# Experience with the gLite Workload Management System in ATLAS Monte Carlo Production on LCG

*Simone Campana, CERN*

*David Rebatto, INFN*

*Andrea Sciabà, CERN*

*CHEP2007 Victoria (Canada)*

**www.eu-egee.org**

Information Society
and Media

**Enabling Grids for E-sciencE**

- **LHC experiments need to generate huge amounts of simulated data**
  - Validate the computing and data model
  - Test the complete software suite
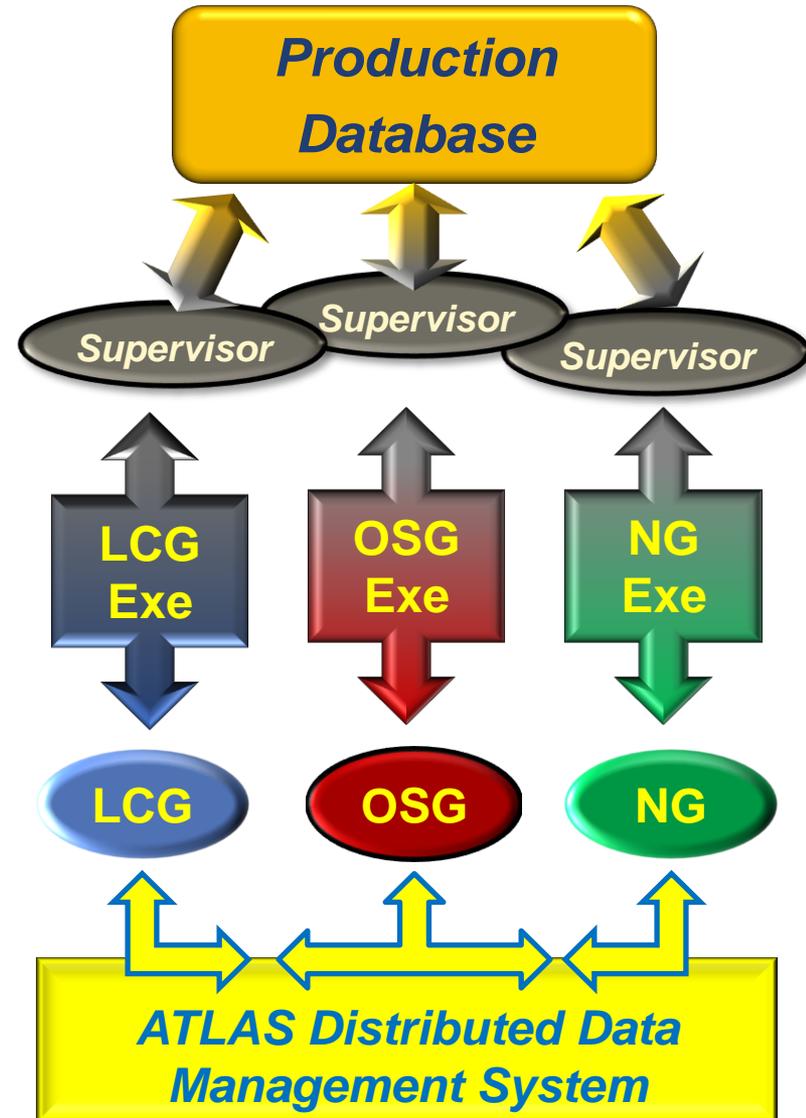  - Develop physics data analysis

- **ATLAS currently benefits of ~10 KSI2K of CPU power and more than one PB of disk space distributed around more than 60 sites**
  - Decentralization and sharing of computing resources
  - Different computing facilities are organized in a hierarchical structure (T0, T1, T2)
    - distinct roles at different levels.
  - In ATLAS Computing Model, MC Production is run at T2s
    - … but currently also at T1s

