

Biomed data challenge : Comparison from end user statistics against RBs statistics

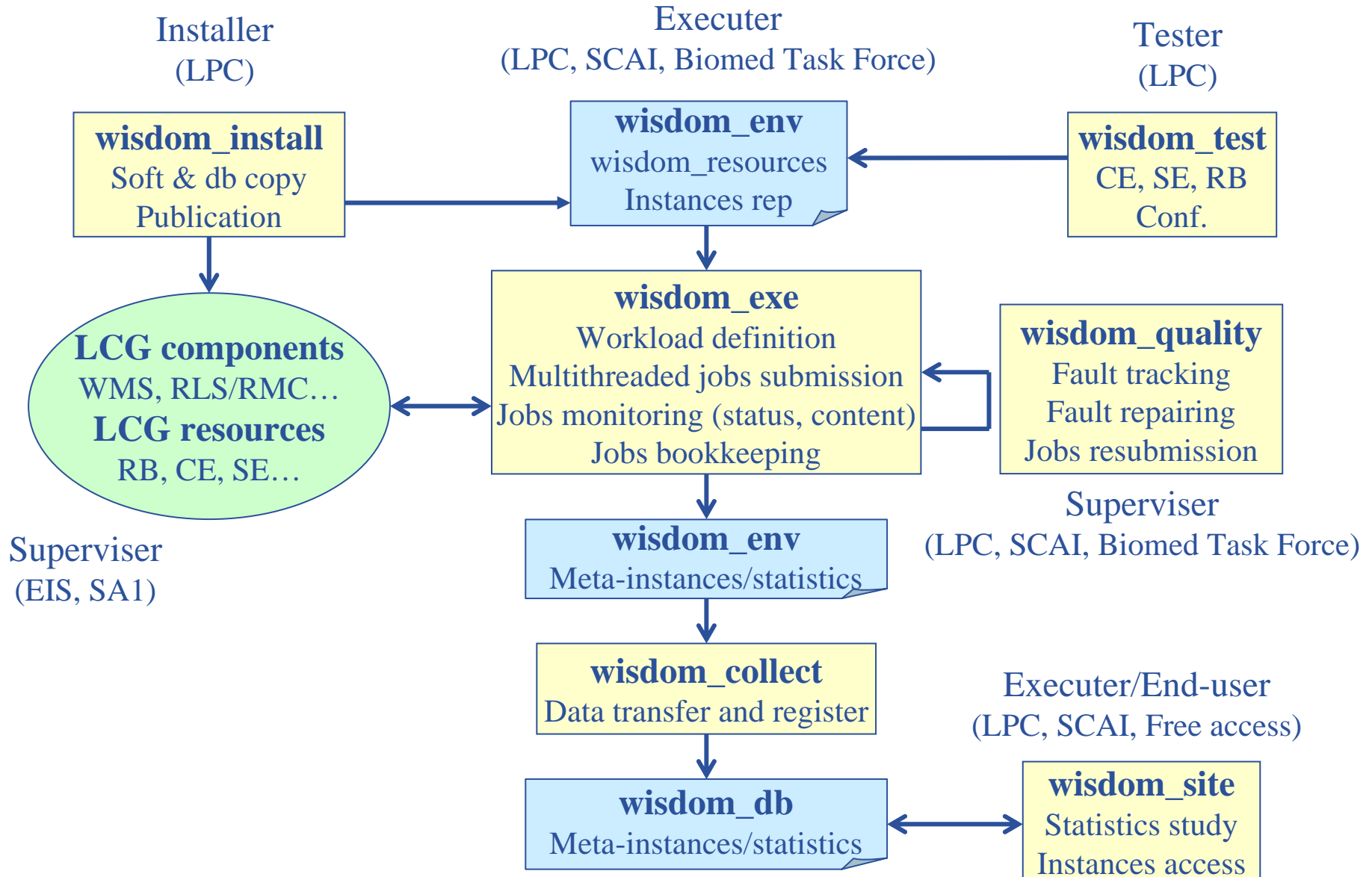
Nicolas Jacq

LPC, IN2P3/CNRS, France

Credit : Geneviève Romier, Cyril Lorphelin

- **End user statistics**
 - Objectives of the statistics collect
 - Nature of the collected statistics
- **Comparison with statistics gathered from RBs**
 - Nature of the collected statistics
 - Statistics comparison
- **Conclusion**

