# Status and plans of the T3 in Geneva

**Szymon Gadomski
DPNC, University of Geneva**

**Swiss ATLAS Grid Working Group
Jan 7th, 2008**

# The cluster in Geneva (1)

# The cluster in Geneva (2)



12 worker
nodes
in 2005

21 in 2006

and 20 in 2007!

# The cluster in Geneva (3)

three nodes
for services
(grid, batch,
storage
abstraction)

direct line from CERN

power and network
cabling of worker nodes

# The hardware in numbers

- 61 computers to manage
  - 53 workers, 5 file servers, 3 service nodes
- 184 CPU cores in the workers
- 75 TB of disk storage
- can burn up to 30 kW (power supply specs)

A part of the hardware is already in production. More about the status later.

# The functionality we need

- our local cluster computing
  - log in and have an environment to work with ATLAS software, both offline and trigger
    - develop code, compile,
    - interact with ATLAS software repository at CERN
  - work with nightly releases of ATLAS software, normally not distributed off-site but visible on /afs
  - We also use CERN afs accounts.
    - We get all advantages of lxplus.cern.ch (examples of ATLAS software…)
    - without disadvantages (much more disk space, fewer users).
  - disk space visible as normal linux file systems (cd, cp, …)
  - use of final analysis tools, in particular ROOT
  - a convenient way to run batch jobs
- grid computing
  - tools to transfer data from CERN as well as from and to other Grid sites worldwide
  - ways to submit our jobs to other grid sites
  - a way for ATLAS colleagues, Swiss and other, to submit jobs to us
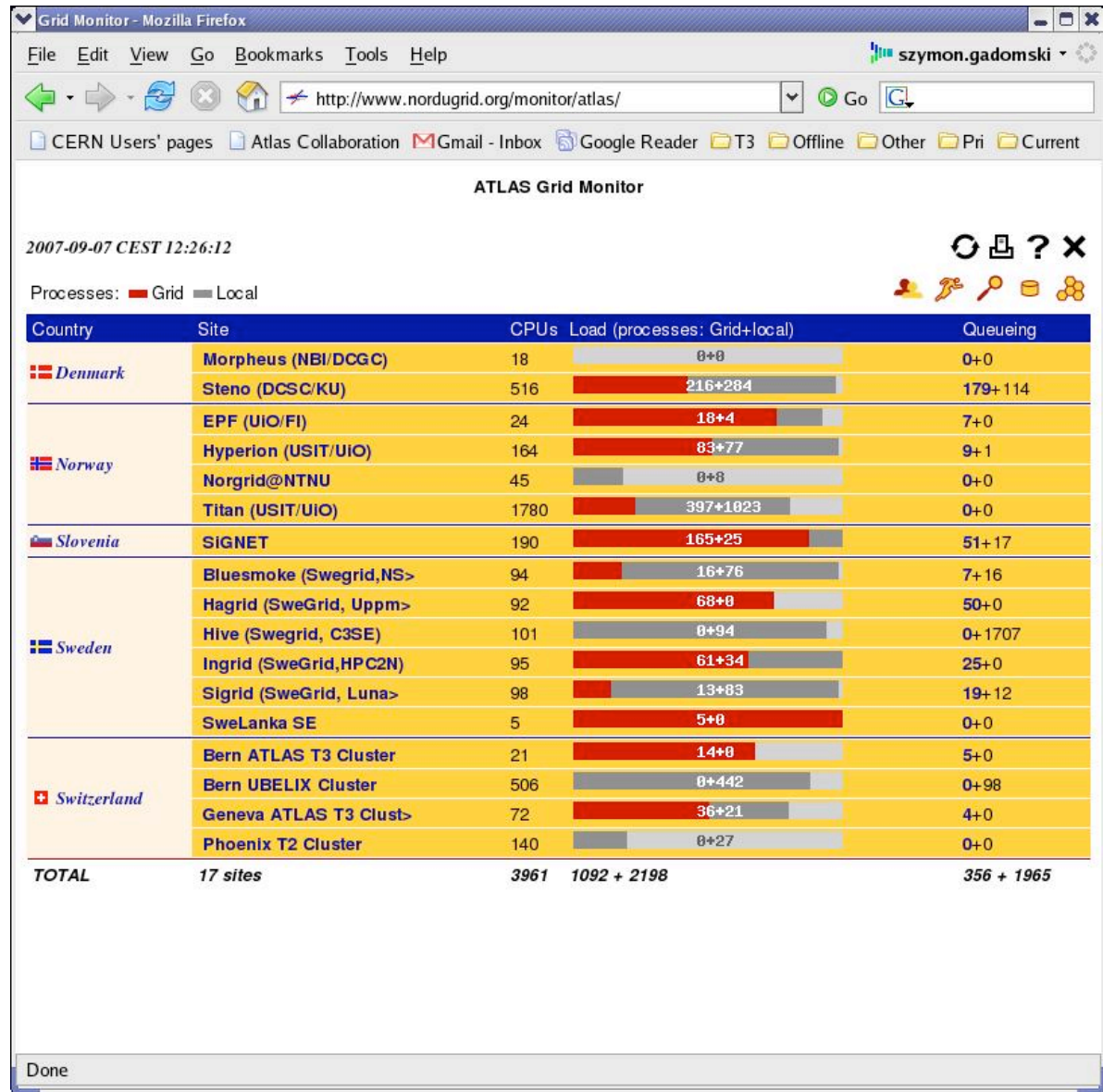
# Our system in production
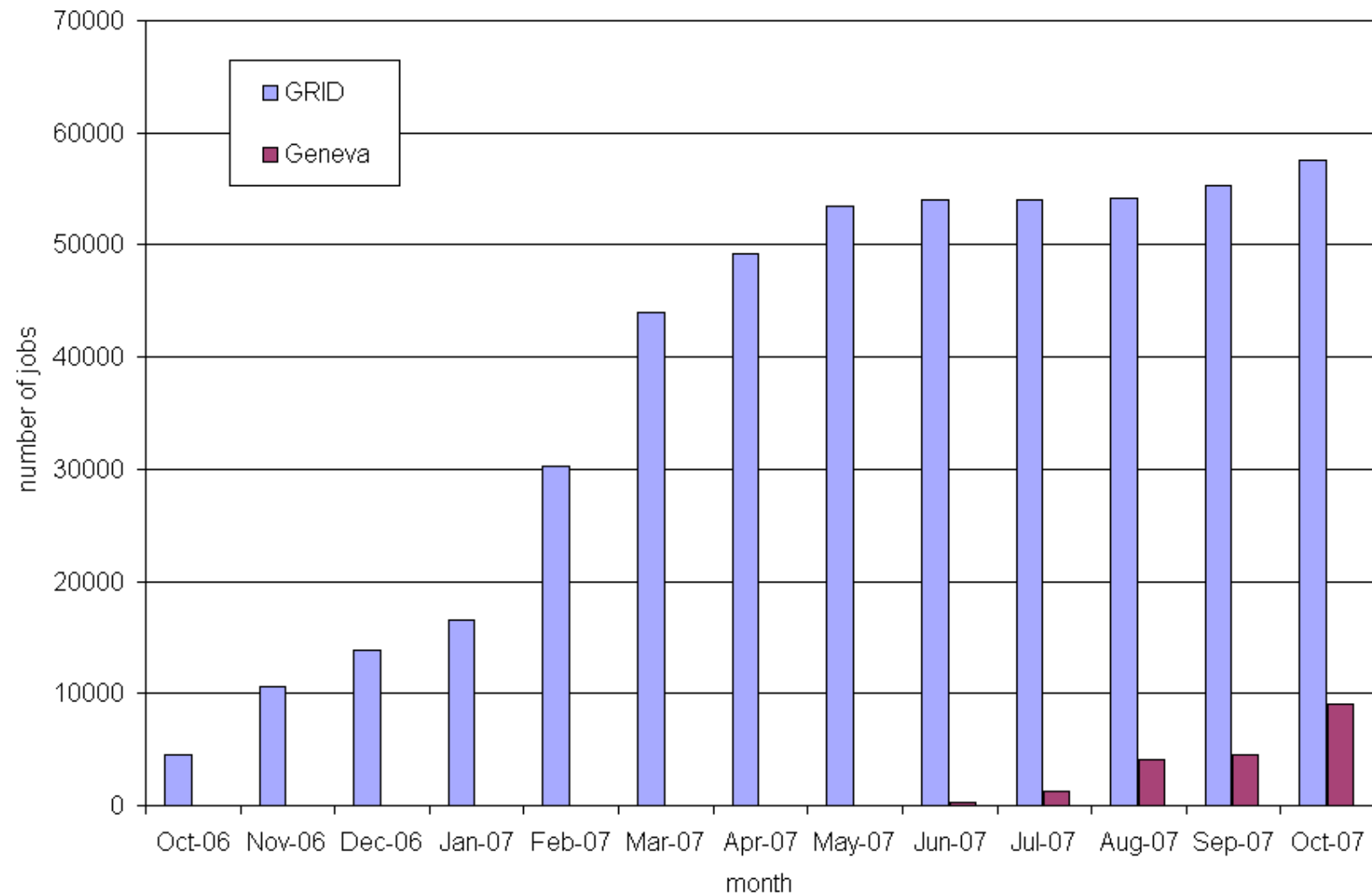
It changes. The status right now is:

