

Experience in development of grid monitoring and accounting systems in Russia

NEC2009, Varna, Bulgaria
11.09.2009

Sergey Belov, Vladimir Korenkov
LIT JINR, Dubna
belov@jinr.ru

Grid monitoring – why it's important

- Allows to keep an eye on parameters of Grid operation in real time
- History view and analysis for major parameters of Grid elements
- Measurement resources' availability and reliability
- Fast and efficient error discovering and fixing
- Planning of huge calculating productions
- Discovering system's bottlenecks
- Failures prediction
- Forward planning of grid-architecture development (services, CPUs, disks, network channels)

What to monitor in Grid - examples

- **Sites**

- Site services – Computing Element, Storage Element, BDII, ...
- Services' certificates lifetime
- Functionality – job submission, file transfers

- **Core grid services**

- Workload management system (WMS)
- Information system, common for a group of sites (top BDII)
- File catalogs (LFC)
- Service for automated renew of user proxy certificates (MyProxy)

- **Resources**

- Computational nodes
- Network channels
- Disk space

- **User jobs**

Monitoring for Russian Data Intensive Grid

- **There are 15 Resource Centers in RDIG:**

- Ru-Moscow-SINP-LCG2, ITEP, JINR-LCG2, Kharkov-KIPT-LCG2, RRC-KI, RU-Moscow-KIAM-LCG2, RU-Phys-SPbSU, RU-SPbSU, Ru-Troitsk-INR-LCG2, ru-IMPB-LCG2, ru-Moscow-FIAN-LCG2, ru-Moscow-GCRAS-LCG2, ru-Moscow-MEPHI-LCG2, ru-PNPI, RU-Protvino-IHEP