- **Grid objective of the data challenge :**
 - Producing a large amount of data in a limited time with a minimal human cost during the data challenge
- **Objectives of the statistics collect**
 - Prove the performance of the grid for our application
 - Prove the large-scale deployment on the grid
 - Test the reliability of the grid
- **Need dedicated indicators to measure performance**
 - But the detailed statistics collect was not a priority
- **Need a specific environment for data challenge statistics**



- **edg-job-status**
 - RB
 - CE
 - Status
 - Times
- **UI environnement**
 - jobID
 - Submission time
 - Failures during the submission
 - Failures during the process (status, get output)
 - Failures in the job output content
- **Job script**
 - CPU time
 - Data transfer time
 - Transferred data size
 - SE and PFN where data is stored

- **Test instances**
 - During the DC
- **Failed instances / jobs**
 - WISDOM or grid failures
 - Abandoned instances or jobs (Human failure, RB overload, license server...)
- **Problems to download final status/time information for long instance**
 - Download by WISDOM is done at the end of the instance => data can be lost or deleted
 - Download by JRA2 is done 2 or 3 days after
- **RetryCount attribute**
 - Only job information of the last node is registered
 - JRA2 also
 - Remark : how to count retry time (=> currently submitted time)
- **Transfer data size and time**
 - Input data : Lcg-cp/globus-url-copy
 - Output data : lcg-cr/lcg-rep/globus-url-copy

- **The aim was to check the user data**
 - Not all RBs
 - We didn't really know what information JRA2 had access
- **But there is also a real added value**
 - Distorsion of the user statistics
 - Aborted job reasons
 - Checking by the duration of the jobs (> 3h)
- **Even if JRA2 statistics depend of the RBs**
 - Unclassified jobs, disk crash...

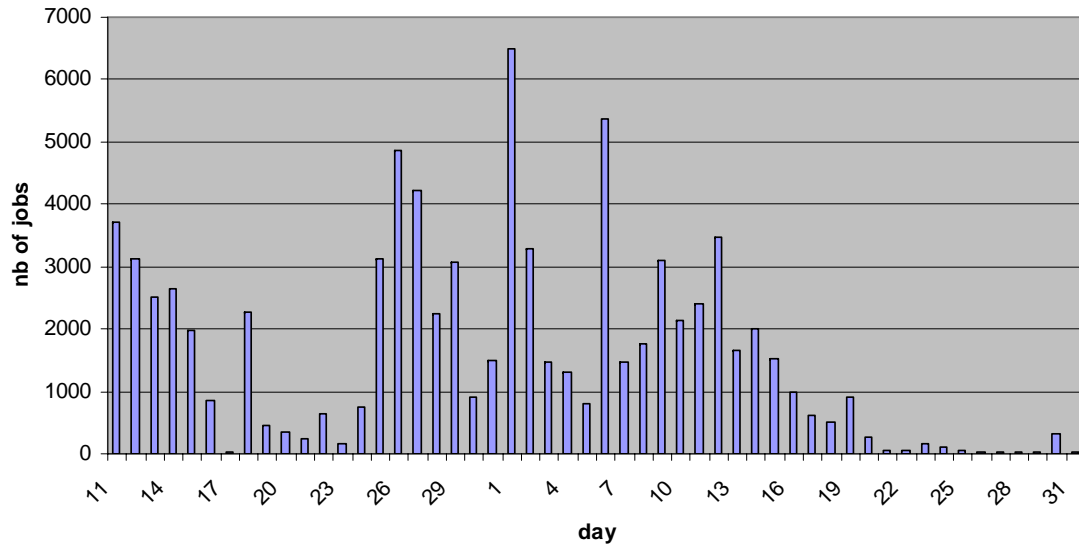
- **Statistics extracted by JRA2 for the user**
 - Number of registered jobs by RB and CE
 - Final status of the jobs
 - Duration of the jobs
 - Aborted job reasons

- **Efficiency from RBs : 81,5%**
- **Efficiency from user : 77%**
- **Global job efficiency for the user: ~46%**