- 41 TB of storage, available with NFS
  - 2 SunFire X4500 (Solaris), 1 old server (SLC3)

- 8 login machines

- 18 batch worker nodes
  - 18 SunFire X2200 ($2\times2$ core)

- 1 grid front-end machine
  - SunFire X4500 with SLC4 (not used as a file server anymore)

- The direct line from CERN goes to three login machines and 2 file servers.
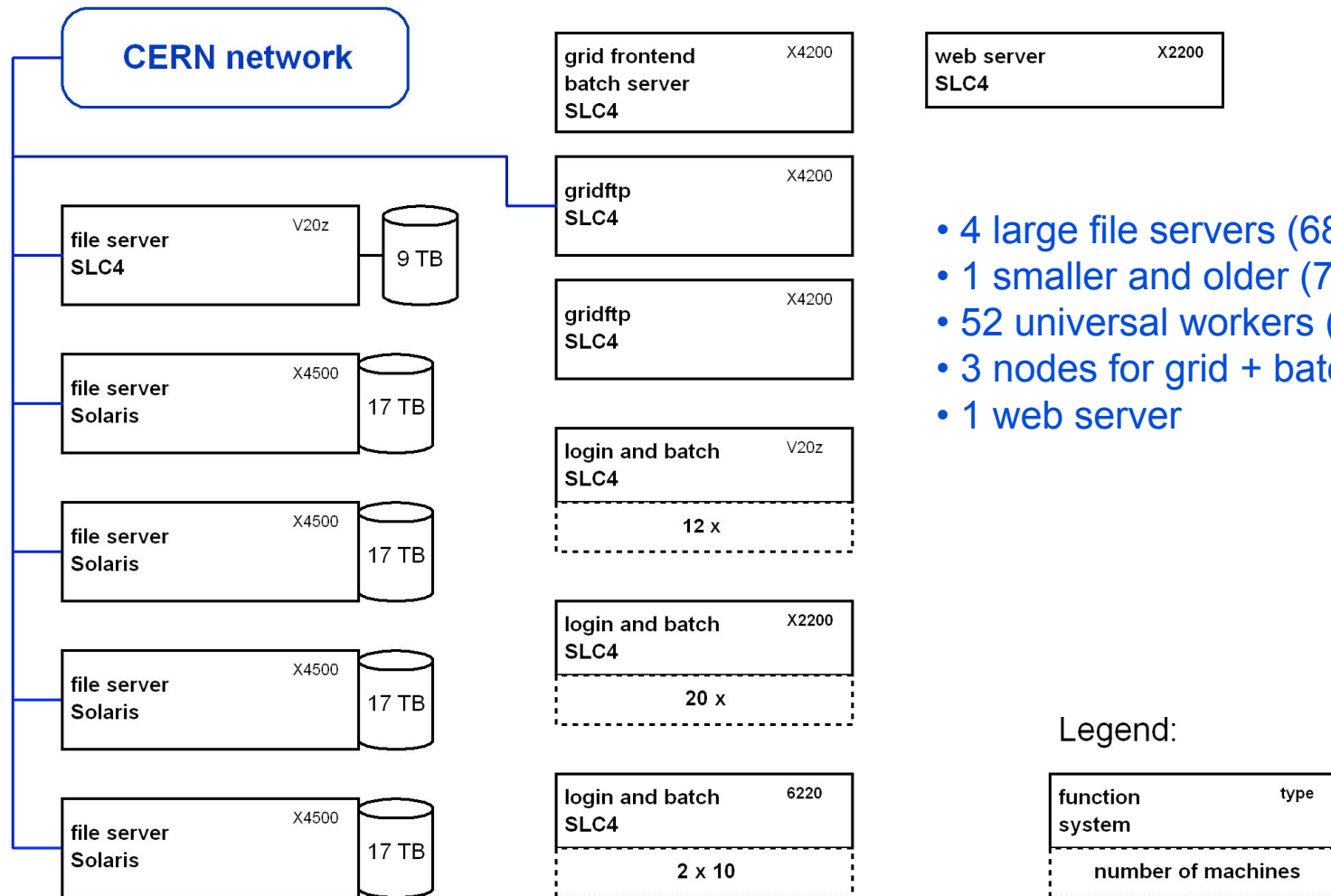
# Our system in the NorduGrid