- **Middleware**

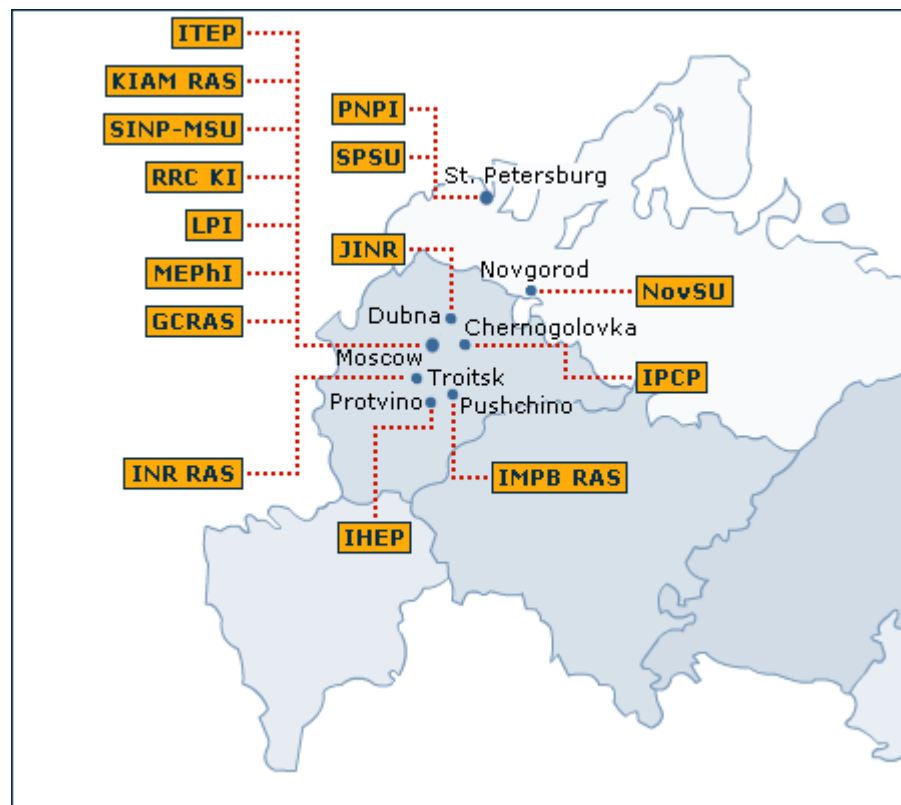
- gLite

- **Resources**

- > 3300 CPUs
- > 7700 kSI2k
- disk space about 1.2 PB

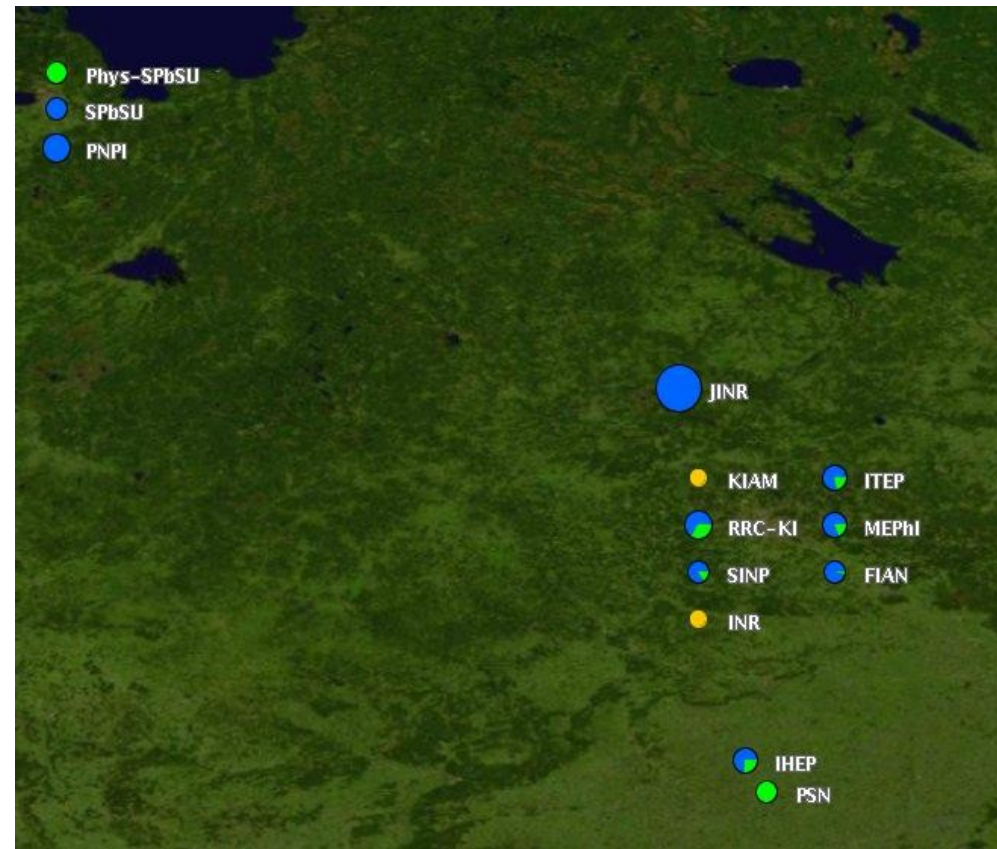
RDIG
Russian Data
Intensive Grid

eGee
Enabling Grids
for E-science