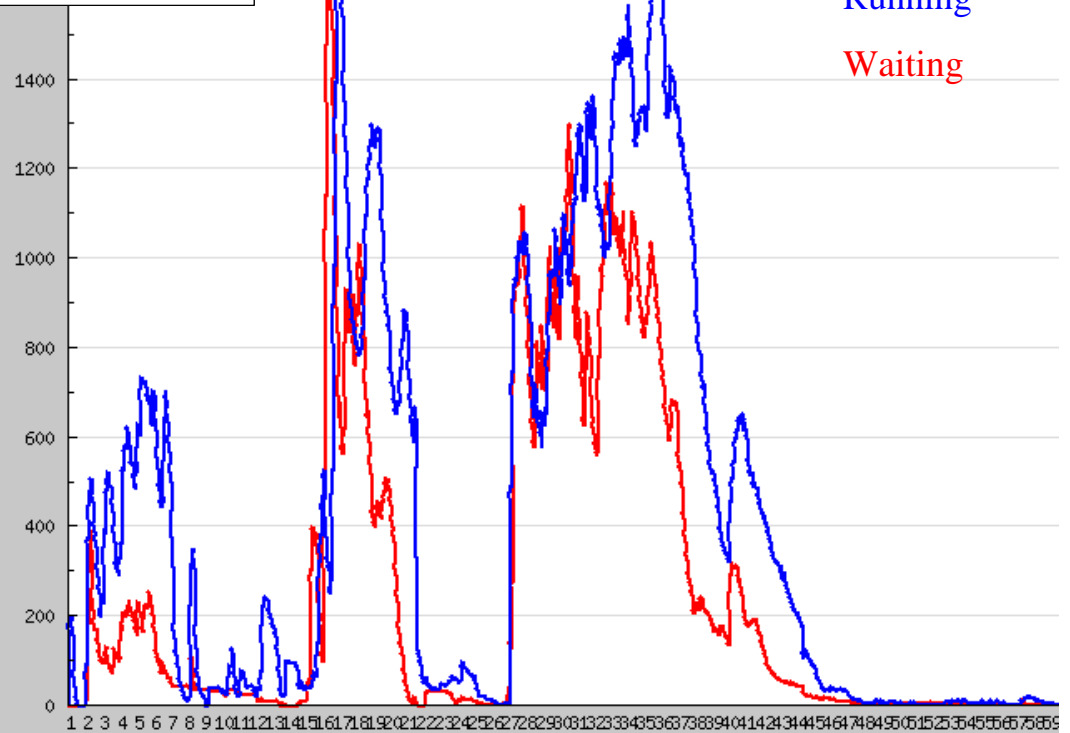
- **Failures reasons not estimated by JRA2**
 - License server failures
 - Some user failures : wisdom environment, human factor
 - Some sites failure : configuration (tar, disk space, configuration for biomed VO)

- **Proportions of the aborted reasons are only estimation**
 - We can't distinguish wisdom failure from grid failure