- Geneva is in NorduGrid since 2005
- In company of Berne and the T2 in Manno

# Batch jobs statistics

# The system as planned for 1st data

**CERN network**

file server
SLC4 — V20z — 9 TB

file server
Solaris — X4500 — 17 TB

file server
Solaris — X4500 — 17 TB

file server
Solaris — X4500 — 17 TB

file server
Solaris — X4500 — 17 TB

grid frontend
batch server
SLC4 — X4200

gridftp
SLC4 — X4200

gridftp
SLC4 — X4200

login and batch
SLC4 — V20z
12 x

login and batch
SLC4 — X2200
20 x

login and batch
SLC4 — 6220
2 x 10

web server
SLC4 — X2200

- 4 large file servers (68 TB)
- 1 smaller and older (7 TB)
- 52 universal workers (184 cores)
- 3 nodes for grid + batch server
- 1 web server

Legend:

function — type
system
number of machines
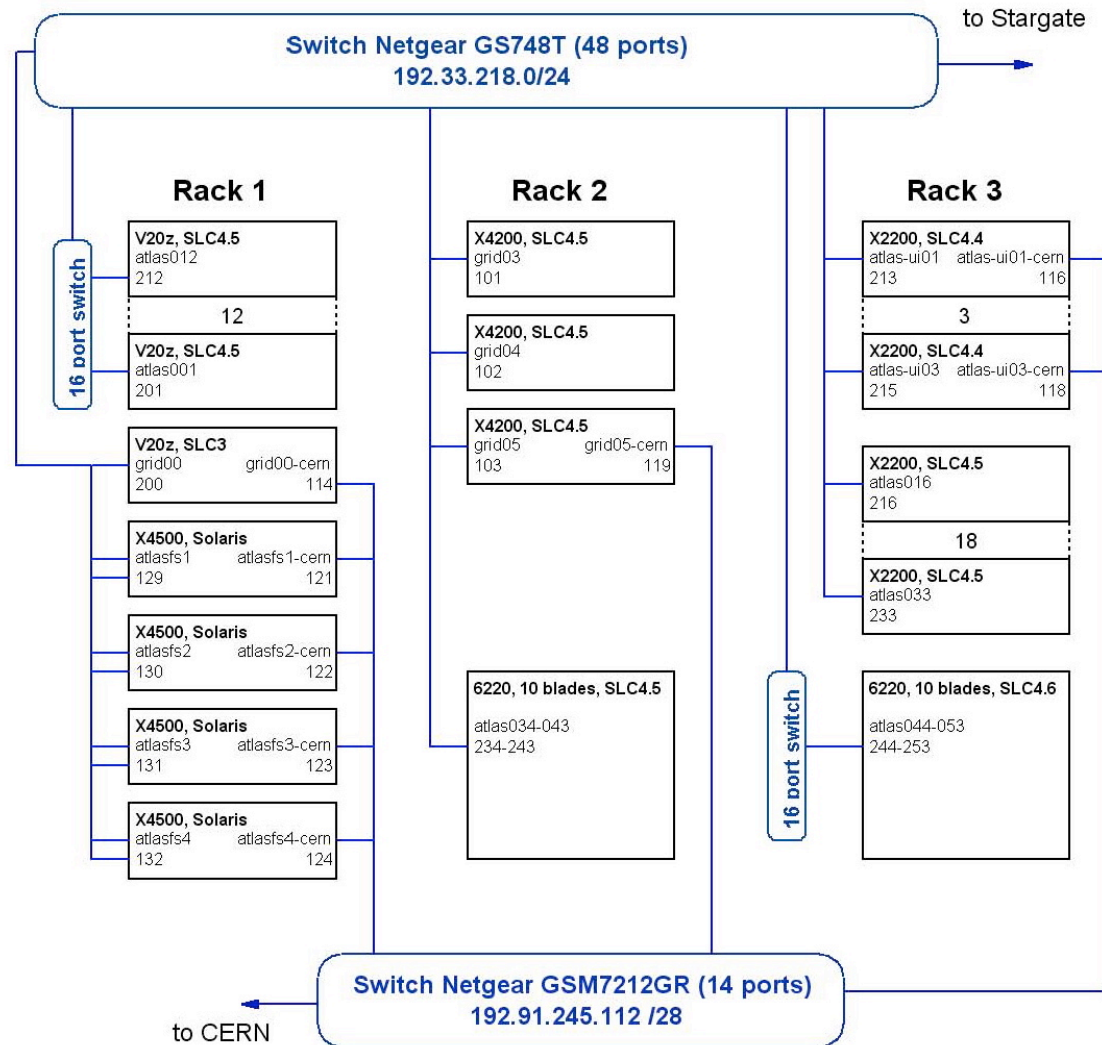
# (Re)installing everything. Why?

