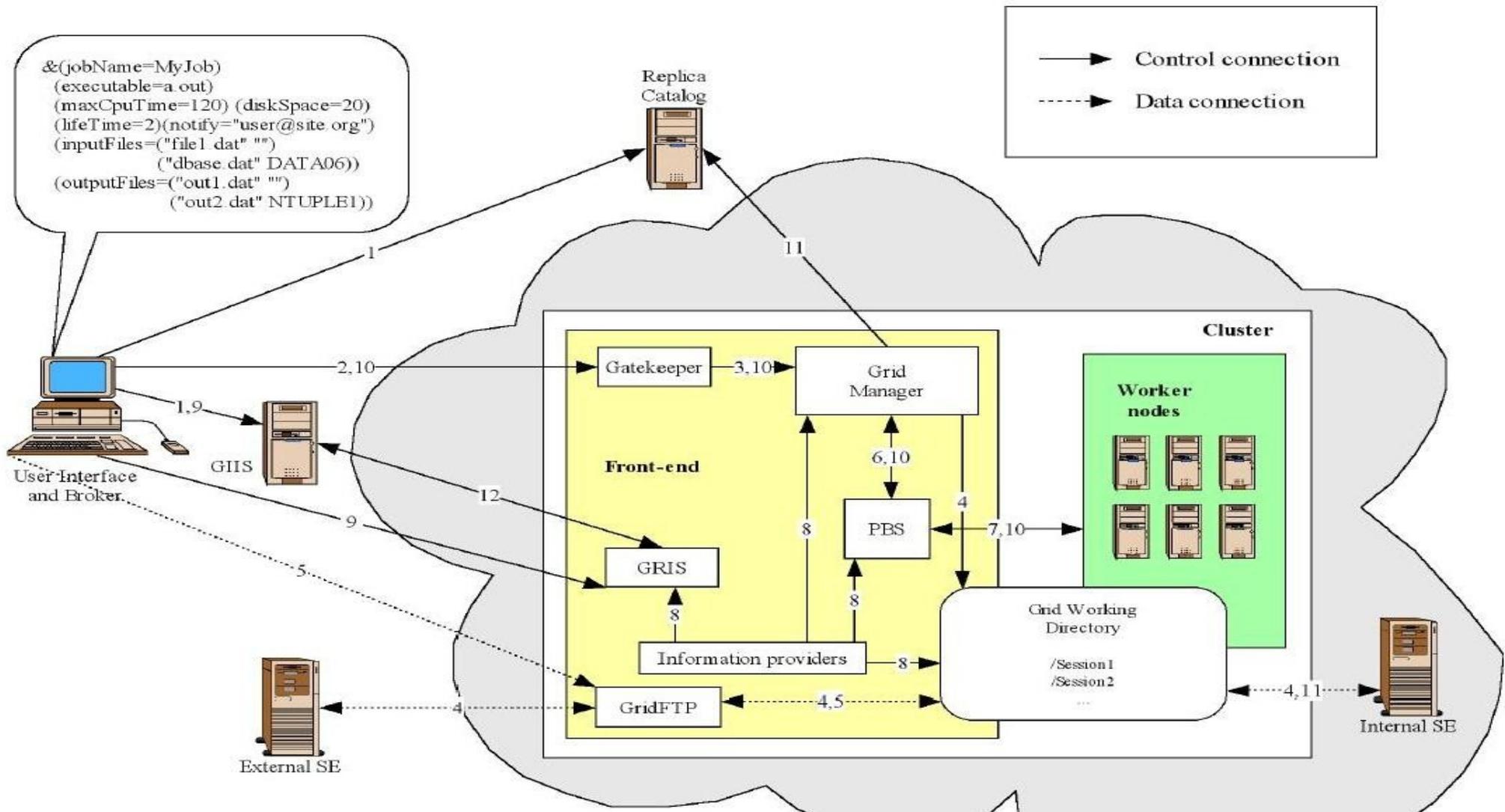
Grid Federation - OSG



OSG by its design allows several flexibilities: [Classical Early binding model]

- Allows temporary storage pools
- GateKeeper apart from job scheduling also responsible for authentication/authorization
- Currently no resource broker (individual VOs if needed can provide their own)
- Recent emphasis is more on dynamic information system

Grid Federation - NorduGrid



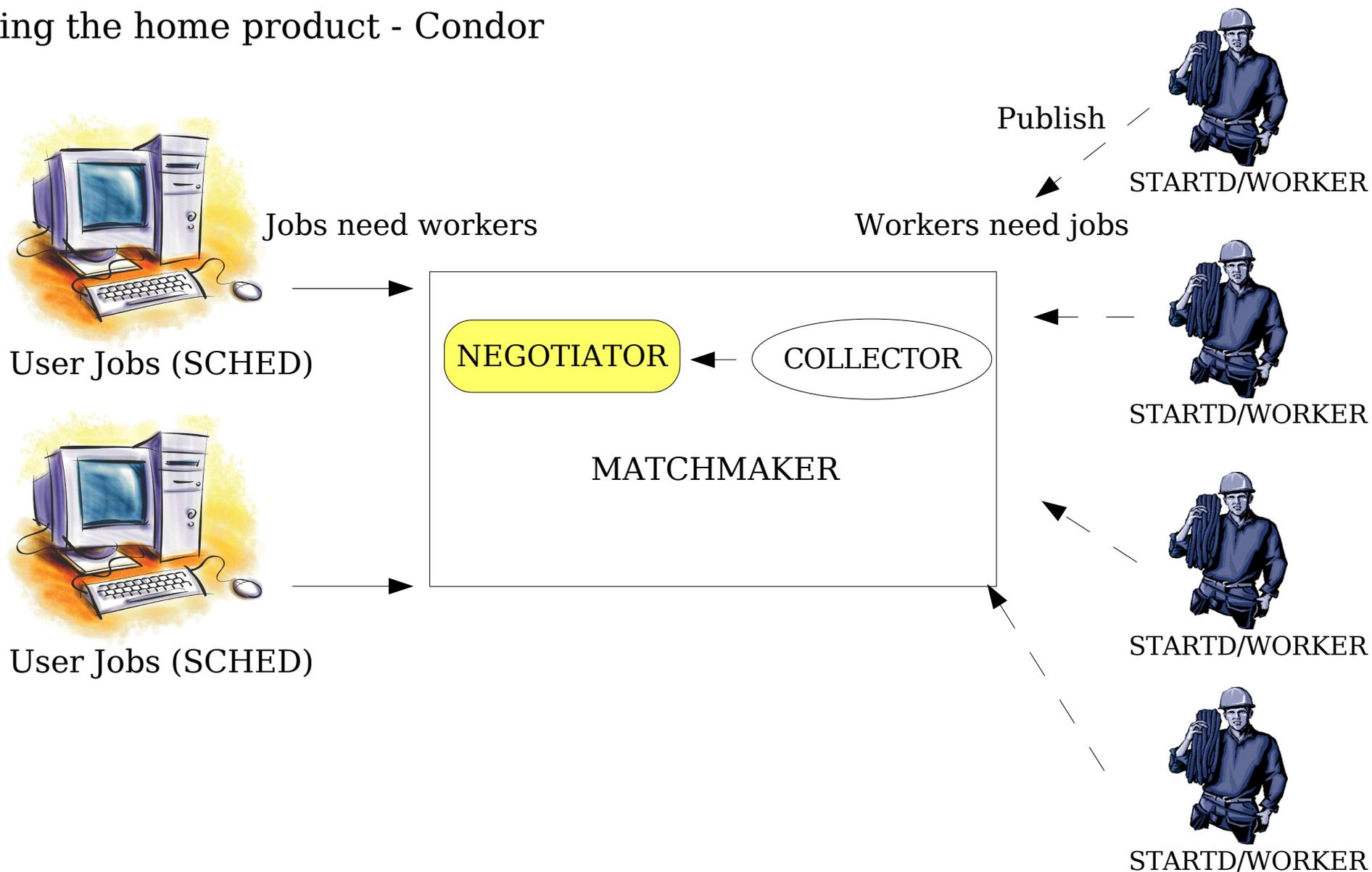
The greatest strength is in the gatekeeper:

- Responsible for job scheduling, StageIn & StageOut of the DATA
- Authentication and Authorization etc.

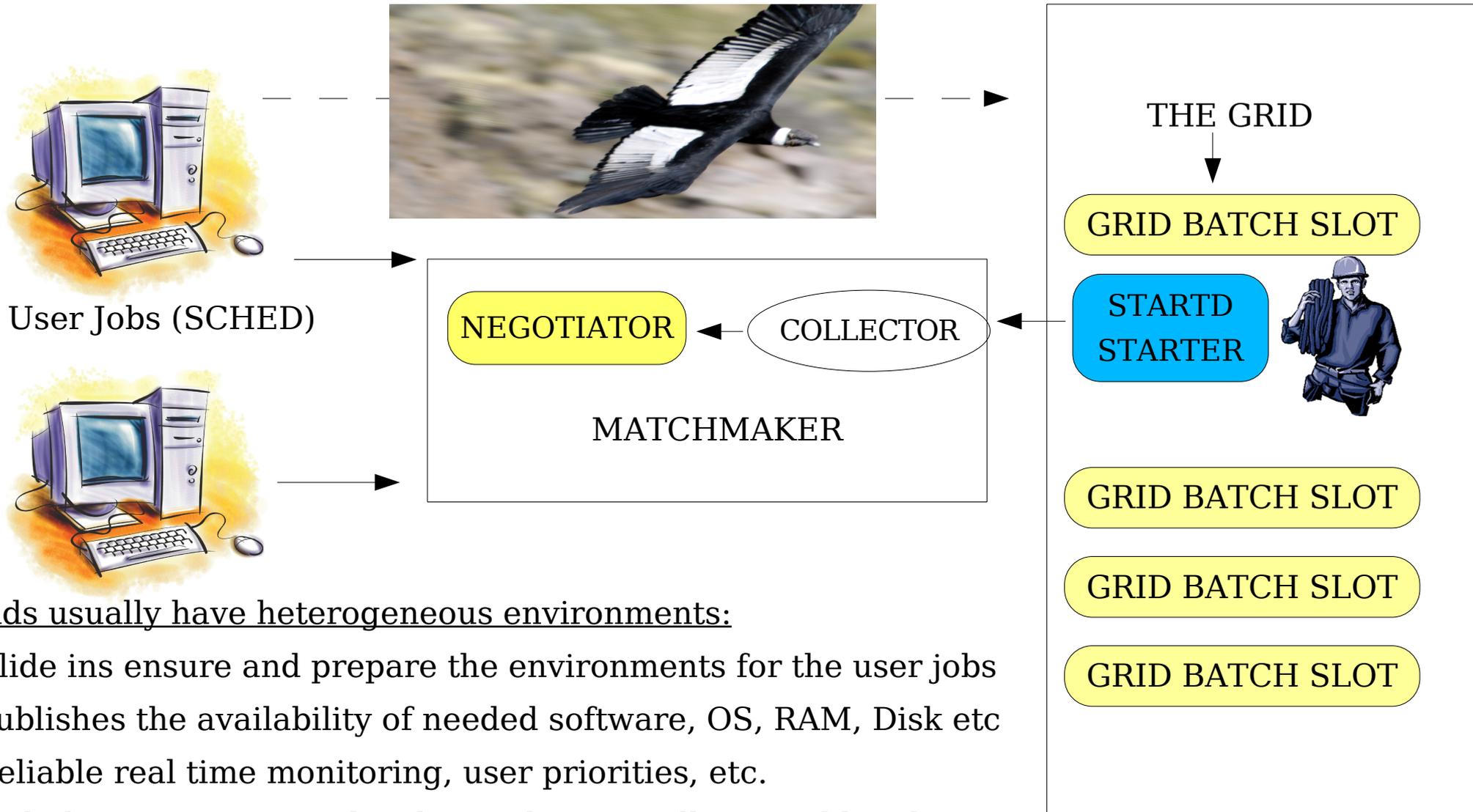
Although highly heterogeneous batch system, the information system is the very best