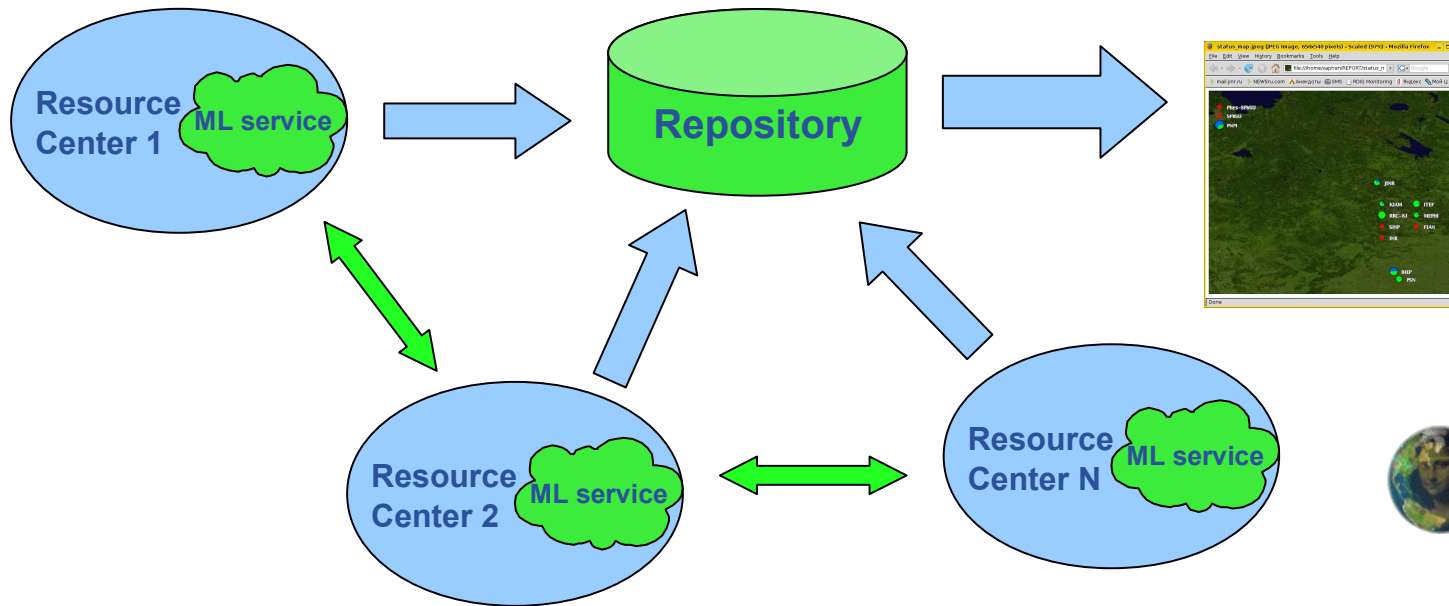


Main monitored values in RDIG

- **CPUs**
 - total
 - working / down
 - free / busy
 - total productivity
- **Jobs (per VO)**
 - running / waiting
- **Storage space**
 - used / available
 - totals for site
- **Network**
 - available bandwidth