- **Specific tools have been developed by each experiment to manage the production workflow**
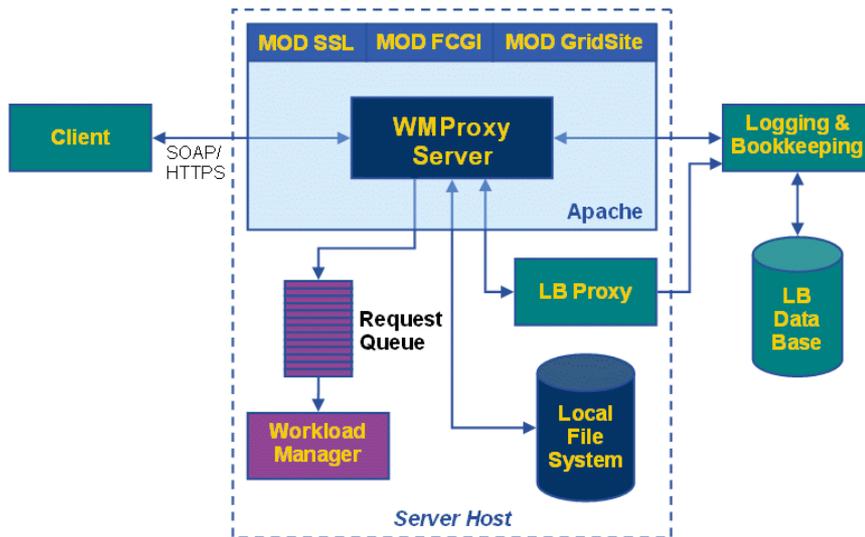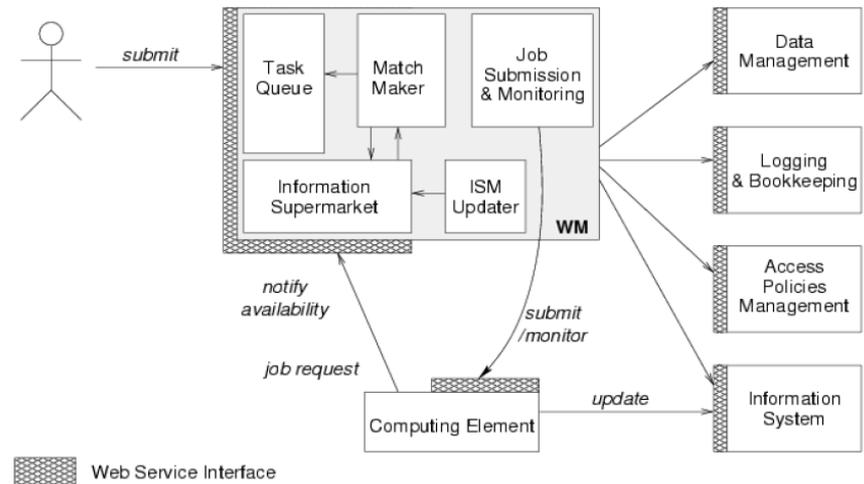  - For ATLAS this is the Production System (Prodsys)

EGEE

OSG

Nordugrid

**eGee**

- **A central database**
    - holds grid neutral definitions of tasks and jobs, together with job states

- **A "supervisor"(also Grid neutral)**
    - pulls jobs from the central database
    - submits jobs to the Grid
    - monitors jobs and checks their outcome

- **An "executor" layer acting as interface to the Grid middleware**
    - EGEE/WLCG
        - <u>Lexor using the gLite WMS</u>
            - *<u>Was the LCG-RB before</u>*
        - Condor-G direct submission
        - CRONUS (Condor glide-ins)