BIOMED-DC
registered jobs July-August



JRA2 statistics : registered jobs, day by day



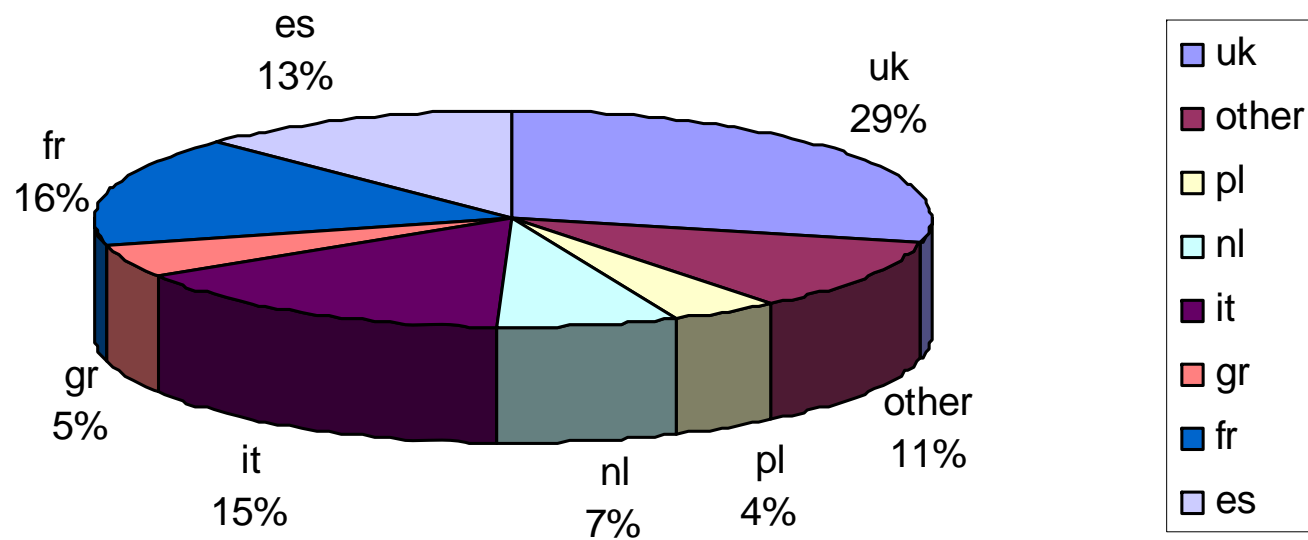
Registered jobs vs time

User statistics : registered jobs, min by min

Estimate number of jobs by RB

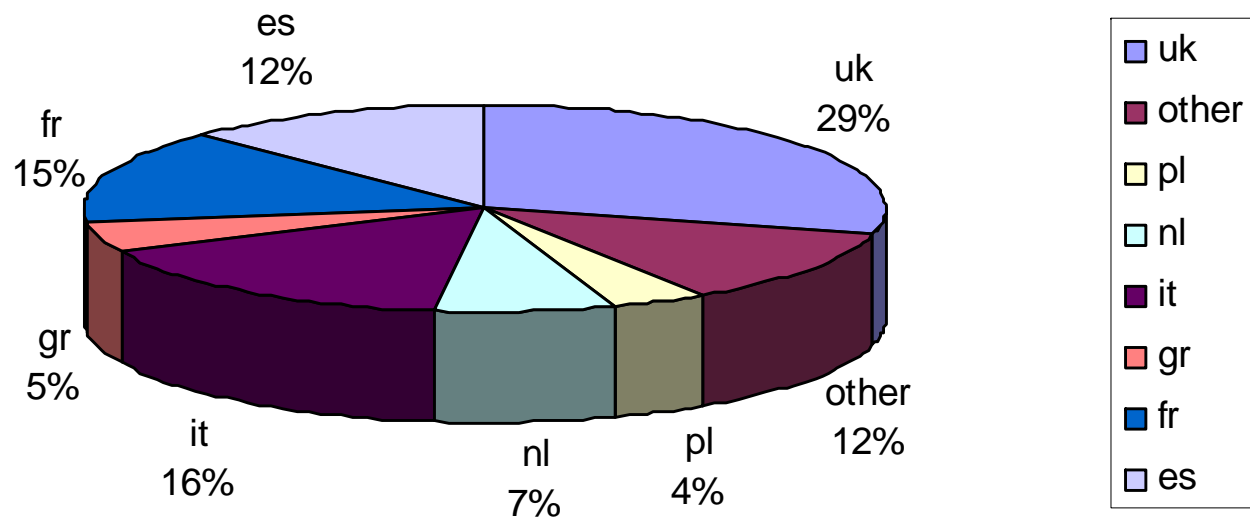
RB	Number of jobs	% user	% jra2	Number of jobs
rb.isabella.grnet.gr	10755	14,8	16,0	13701
bosheks.nikhef.nl	10684	14,7	14,9	12768
rb01.pic.es	9052	12,4	12,7	10877
lcg00124.grid.sinica.edu.tw	8221	11,3	11,5	9799
lcgrb01.gridpp.rl.ac.uk	7615	10,5	11,0	9436
rb1.egee.fr.cgg.com	7608	10,5	10,1	8684
node04.datagrid.cea.fr	4869	6,7	6,2	5302
grid09.lal.in2p3.fr	4586	6,3	5,6	4832
lappgrid07.in2p3.fr	4030	5,5	5,5	4723
egeerb.ifca.org.es	3519	4,8	4,2	3551
edt003.cnaf.infn.it	919	1,3	1,1	919
egee-rb-01.cnaf.infn.it	893	1,2	1,1	966
Total	72751	100	100	82093