The Condor

Using the home product - Condor



Glide-ins are the **STARTD/WORKER** submitted as a grid job

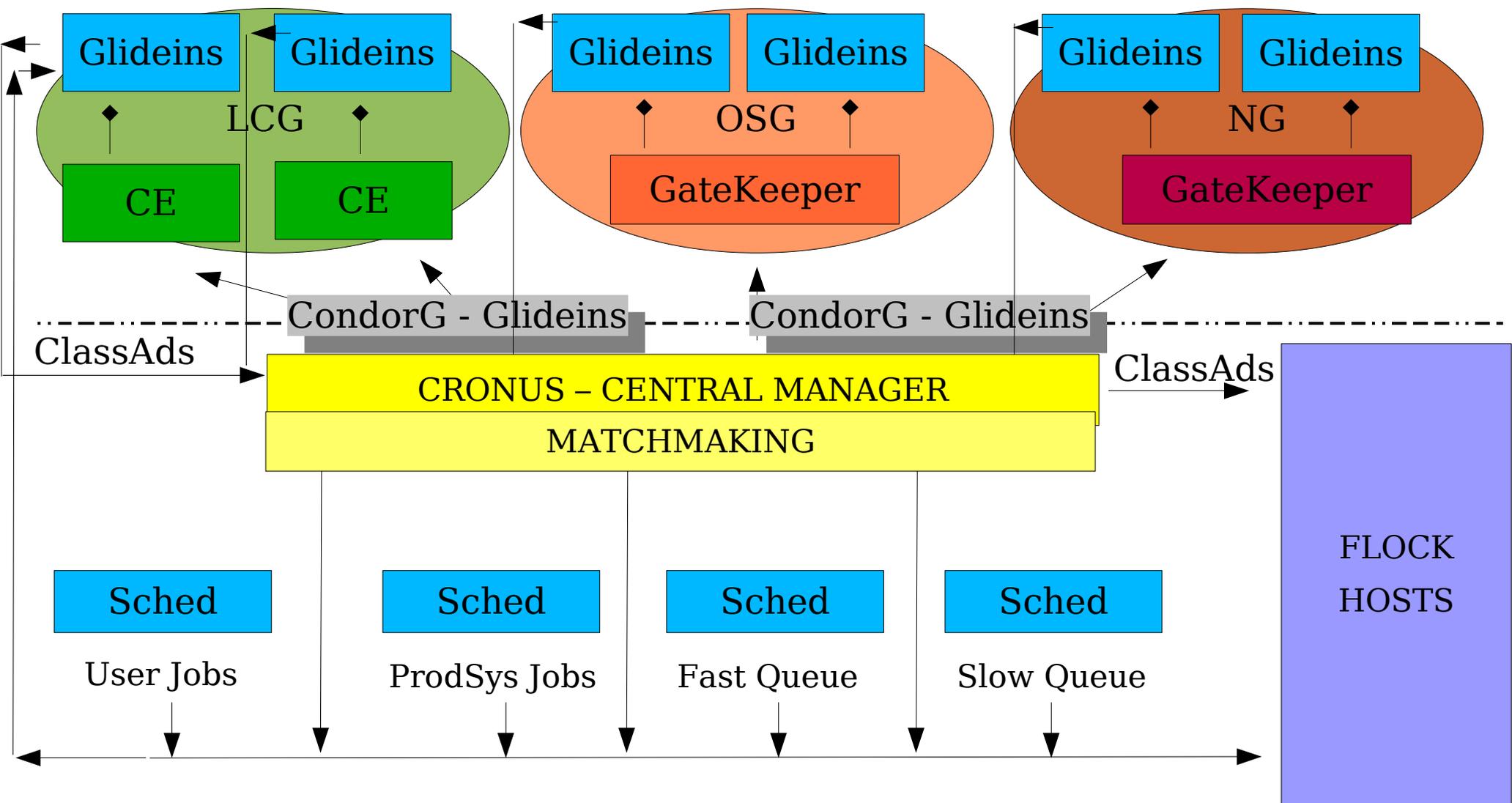


Grids usually have heterogeneous environments:

- Glide ins ensure and prepare the environments for the user jobs
- Publishes the availability of needed software, OS, RAM, Disk etc
- Reliable real time monitoring, user priorities, etc.
- Includes group quotas, live log updates – Fully virtual batch system

⇒ **See the talk by Igor Sfiligoi (FNAL)**

The concept of late binding is intrinsic to Condor via the ClassAds

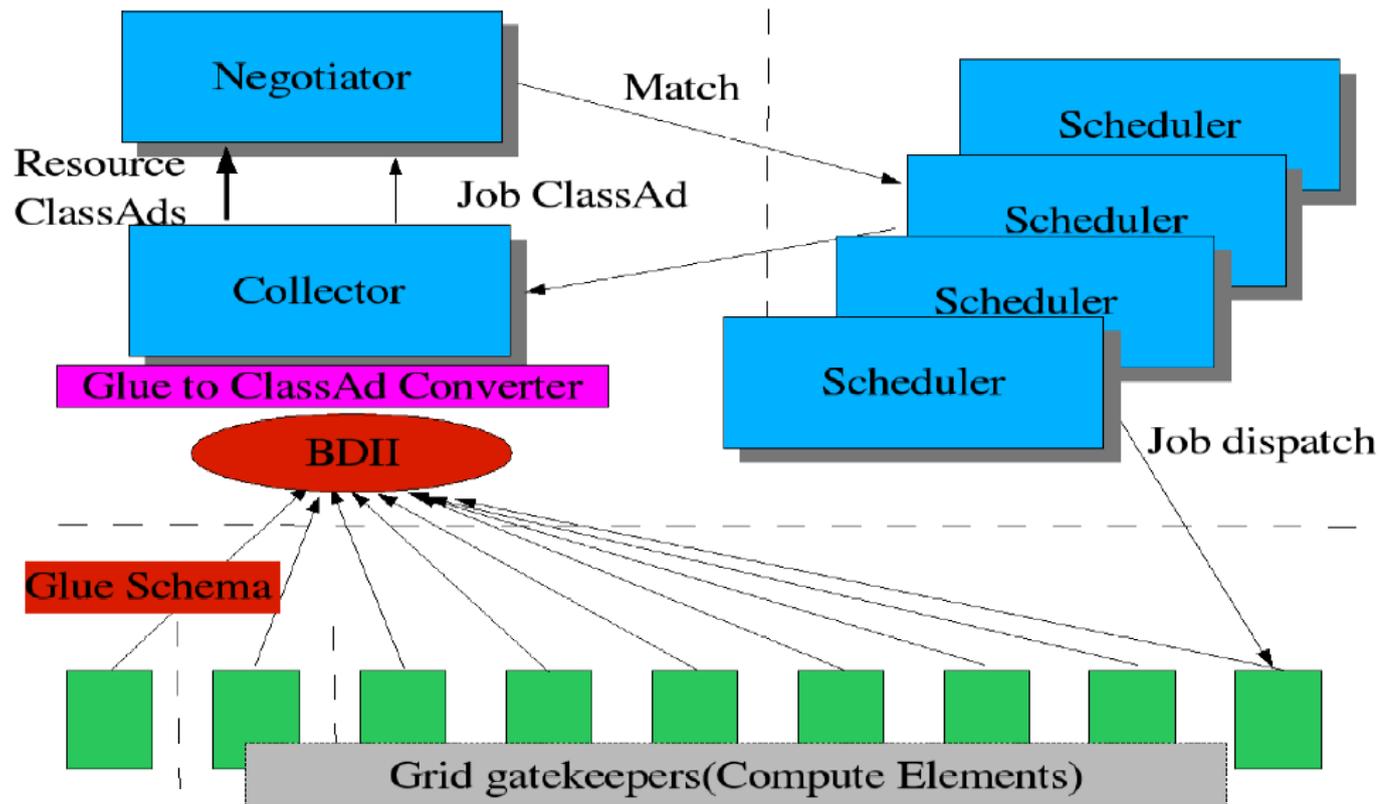


Only one communication language among all - ClassAds

Private Resources

Glide-Ins using Condor-G or WMS Resource Broker

- Condor Startup Glide-Ins are submitted using Condor-G
 - No need to re-invent the wheels !!!
- Glide-in Startup jobs ensures the needed resources and the ATLAS Environment