*Production Database*

*Supervisor* *Supervisor* *Supervisor*

**LCG Exe** **OSG Exe** **NG Exe**

**LCG** **OSG** **NG**

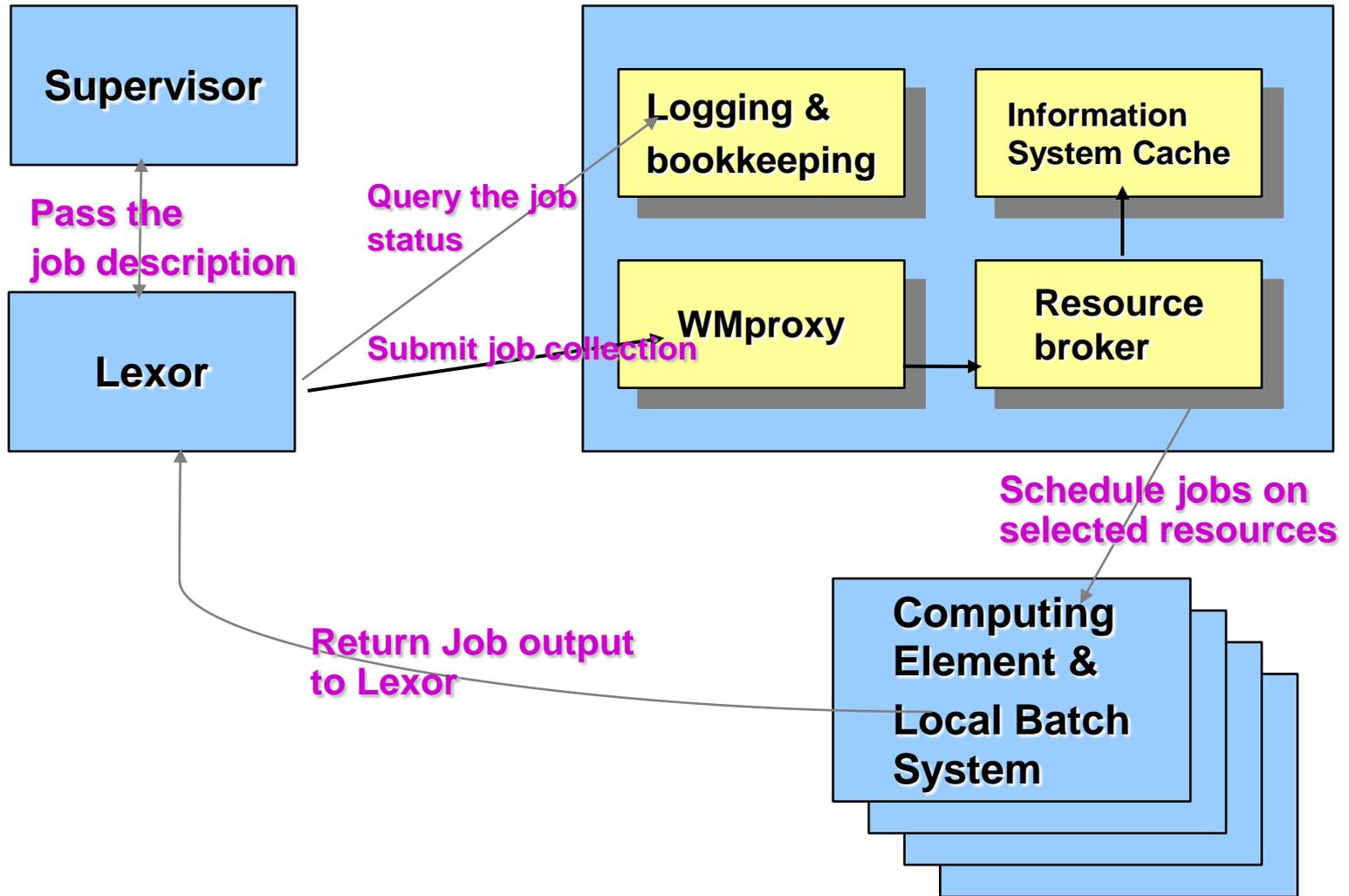*ATLAS Distributed Data Management System*

- **The service to submit and manage jobs**
    - **Task queue**: holds jobs not yet dispatched
    - **Information SuperMarket**: caches all information about Grid resources
    - **Match Maker**: selects the best resource for each job
    - **Job Submission & Monitoring**
    - Interacts with Data Management, Logging & Bookkeeping, etc.



- **WMProxy service optimizes job management and stands between the user and the real WMS**
    - Service Oriented Architecture (SOA) compliant
        - Implemented as a SOAP Web service
    - Validates, converts and prepares jobs and sends them to the WM
    - Interacts with the L&B via LBProxy (a state storage of active jobs)
    - Implements most new features

- **The gLite WMS offers several advantages over the old LCG WMS**
  - Bulk submission
    - Collections: sets of <u>independent</u> jobs
    - New, much more reliable implementation as a compound job submission
  - Job sandboxes
    - Shared input sandboxes for a collection
    - Download/upload of sandboxes via GridFTP, http, https
  - Faster match-making
    - "bulk" matchmaking and ranking for collections
  - Internal task queue
    - If a job cannot match right away it is kept for some time until it matches
  - Resubmission of failed jobs
    - a job is resubmitted right away after a middleware/infrastructure-related failure
    - greatly improves the job success rates
  - A limiter mechanism which prevents submission of new jobs if the load exceeds a certain threshold
    - Leads to "artificial", but desired, limitations of the job submission rate
    - Improves the stability of the system
  - Last but not least, the gLite WMS is actively developed and maintained, while the LCG RB is "frozen"

**eGee**

**Supervisor**

**Pass the
job description**

**Lexor**

**Query the job
status**

**Submit job collection**

**Logging &
bookkeeping**

**Information
System Cache**

**WMproxy**

**Resource
broker**

**Schedule jobs on
selected resources**

**Return Job output
to Lexor**

**Computing
Element &**

**Local Batch
System**

- **The gLite WMS has been one source of inefficiency of Lexor in the past**

- **Lexor is using gLite WMS since Summer 2006**
  - gLite WMS not really ready for production
  - A lot of manual intervention was needed
    - Period September-November 2006 has been very critical
- **Main problems**
  - Bulk submission not completely reliable
    - jobs remaining "stale" (in the same state) forever: "zombie" jobs
  - Memory usage growth
    - linear it time under continuous job submission