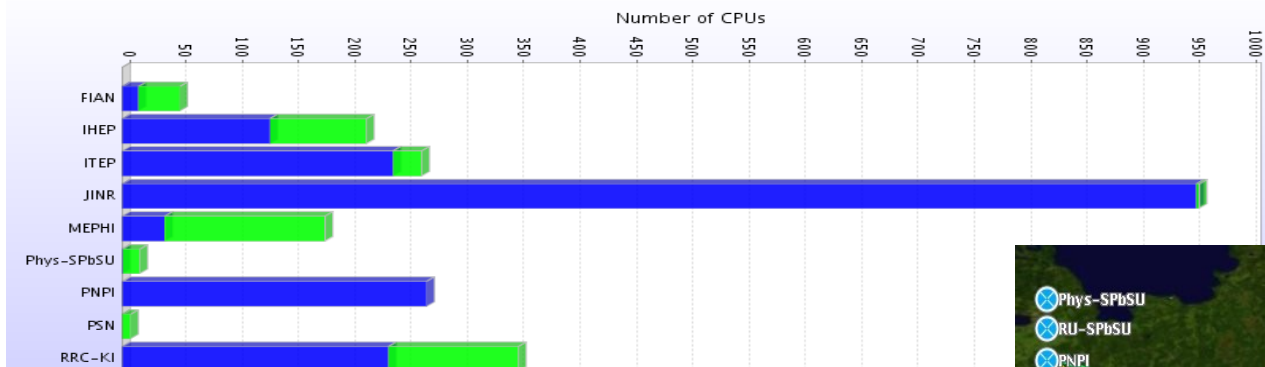
RDIG monitoring architecture



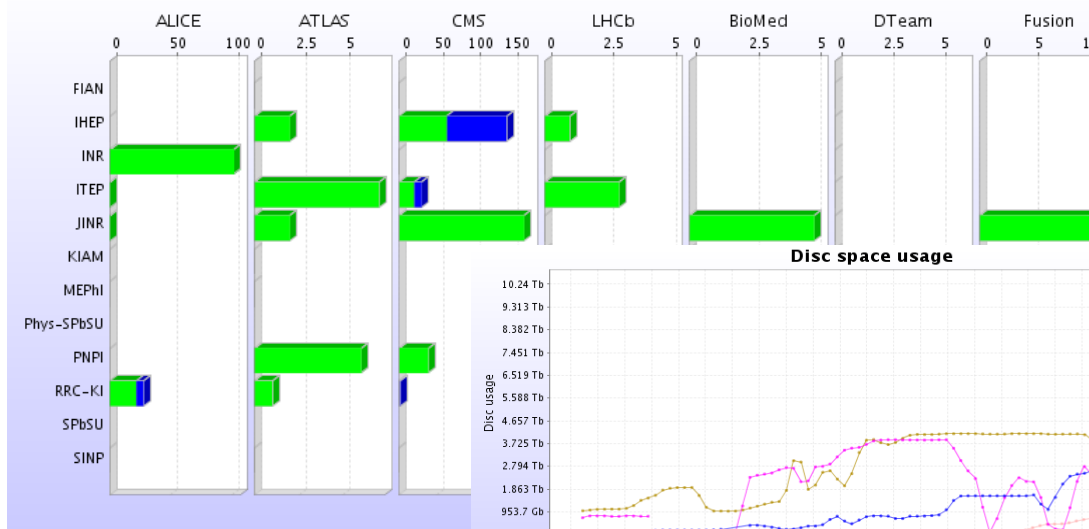
- **Monitoring package is installed at each Resource Center**
- **Collecting data locally from**
 - LRMS (torque)
 - Site BDII
 - Network bandwidth tests
- **Communication between services gives possibility to run network monitoring**

Some screenshots

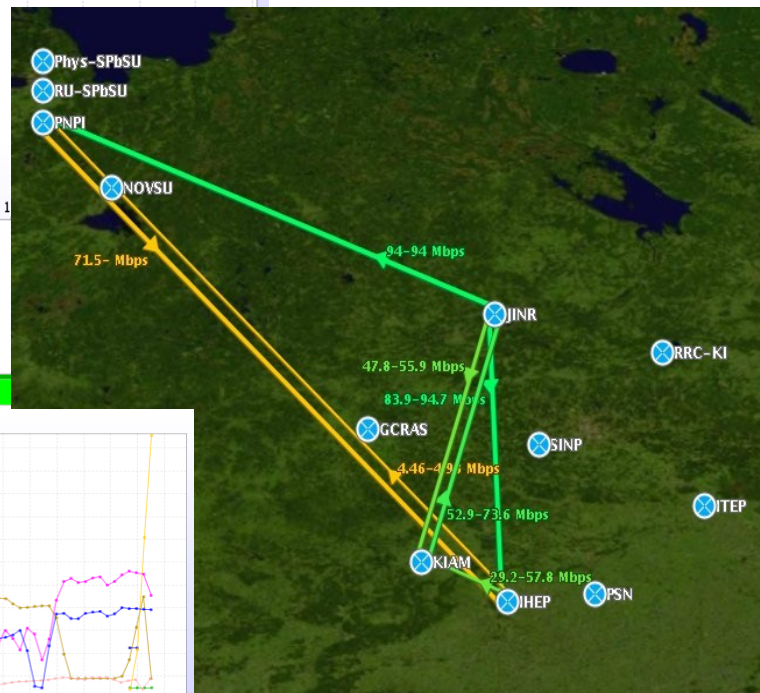
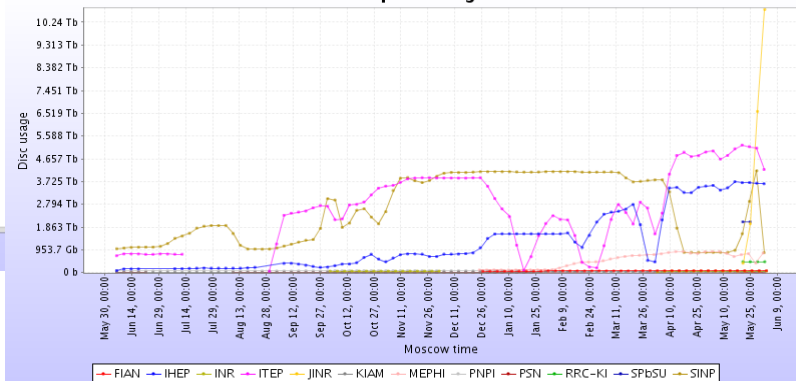
Farms usage