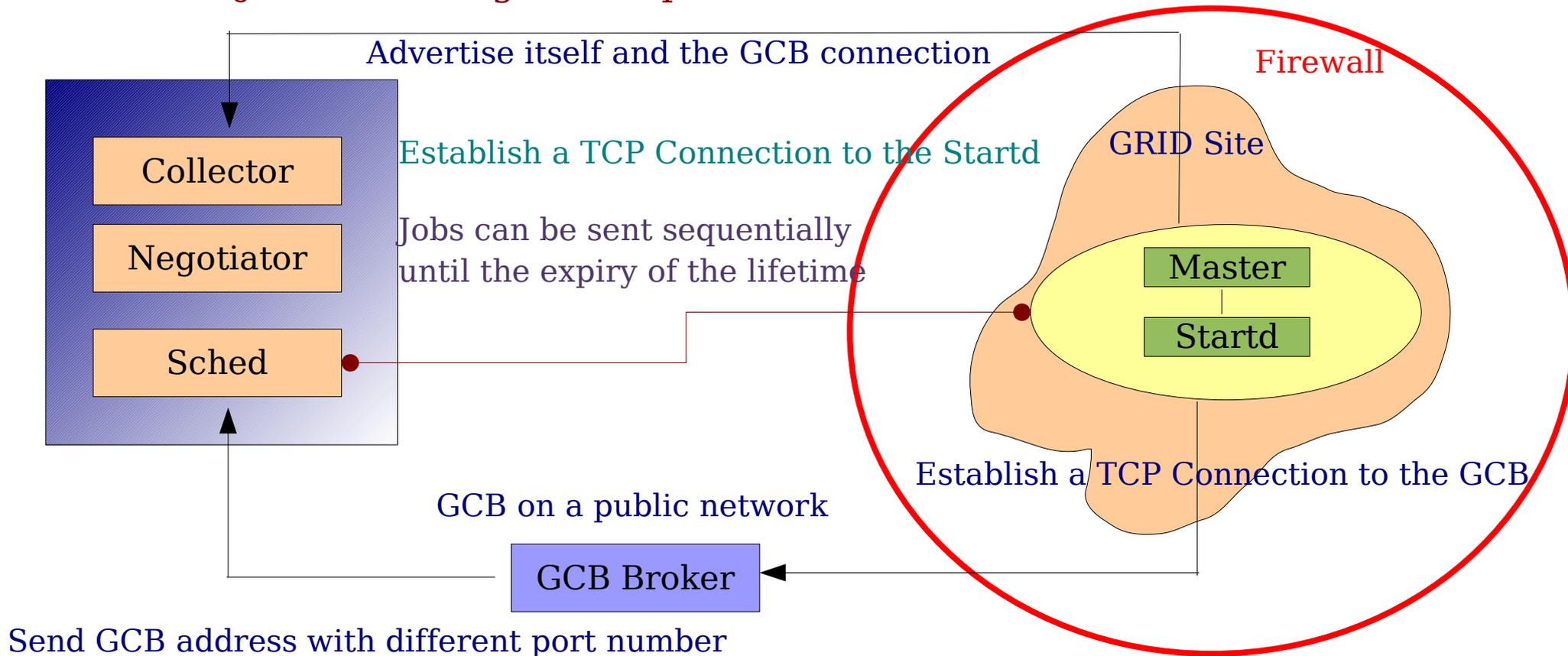


LCG-BDII: atlas-bdii.cern.ch; OSG-BDII: data.grid.iu.edu; NG-BDII: arc-bdii.cern.ch

Collector uses all informations systems across various grid flavours

GCB is a Condor proxy service to the Startds :

- Part of Condor distribution, fully integrated in Condor
- Just few configuration parameters.

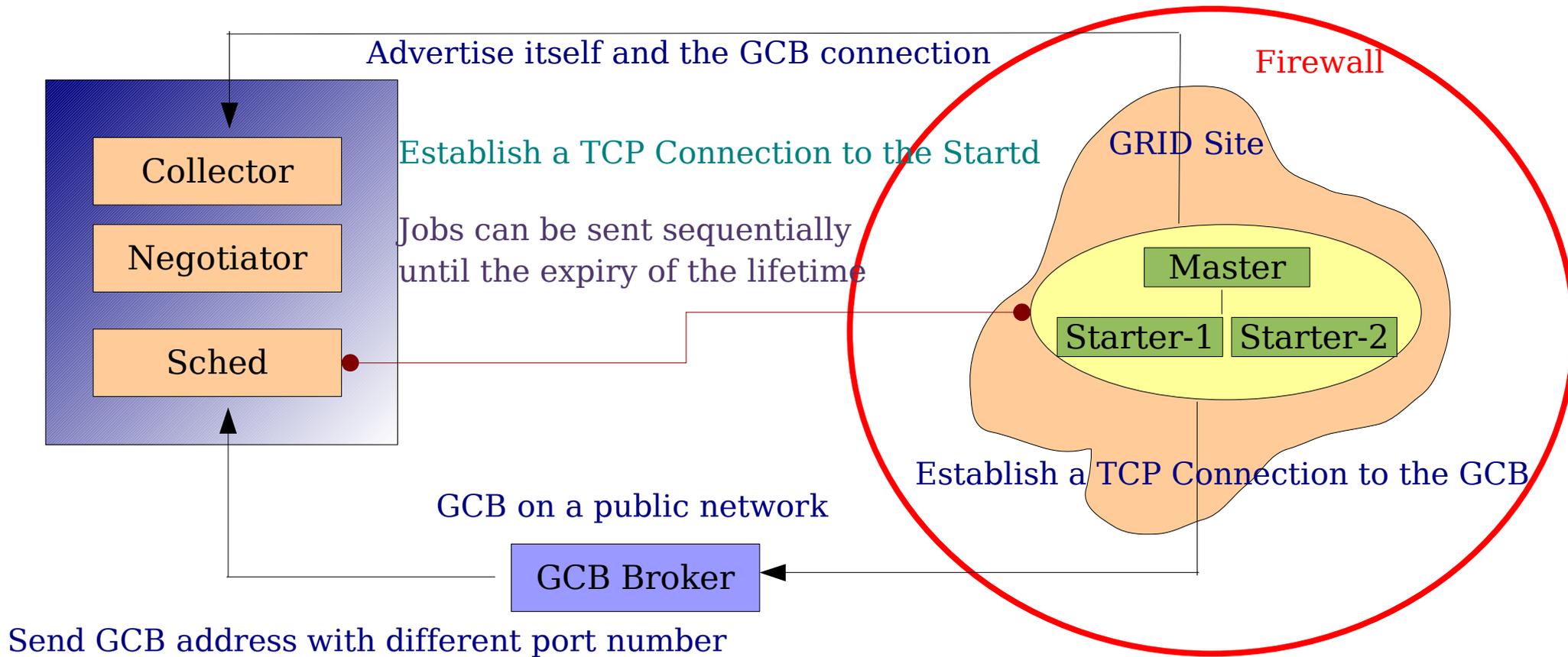


For scalability reasons – use multiple GCB instances

Multiple Starters (2)

STARTER-1 : Responsible for “standard” production job

STARTER-2 : Responsible for instant DATA transfer/subscription/error recovery jobs



Example Virtual Cluster: Uniformity across various flavours of Grid – OSG, EGEE ..