- **New activity of testing and debugging of gLite WMS started in January 2007**
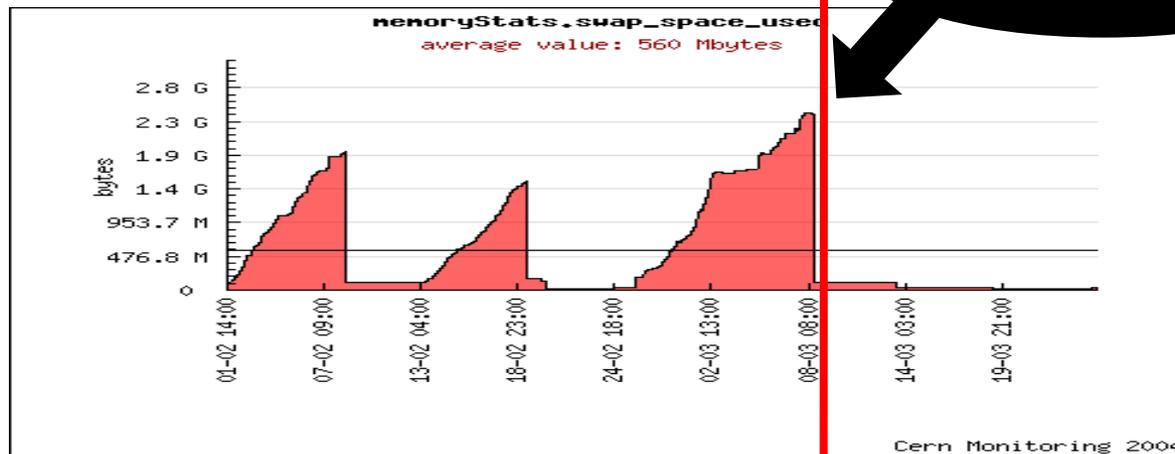  - Starting from Experiment' s requirements and acceptance criteria

| | CMS | ATLAS |
|---|---|---|
| **Performance** | | |
| 2007 | 50K jobs/day | 20K production jobs/day + analysis load |
| 2008 | 200K jobs/day (120K to EGEE, 80K to OSG)<br><br>Using <10 WMS entry points | 100K jobs/day through the WMS;<br><br>Using <10 WMS entry points |
| **Stability** | | |
| | | <1 restart of WMS or LB every month under load |

**eGee**

- **Based on the experiment requirements, some criteria have been defined to decide if the gLite WMS satisfies the requirements**
  - At least 10000 jobs/day submitted for at least five days
  - No service restart required for any WMS component
  - The WMS performance should not show any degradation during this period
  - The number of zombie jobs should be less than 0.5% of the total

**Enabling Grids for E-sciencE**

- **The testing of the gLite WMS is mainly done by the Experiment Integration and Support team of WLCG**
  - Collaboration between Experiment Integration Support Team, JRA1 (EGEE developers), SA1 (EGEE operations), SA3 (EGEE integration and testing)
  - Bugs discovered, fixed and patched bypassing normal certification procedures
    - WMSes continuously tested, patched and re-deployed
    - Pragmatic approach: very quick turnover
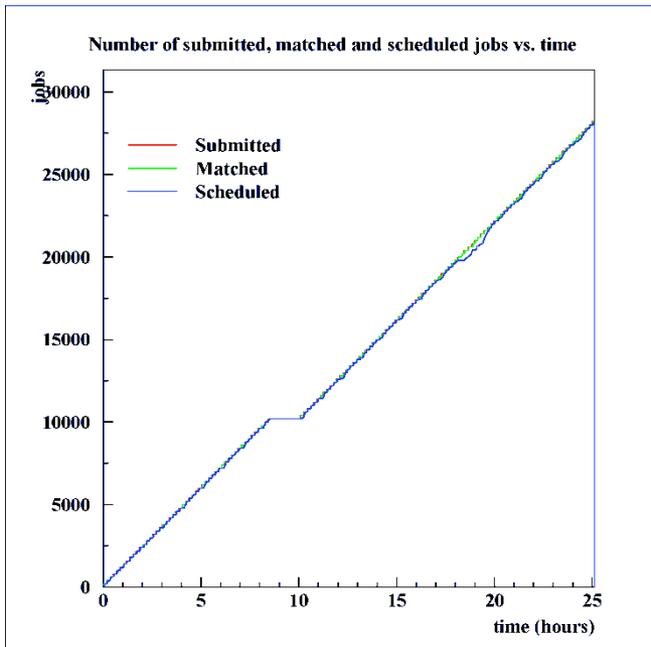  - Huge improvements in stability and performance