VO jobs running on RDIG farms



Disc space usage



<http://rocmon.jinr.ru:8080>

CMS PhEDEx file transfer statistics (I)

- **Is in use for RDIG sites**
- **Collects and stores the most useful information about transferred files**
 - source site
 - destination site
 - file size
 - transfer and validation time
 - transmission speed
 - start and stop time



How it works

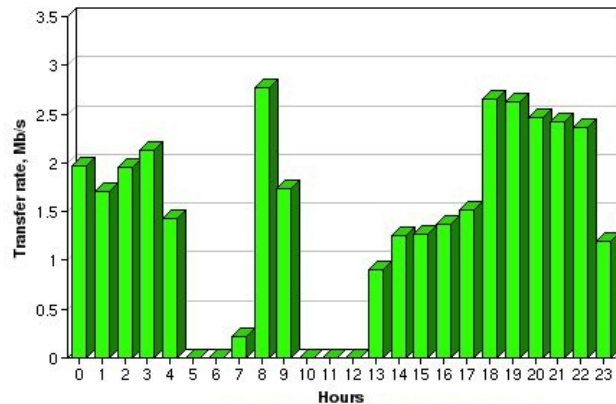
- **Parse PhEDEx log files**
- **Store information on each transmission to the database**
- **Aggregate transfers statistics to speed up information retrieve**
- **Dynamic representation on the Web**
 - hourly and daily statistics by sites
 - summary successful and failed transfers statistics (files count, volumes)
 - It's very easy to organize on-line file transfers monitoring

CMS PhEDEx file transfer statistics (III)

CMS PhEDEx file transfer statistics

Select interval: from to
 Get data for start date only
 Destination site:
 Source sites: ASGC CNAF FNAL FZK IN2P3 PIC RAL [\(check all\) \(uncheck all\)](#)

Transfers to SINP (08-11-2006)

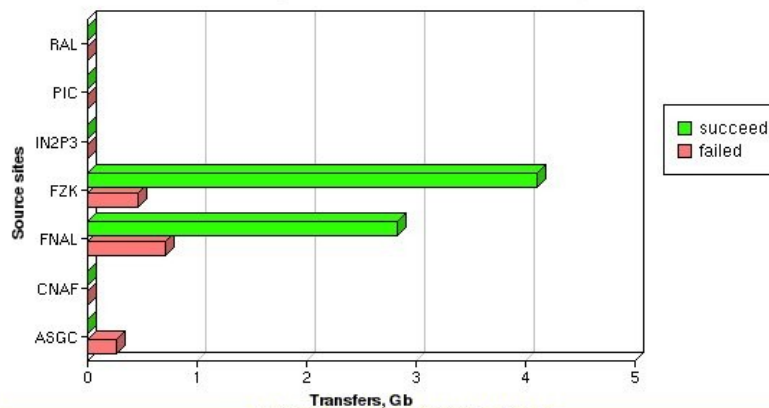


ChartDirector (unregistered) from www.adisofteng.com

File transfers statistics
destination: SINP
date: 08-11-2006

Source site	Succeed transfers	Failed transfers	Total count	Success ratio, %
ASGC	0	0	0	0
CNAF	0	0	0	0
FNAL	5	2	7	71
FZK	1	0	1	100
IN2P3	0	0	0	0
PIC	0	0	0	0
RAL	0	0	0	0
Total	6	2	8	75