FileSystemDomain = "bigmac-lcg-ce.physics.utoronto.ca"
FileSystemDomain = "cclcgceli02.in2p3.fr"
FileSystemDomain = "ce101.cern.ch"
FileSystemDomain = "ce102.cern.ch"
FileSystemDomain = "cern.ch"
FileSystemDomain = "cit-gatekeeper.ultralight.org"
FileSystemDomain = "epgce1.ph.bham.ac.uk"
FileSystemDomain = "ft.uam.es"
FileSystemDomain = "gate02.grid.umich.edu"
FileSystemDomain = "grid003.ft.uam.es"
FileSystemDomain = "grid-ce0.desy.de"
FileSystemDomain = "grid-ce1.desy.de"
FileSystemDomain = "grid-ce2.desy.de"
FileSystemDomain = "hamptonu.edu"
FileSystemDomain = "ifaece01.pic.es"
FileSystemDomain = "lancs.ac.uk"
FileSystemDomain = "lcg002.ihep.ac.cn"
FileSystemDomain = "lcg2ce.ific.uv.es"
FileSystemDomain = "lcgce01.gridpp.rl.ac.uk"
FileSystemDomain = "lcgce02.nic.ualberta.ca"
FileSystemDomain = "lcg-ce0.ifh.de"
FileSystemDomain = "qmul.ac.uk"
FileSystemDomain = "teraport.edu"
FileSystemDomain = "tier2-osg.uchicago.edu"
FileSystemDomain = "tp-osg.uchicago.edu"
FileSystemDomain = "vampire.accre.vanderbilt.edu"
FileSystemDomain = "vanderbilt.edu"
FileSystemDomain = "zeus02.cyf-kr.edu.pl"

Jobs Running on EGEE

CRONUS

Jobs Running at OSG

- Fully inter-operable Virtual batch System

```
MyType = "Machine"  
TargetType = "Job"  
Name = "8505@ccwali04.in2p3.fr"  
Machine = "ccwali04.in2p3.fr"  
Rank = 0.000000  
CpuBusy = FALSE  
IS_GLIDEIN = TRUE
```

MACHINE ClassAd

```
DaemonStopTime = DaemonStartTime + 172740  
ATLAS_SOFT = "/afs/in2p3.fr/group/atlas/sw/software"  
ATHENA_VERSION = "10.0.1,10.0.4,11.0.3,11.0.41,11.0.42,11.0.5,11.2.0,12.0.1,12.0.2,12.0.3,12.0.31,12.2.0"  
GLOBAL_SE = "LYON"
```

CRONUS

Job ClassAd

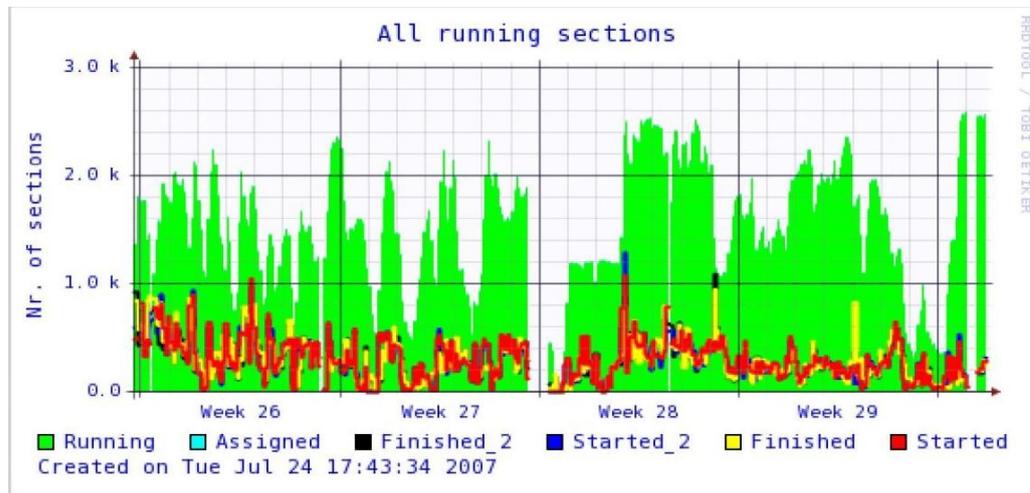
```
universe = vanilla  
executable = EvGenproduction.sh  
transfer_files = always  
output = job/out.mc11test.005030.Jimmy_jetsJ1.evgen.EVNT.v12000301._00002  
error = job/err.mc11test.005030.Jimmy_jetsJ1.evgen.EVNT.v12000301._00002  
log = job/SubmissionLog2  
stream_output = False  
stream_error = False  
notification = never  
Requirements = (StringlistMember("12.0.31",ATHENA_VERSION) && ( GLOBAL_SE == "LYON" )) && ( Arch == "INTEL")  
WhenToTransferOutput = ON_EXIT  
OnExitRemove = TRUE  
arguments = "mc11test.005030.Jimmy_jetsJ1.evgen.v12000301_tid3326 \  
5030 5000 5000 2 DC3.005030.Jimmy_jetsJ1.py \  
mc11test.005030.Jimmy_jetsJ1.evgen.EVNT.v12000301._00002.pool.root NONE NONE NONE"
```

Glide-In Scalability

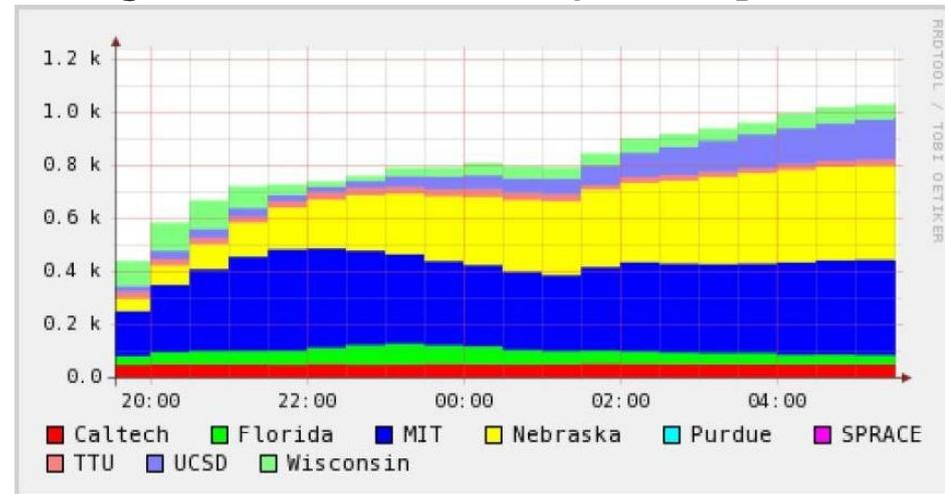
ATLAS - CRONUS using ~ 5200 CPUs in parallel [http://lxb2170.cern.ch/condor_view/]