**Enabling Grids for E-sciencE**

- **Memory usage**
  - Grows linearly in the gLite WMS 3.0
    - Hard to maintain the service: restarts, reboots …
  - OK in gLite WMS 3.1

**Switch to 3.1**



- **The problem of "zombie" jobs have been identified**
  - Collection handling via Condor DAGMAN

- **DAGMAN has been removed for collection handling**
  - Still there for Acyclic Job Diagrams
  - Collections are handled via a native mechanism
  - The latest test have shown NO jobs stale at all

**Enabling Grids for E-sciencE**

- **115000 jobs submitted in 7 days**
  - ~16000 jobs/day well exceeding acceptance criteria
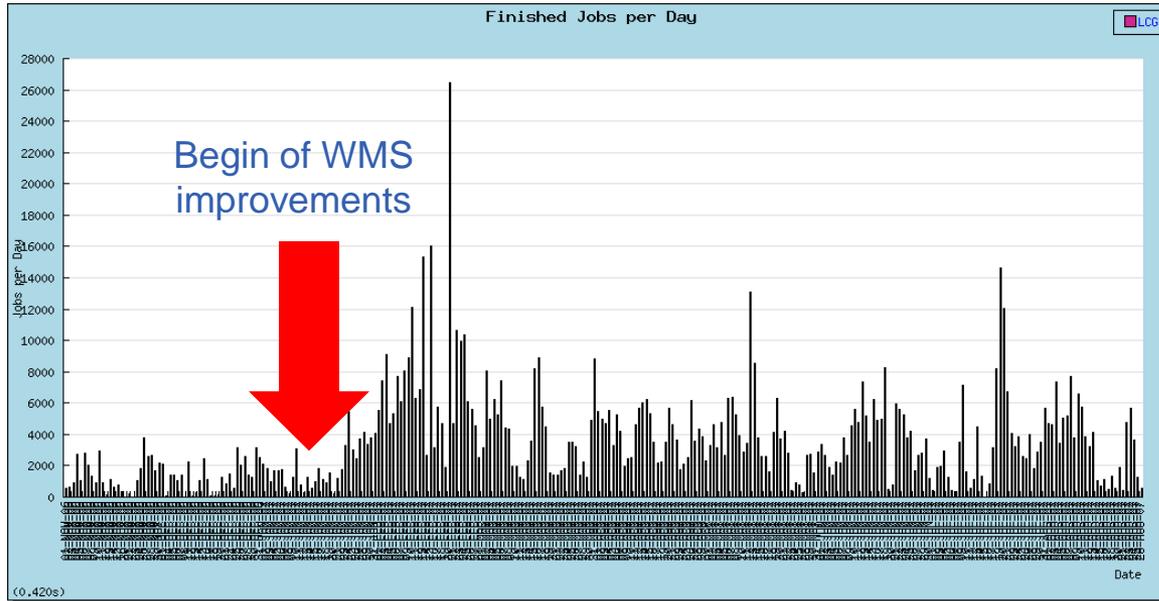  - The "limiter" prevented submission when load was very high (>12)
- **All jobs were processed normally but for 320**
  - ~0.3% of jobs with problems, well below the required threshold
  - Recoverable using a proper command by the user



Number of submitted, matched and scheduled jobs vs. time



Number of jobs in each status vs. time

- **The WMS dispatched jobs to computing elements with no noticeable delay**
- **Acceptance tests were passed**

Finished Jobs per Day

Begin of WMS improvements

Big ramp-up after WMS improvements

Reached 20000 jobs/day

Oscillatory behavior
(also depending on the type of production activity being run)

Job efficiency:
Roughly 60%

See later slides for error categorization and impact



Jobs per Day

**Enabling Grids for E-sciencE**



WCT per day:

Somehow different from Job/Day distribution

Different type of jobs in different periods



WCT efficiency:

Much higher than job efficiency

We focsed in reducing error with large waste of WCT

LCG Jobs Errors 1/1/2007 - 27/8/2007

stageOut

Executor

gLiteWMS

StageIN

StageIN

ATLAS SW

Legend:
- NONE
- WRAPLCG_STAGEOUT_LFCMKDIR
- WRAPLCG_WNCHECK_SWENV
- EXELEXOR_GLITE_WMS
- WRAPLCG_WNCHECK_SWMISS
- WRAPLCG_STAGEOUT_MISS
- WRAPLCG_STAGEOUT_LCGCR
- EXELEXOR_GETOUT_OBJLOAD
- WRAPLCG_WNCHECK_PROXY
- EXELEXOR_GLITE_MAXRETRYCOUNT
- EXELEXOR_GLITE_PROXYEXPIRED
- WRAPLCG_STAGEIN_NOREPLICAS
- EXELEXOR_UNSPEC
- WRAPLCG_STAGEIN_SIZE
- WRAPLCG_STAGEIN_MD5SUM
- EXE_CONTROL_STALESTATUS
- WRAPLCG_JTINST_PACUNTAR
- WRAPLCG_WNCHECK_SWREAD
- WRAPLCG_STAGEIN_LCGCP
- WRAPLCG_JTINST_PACDWNLD
- TRFERROR
- EXELEXOR_GETOUT_UNTAR
- WRAPLCG_JTINST_GET
- EXELEXOR
- WRAPLCG_STAGEIN_FAILLOOKUP
- WRAPLCG_STAGEIN_MISSINGIN
- WRAPLCG_STAGEOUT_POOLGUID

(0.179s)

**gLite WMS: ~22%**          **Data Management: 36%**          **ATLAS SW: 8%**

**eGee**

Enabling Grids for E-sciencE



LCG Jobs Errors 1/7/2007 - 27/8/2007

gLiteWMS

Executor

StageIN

gLiteWMS

ATLAS SW

7.3%
0.1%
0.2%
11.7%
3.2%
1.9%
0.1%
1.3%
2.2%
0.0%
.0%
.2%
47.8%
7.7%
0.2%
1.0%
11.4%
0.0%
0.0%
0.0%

(0.162s)

Legend:
- NONE
- WRAPLCG_WNCHECK_SWENV
- EXELEXOR_GLITE_WMS
- WRAPLCG_WNCHECK_SWMISS
- WRAPLCG_STAGEOUT_MISS
- WRAPLCG_STAGEOUT_LCGCR
- EXELEXOR_GLITE_MAXRETRYCOUNT
- EXELEXOR_GETOUT_OBJLOAD
- EXELEXOR_GLITE_PROXYEXPIRED
- WRAPLCG_WNCHECK_PROXY
- EXE_CONTROL_STALESTATUS
- WRAPLCG_STAGEIN_MD5SUM
- WRAPLCG_STAGEIN_SIZE
- WRAPLCG_STAGEIN_LCGCP
- WRAPLCG_WNCHECK_SWREAD
- WRAPLCG_JTINST_PACDWNLD
- WRAPLCG_JTINST_PACUNTAR
- RFERROR
- WRAPLCG_STAGEIN_MISSINGIN
- WRAPLCG_STAGEIN_FAILLOOKUP
- EXELEXOR_GETOUT_UNTAR
- WRAPLCG_JTINST_GET
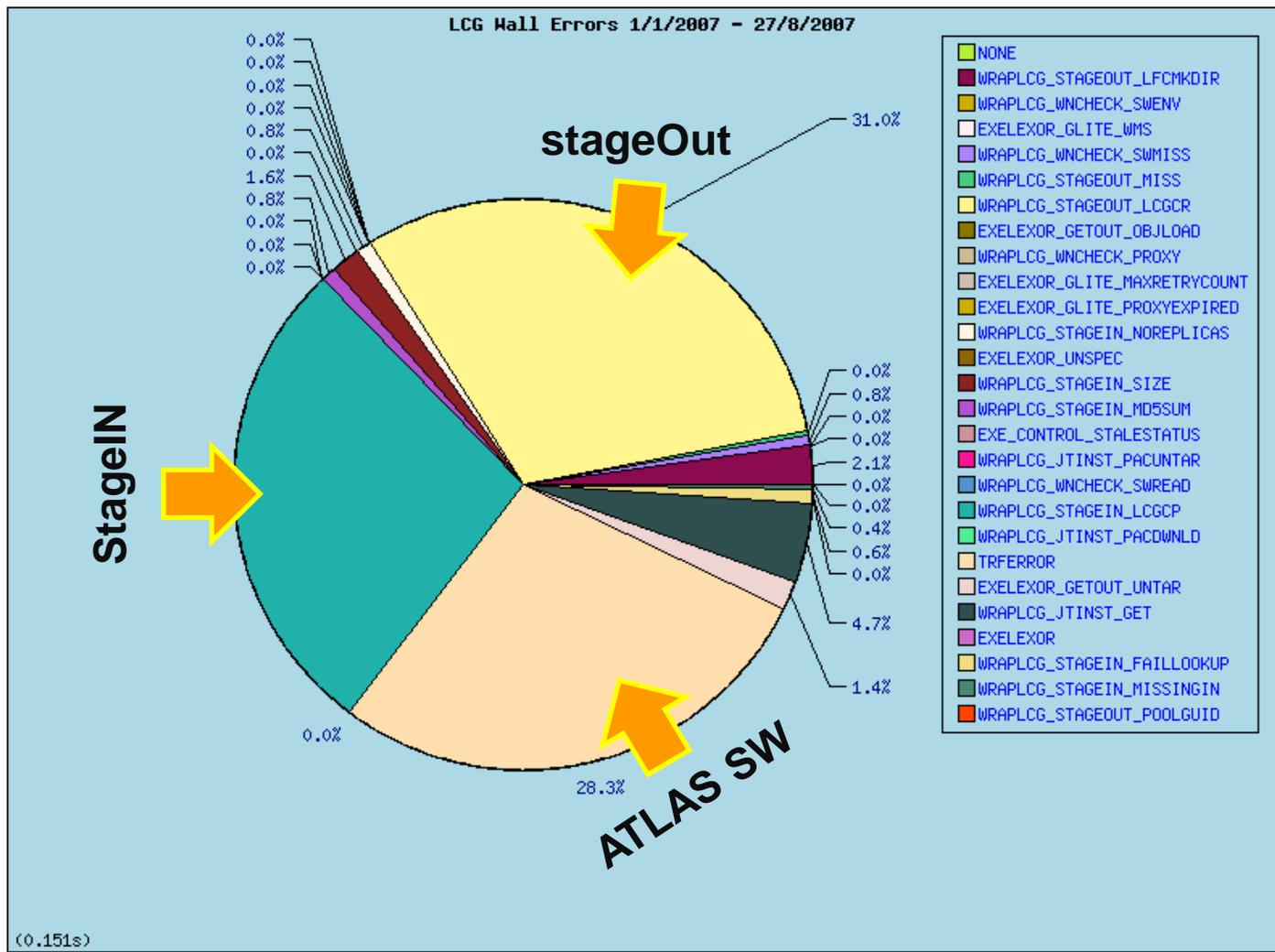- WRAPLCG_STAGEOUT_POOLGUID