Transfers to SINP (08-11-2006)



ChartDirector (unregistered) from www.adisofteng.com

File transfers statistics
destination: SINP
date: 08-11-2006

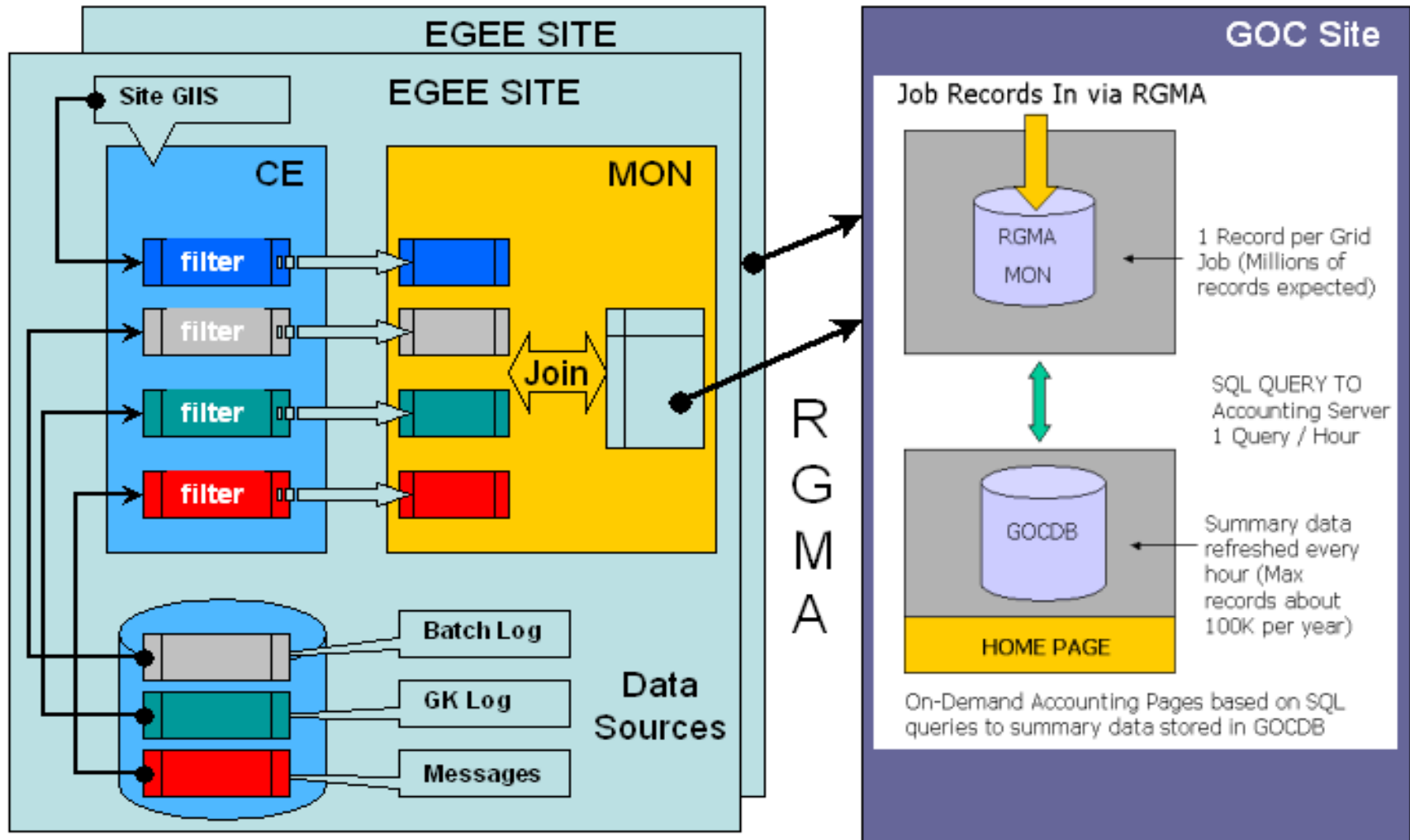
Source site	Successful transfers, Gb	Failed transfers, Gb
ASGC	0.000	0.262
CNAF	0.000	0.000
FNAL	2.827	0.716
FZK	4.102	0.458
IN2P3	0.000	0.000
PIC	0.000	0.000
RAL	0.000	0.000
Total	6.929	1.436