CDF-CAF ~ 2.5 K Jobs in parallel

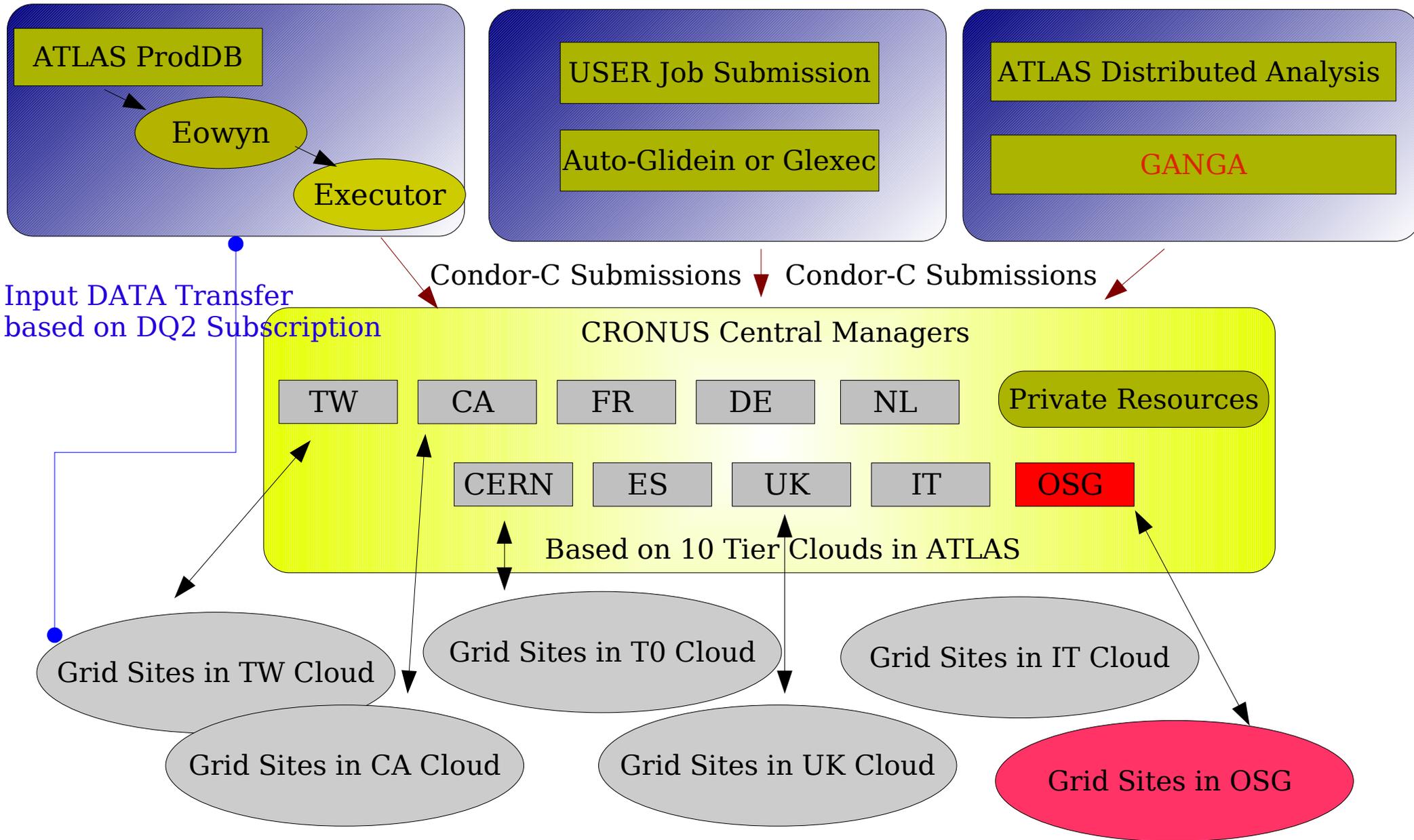


CMS-glideinWMS ~ 1.2 K Jobs in parallel

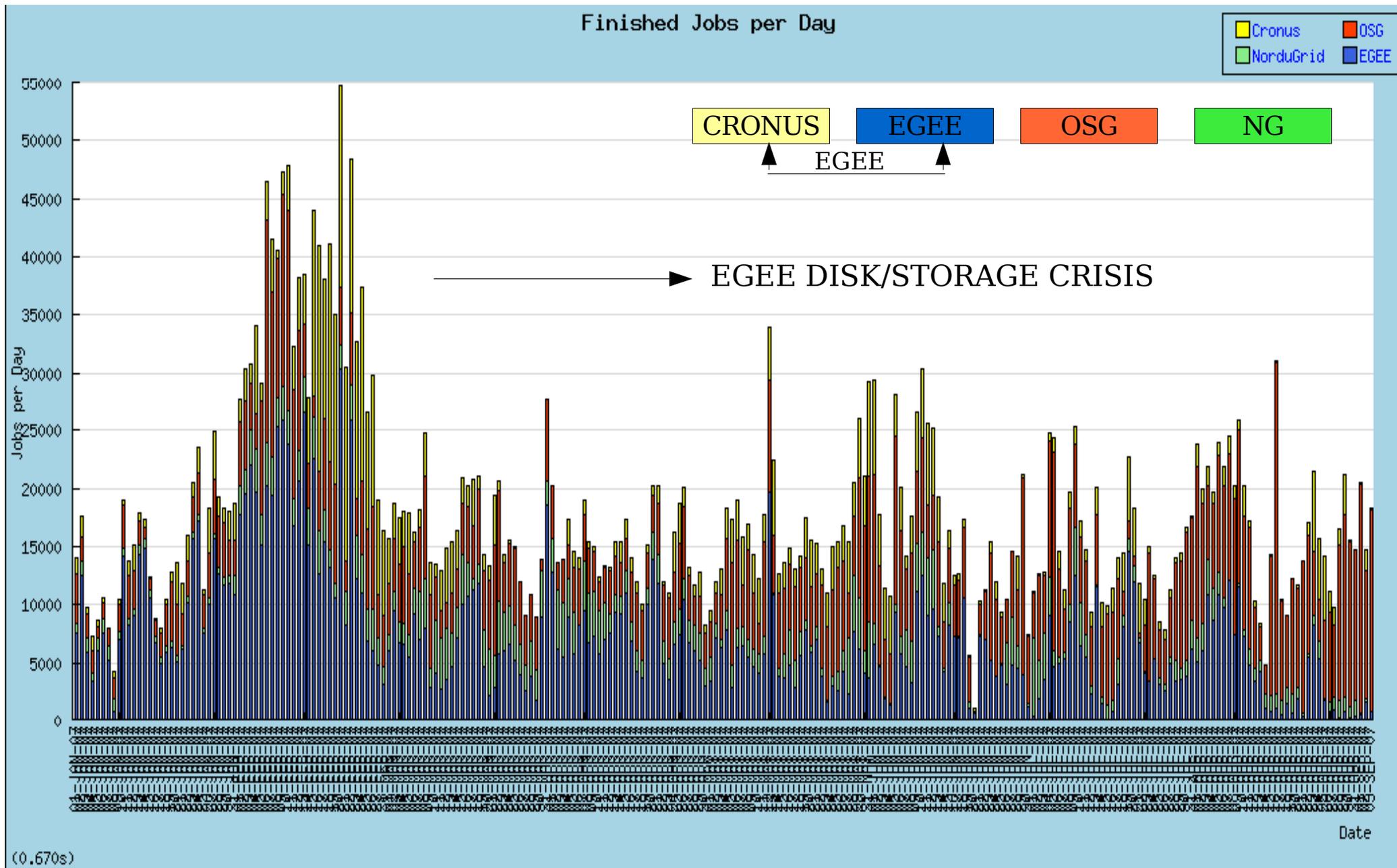


CRONUS - ATLAS Production Executor

Auto-Glideins – A user uses his own GlideIn [submitted automatically with the job]