**gLite WMS: ~13%**          **Data Management: 47%**          **ATLAS SW: 11%**

gLite WMS category includes also site specific issues and  problematic job distribution (with subsequent proxy expiration).

**Enabling Grids for E-sciencE**



LCG Wall Errors 1/1/2007 - 27/8/2007

stageOut — 31.0%

StageIN

ATLAS SW

0.0%
0.0%
0.0%
0.0%
0.8%
0.0%
1.6%
0.8%
0.0%
0.0%
0.0%

0.0%
0.8%
0.0%
2.1%
0.0%
0.0%
0.4%
0.6%
0.0%
4.7%
1.4%

0.0%
28.3%

(0.151s)

Legend:
- NONE
- WRAPLCG_STAGEOUT_LFCMKDIR
- WRAPLCG_WNCHECK_SWENV
- EXELEXOR_GLITE_WMS
- WRAPLCG_WNCHECK_SWMISS
- WRAPLCG_STAGEOUT_MISS
- WRAPLCG_STAGEOUT_LCGCR
- EXELEXOR_GETOUT_OBJLOAD
- WRAPLCG_WNCHECK_PROXY
- EXELEXOR_GLITE_MAXRETRYCOUNT
- EXELEXOR_GLITE_PROXYEXPIRED
- WRAPLCG_STAGEIN_NOREPLICAS
- EXELEXOR_UNSPEC
- WRAPLCG_STAGEIN_SIZE
- WRAPLCG_STAGEIN_MD5SUM
- EXE_CONTROL_STALESTATUS
- WRAPLCG_JTINST_PACUNTAR
- WRAPLCG_WNCHECK_SWREAD
- WRAPLCG_STAGEIN_LCGCP
- WRAPLCG_JTINST_PACDWNLD
- TRFERROR
- EXELEXOR_GETOUT_UNTAR
- WRAPLCG_JTINST_GET
- EXELEXOR
- WRAPLCG_STAGEIN_FAILLOOKUP
- WRAPLCG_STAGEIN_MISSINGIN
- WRAPLCG_STAGEOUT_POOLGUID