Why another monitoring in WLCG?

- **There are already:**
 - GStat
 - SAM tests
 - GridView
 - Dashboard
- **Why making one more tool?**
 - customizable parameters to gather
 - fit to regional managers demands
 - can combine all in one place
 - and many more other reasons

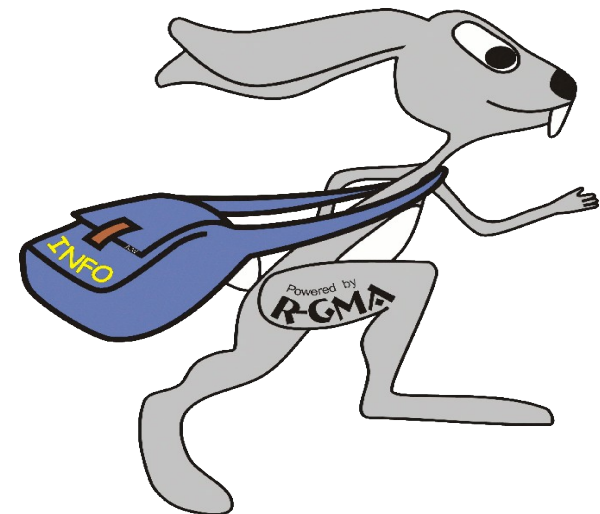
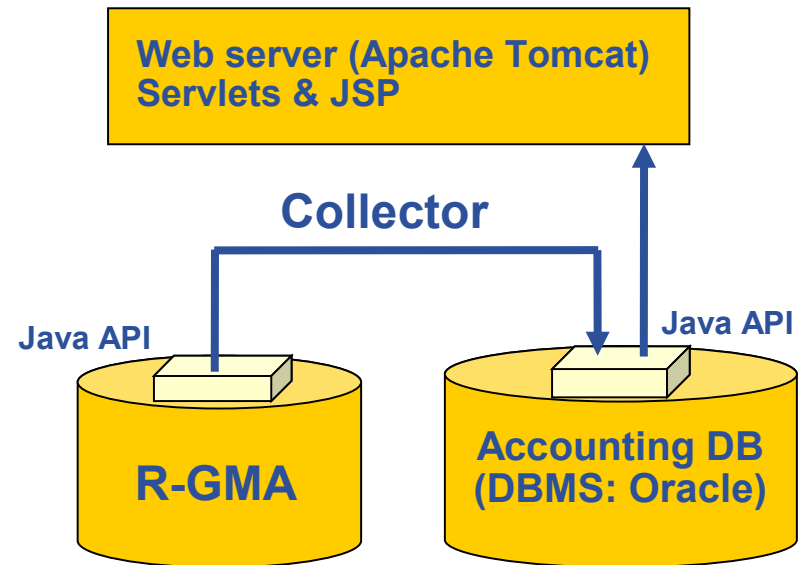
- Aim – to collect and store data on Grid resources usage (per user, per VO)
- Target users – VO, regional and site managers
- Could be used as a base for billing system in commercial projects
- CPU time is normalized to the system speed
 - now in hours*kSI2000
 - transition to HEPSPEC06 is upcoming