- three generations of hardware
  - 2005, SLC3, system getting obsolete
  - 2006, setup in record time, can't be propagated anyway
  - 2007, recently delivered, need to be installed
  - we want the same version of SLC on all the workers
- not only a pure batch job facility anymore
- grid and batch services on dedicated nodes, no longer on file servers
- file servers under Solaris
  - SunFire X4500 is designed for this, zfs is attractive
  - I can get help at the Uni to have this done. I still have the remaining 57 machines to manage.

# Towards the final system

- **What is done already**
  - Experience with the system in production
    - File servers used also as batch servers and grid front-ends
      - Both crashed last June! File servers will not do anything else anymore.
    - Power was lost three times last Summer.
  - Stable operation since September
    - >30 users rely on the cluster for daily work, interactive + batch
    - NorduGrid jobs are running ~all the time
  - All the new hardware is in the racks
  - The networks are done
  - Work on power supply is done (during the Christmas break).
  - The new "final" batch system is in place
    - Batch server on a SunFire X4200 machine
      - TORQUE and pbs_sched
    - 10 worker nodes on new blades (Sun Blade 6220)
    - Independent hardware, see the same file servers as the system in production

# Networking for the final system

In place since last November.

# Towards the final system (2)

- **Next to do: finish "sysadmin" work**
  - Change scheduler to Maui on the new batch server
    - learn how to configure it…
  - Get machine certificates
    - Can we use GlobalSign, rather than SwissSign? CHIPP gives us this option now! Certificates are valid for three years…
  - Setup NorduGrid on a new front-end machine
  - Operate the new system in parallel with the one currently in production for some time, check that all is stable.
  - Setup a web server (my 1st…)
  - Setup monitoring of network and CPU use.
    - learn about Ganglia
  - Install all remaining hardware
    - 10 remaining blades, 5 currently unused machines from 2005
    - two more X4200 to be used for gridftp
  - Migrate worker nodes and login machines to the new system.
  - Retire grid02 as grid+batch

# Towards the final system (3)

- **Longer term things to do**
  - **Sysadmin follow-up**
    - reduce dependence on grid00 (SLC3)
    - other optimizations and fixes
  - **Data transfer exercises**
    - **Official data path: CERN > FZK > CSCS > (BE+GE)**
      - **Data can be pushed up to CSCS.**
      - **GE + BE can pull it from there**
    - **Backup path: CERN > GE > BE**
      - **Pull only. Small quantities of data.**
    - **Consider GE as a DDM service site**
      - **Push data like this: CERN > FZK > (CSCS + GE)**
      - **Pull data GE > BE**
      - **We need a file catalogue listing data in Geneva. Can it be at CSCS?**
  - **Tools to submit many batch jobs at once**
    - **GridPilot?**
    - **Ganga**
  - **Interactive use of many machines**
    - **PROOF**

# Summary

- **The ATLAS cluster in Geneva is a large Tier 3**
  - **will have 184 worker's CPU cores and 75 TB soon**
- **A part of the system is in production**
  - **a Grid site since 2005, runs ATLAS simulation like a Tier 2, plan to continue that.**
  - **since Spring in constant interactive use by the Geneva group and friends, plan to continue and to develop further**
- **Expect to be busy installing machines and services for another few weeks**
  - **a long of history, three generations of hardware**
  - **need a uniform system with a rational architecture**
  - **we have learned from experience**
- **After that, I can start participating in data transfer exercises.**
  - **I will try to get some help for that in the Geneva group.**