**gLite WMS: negligible    Data Management: ~60%    ATLAS SW: 28%**

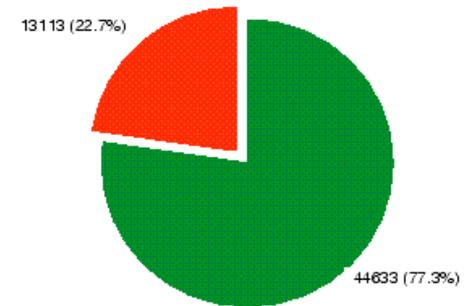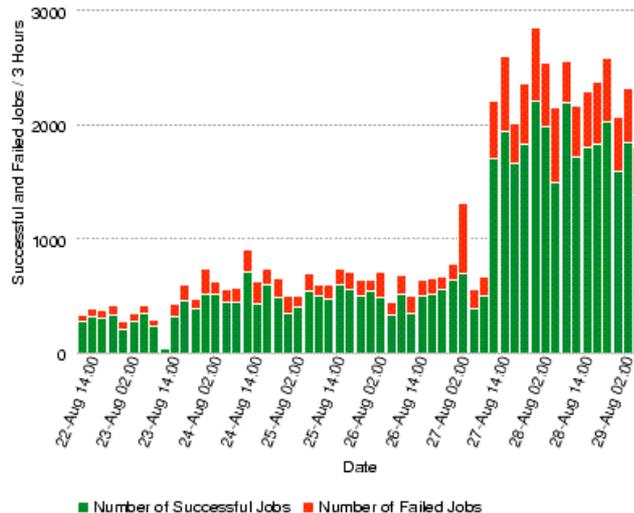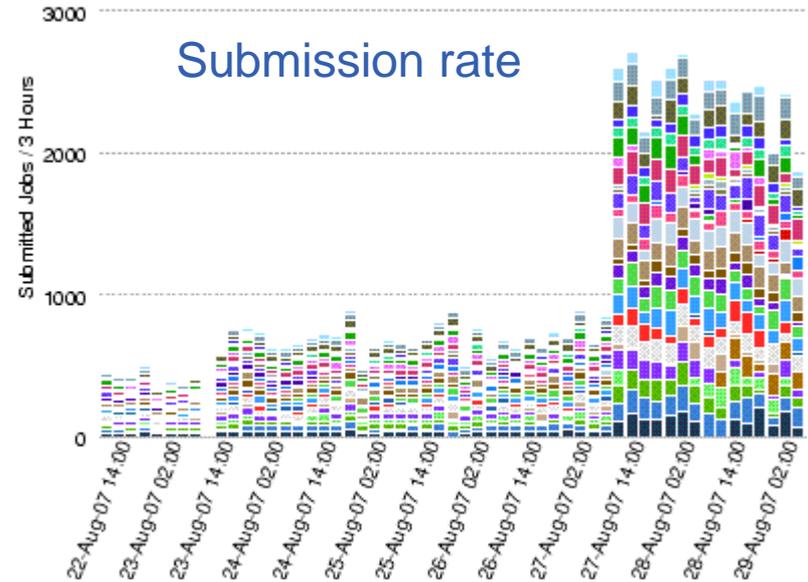| | Lexor | ATLAS Production (tot) |
|---|---|---|
| Finished Jobs | 942439 | 4317243 |
| Failed Jobs | 746252 | 2208135 |
| **Job Efficiency** | **55.8** | **66.2** |
| Finished WCT | 17567932826 | 79018496695 |
| Failed WCT | 2550170383 | 11910469832 |
| **WCT Efficiency** | **87.32** | **86.9** |

Lexor Job Efficiency is somewhat lower than ATLAS overall job efficiency

**Automatic resubmission can cope with this**

In terms of WCT (which means how much resources you are wasting for failed jobs) Lexor is quite efficient (more than ATLAS overall average)

- **CMS supports submission of analysis jobs via WMS**
  - Using two WMS instances at CERN with the latest certified release
  - For CSA07 the goal is to submit at least 50000 jobs/day via WMS
  - The Job Robot (a load generator simulating analysis jobs) is successfully submitting more than 20000 jobs/day to two WMS

Submission rate

Success rate

**Enabling Grids for E-sciencE**

- **Most reliability problems in gLite WMS are understood**
  - A few minor issues still being investigated

- **Several features of gLite WMS still to be considered**
  - Job Perusal: real time access to job stderr and stdout
  - Reputability Ranking: exclusion of resources causing large job failures
  - Job Provenance in Logging and Bookkeeping

- **Some improvements being discussed with developers**
  - e.g. stochastic ranking expression

- **The advantages compared to the LCG Resource Broker are very significant**

- **The achieved improvements had a big impact on many production activities**
  - e.g. for the ATLAS Monte Carlo production

- **All the LHC experiments are ready to use it**
  - Either they are already using it, or have finished the testing phase

- **Experimental services approach has shows to be extremely effective**
  - Adopted also for other components, i.e. CondorCE and CREAM