Estimate number of jobs by CE

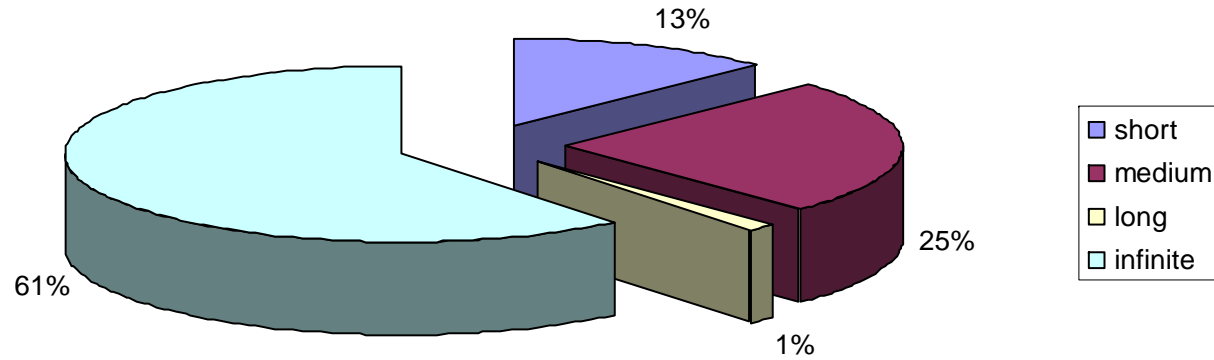


User statistics

JRA2 statistics

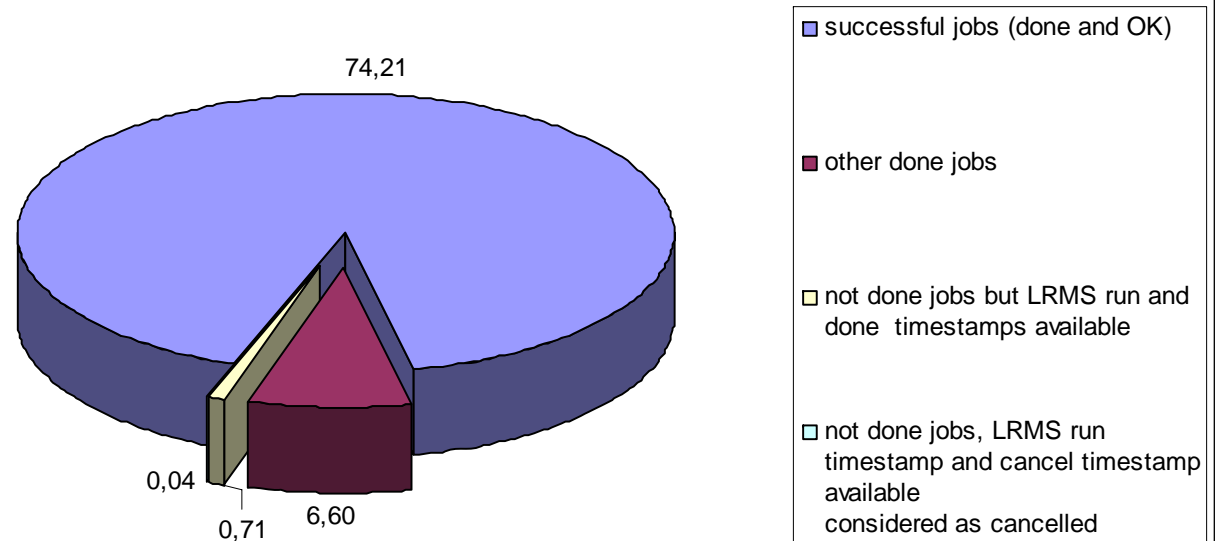


DC duration distribution for successful jobs



Duration statistics by JRA2

DC total duration (years)



User statistics : total running time is 67,5 CPU years

- **JRA2 statistics are coherent with the user statistics. They are probably more precise.**
- **But JRA2 statistics depend of the user to retrieve and validate the information**
- **For a next data challenge in the same conditions**
 - The user doesn't need to keep jobs information if LB data are saved and if JRA2 access RBs
 - The user needs to save the full list of jobIDs
 - The user needs to retrieve data transfer information, CPU time, failures on the UI and on the node
- **Is it possible to retrieve data transfer statistics during a data challenge ?**
 - Size, time, failures
 - Global, by SE, by CE
- **Is it possible to retrieve CE statistics ?**
 - CPU time, memory...
 - Configuration failures