ATLAS Production using CRONUS



Submission to NorduGrid – First time using the client method

[Thanks to the Condor team, Oxana and Alex Read]

Matchmaking at CRONUS using the information from arc-bdii ldap server

Runtime ENV using the Matchmaker

Client – GHAP modified to accommodate stageIn/stageOut by the gatekeeper

No need to change anything from the NorduGrid side

Advantage: Treat exactly like any other gatekeeper

- Syntax/protocol same as used in the NorduGrid RSL

- Runtime ENV executed by the gatekeeper

 - (runtimeenvironment=APPS/HEP/ATLAS-12.0.6.2)

 - Or uploaded for the user by the gatekeeper

 - rls://atlasrls.nordugrid.org:39281/NGExecWrapper7

- Keeping the uniformity in language (ClassAds)

StageIn/StageOut can be local or SURL OR provided by the rls

- rls://atlasrls.nordugrid.org:39281/mc11.003035.J2_Pt_35_70.recon.ESD.v11000301._00089.pool.root

or

- srm://srm.ndgf.org/pnfs/ndgf.org/data/atlas/mc11/

 - mc11.003035.J2_Pt_35_70.recon.ESD.v11000301._00089.pool.root

or

- transfer_output_files=mc11.003035.J2_Pt_35_70.recon.ESD.v11000301._00089.pool.root

Major technological development in terms of interoperability of Grid Federations

Cronus Operations

EXECUTOR	TYPE	FINJOBS	FINCPU	FINWALL	FAILJOBS	FAILCPU	FAILWALL	SUBMITTED	RUNNING	JOBEFF	WALLEFF
CondorG		883658	12836641416	14056799282	711973	994891364	1853968628	392420	247080	55.37	88.34
LCG-DQ		946647	15294451105	17684032852	757858	1028625168	2698145564	1200175	332482	55.53	86.76
panda		1516360	32743667765	30681711168	277820	1357691400	5318389455	1360281		84.51	85.22
Dulcinea		449981	5952315750	6004247520	88936	619561287	631320241	64441	102778	83.49	90.48
Cronus		645607	10968756837	12686539686	412060	728401272	1942276543	231239	267579	61.04	86.72

Although Cronus is very new and ~2 FTE manpower with respect to other executors:

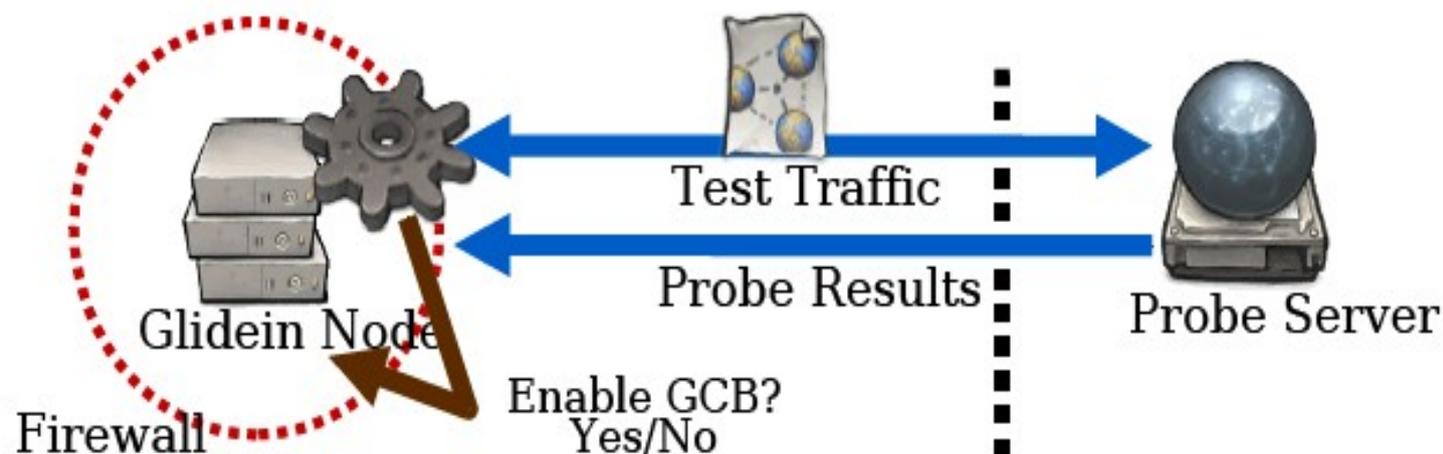
- Finished ~645K full simulation jobs this year
- WallTime Eff. ~ 87%
- Major losses are due to TRFERROR (ATLAS Software) and LCG-CP/LCG-CR
- Scaled up to ~5200 CPUs in parallel (No Jobs left thereafter)
- Performance loss due to network connectivity (GCB) – See next

Panda (OSG) and Dulcinea (Nordugrid) does data handling outside the executors

EGEE Executors needs to have a strong binding with the DDM

Network Probe

- Major concern – loss of connection between startd and the submit host
- Contact Condor servers @ Wisconsin to determine network information.
- Only enable GCB if needed.
- Source code is available!



Recently implemented on Cronus - Glideins

Recent event with CE's forgotten about the glideins:

How do we make sure that glideins are actually doing work and not wasting cycles ?

- *Must handle severed network connections*

New expressions allow Condor daemons to shutdown by themselves

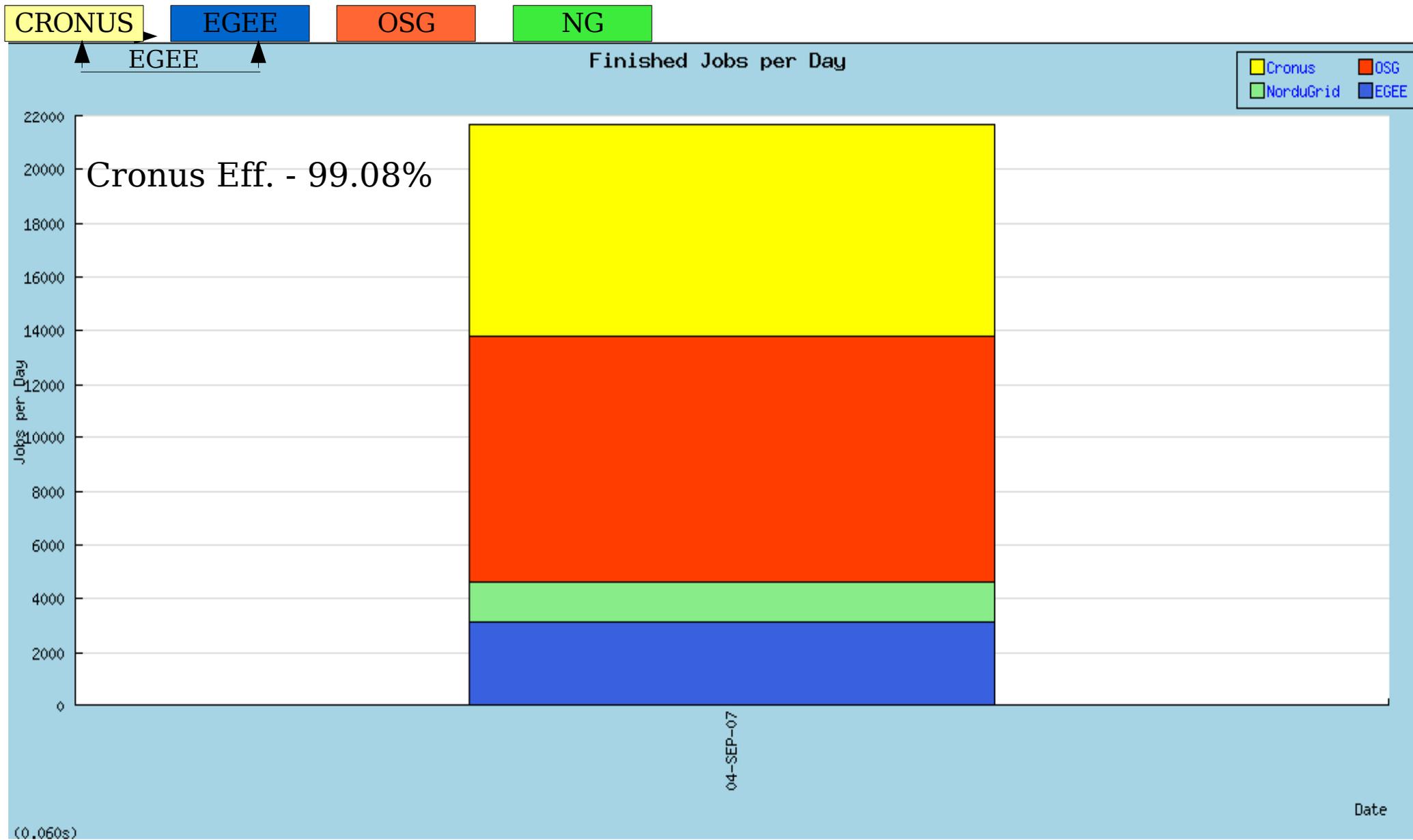
- without going through the Submit host or GCB (glidein proxies)

```
STARTD.DAEMON_SHUTDOWN = (State == "Claimed" || State == "Unclaimed")  
&& Activity == "Idle" && ((CurrentTime - EnteredCurrentActivity) > 1500)
```

Next Step – Provide the ratio of CPUTime/Walltime by the glideins

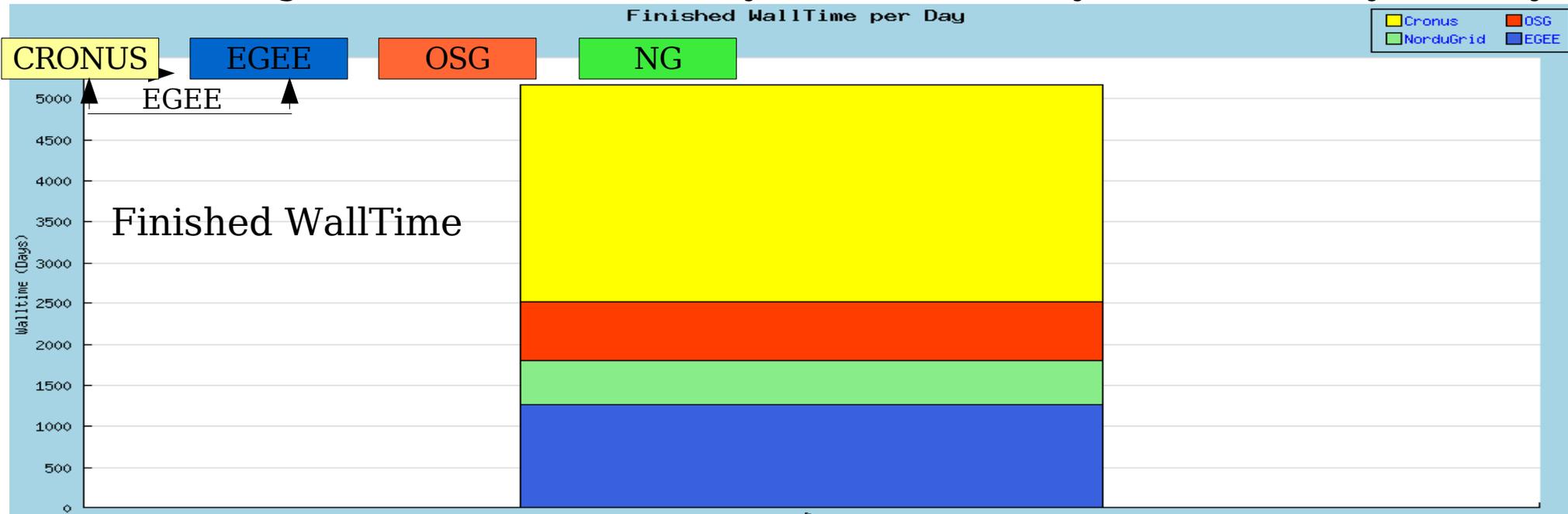
Cronus Operations

When Storage Services behave nicely- (ATLAS Prodsys results from yesterday)



Cronus Operations

When Storage Services behave nicely – (ATLAS Prodsys results from yesterday)



#	PROCESSINGDAY	EXECUTOR TYPE	FINJOBS	FINCPU	FINWALL	FAILJOBS	FAILCPU	FAILWALL	SUBMITTED	RUNNING	JOBEFF	WALLEFF
1	04-SEP-07	CondorG	395	11854481	12198168	264	1965456	2092349	4	75	59.93	85.35
2	04-SEP-07	Cronus	7856	213849839	229412633	471	1314652	2111103	1281	509	94.34	99.08
3	04-SEP-07	Dulcinea	1464	42429480	46131780	420	2711340	4316460	467	491	77.7	91.44
4	04-SEP-07	LCG-DQ	2763	91914792	97241512	574	772723	3722989	1059	1669	82.79	96.31
5	04-SEP-07	panda	9183	58484331	62172657	652	1577114	5670578	14403		93.37	91.64

Summary and Outlook

The pre-production mode using CRONUS was found to be successful !!!!

Plan to run two VMs (Virtual Machines on each CPU) in Cronus:

- One for User/Production Job
- Second to be used completely for the DATA managements

CMS glideinWMs can use Cronus information System across Grid Federations

For the first time interoperability between EGEE, OSG and NorduGrid using Client
Collaboration with Open Lab on Virtualizations of glideins, Network etc.

Successfully used more than ~5200 CPUs in parallel using EGEE infrastructure

Network severity can now be handled via GCB, KeepAlives from the WNs, etc.

DaemonCore automatically:

- decides when to shutdown in case of connection issues between Submit and WNs
- ensures the full usage of the CPUTime

Auto-Glideins ensures - a user uses his own glideins for his own jobs

EGEE Executors needs to have a strong binding with the DDM

**CRONUS - A fully virtual batch system which can be inter operable on the grid
Federations like: EGEE, OSG and NorduGrid**

Summary and Outlook



CRONUS – means time

According to Greek mythology

- CRONUS was the father of ZEUS who sent ATHENA to be in charge of ATLAS.