EGEE/RDIG accounting: technical



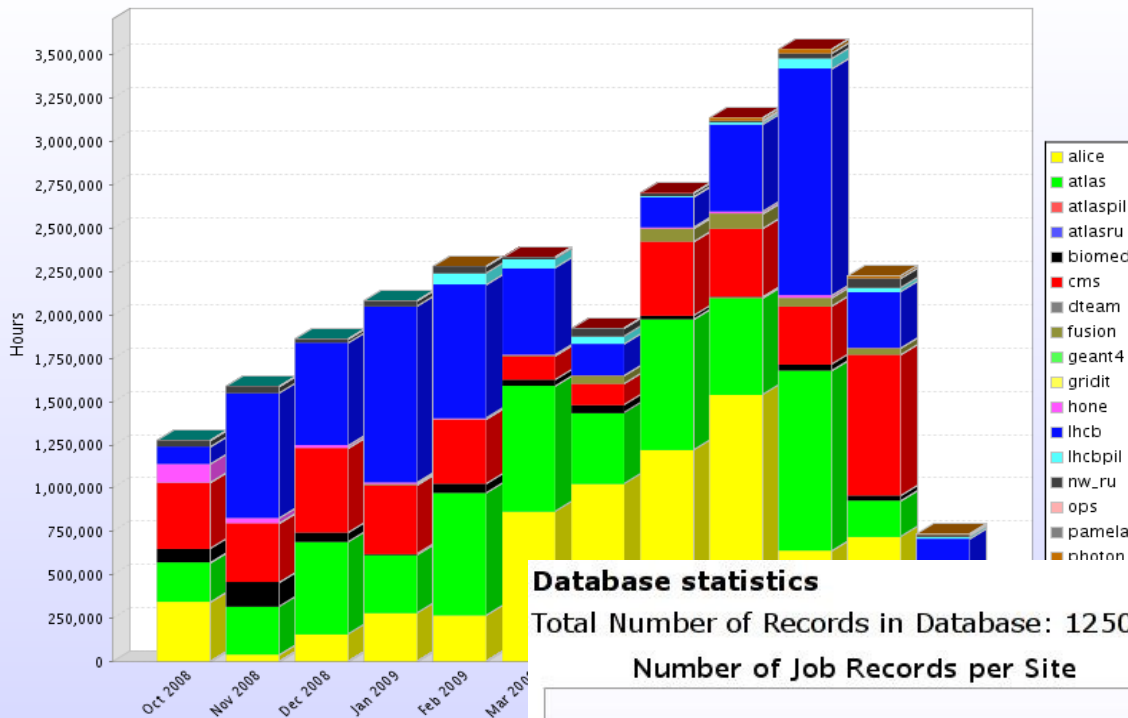
RDIG accounting architecture

- Collector runs daily to renew data in local DB with R-GMA information
- Interface was implemented on Java and works under Apache Tomcat servlet container
- User online requests are serving dynamically using the Accounting DB



RDIG accounting screenshots

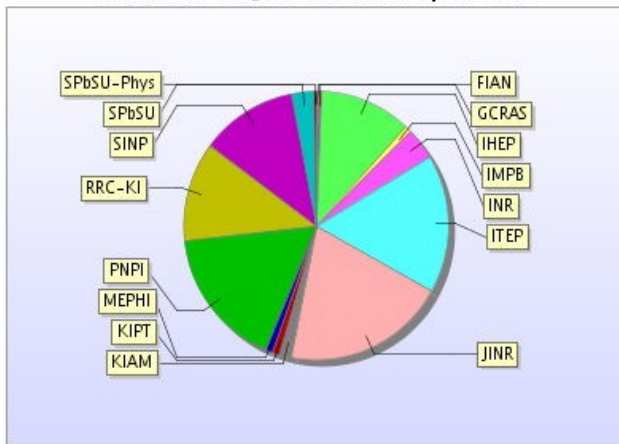
Normalised CPU time (SpectInt2000*hour = 1000)



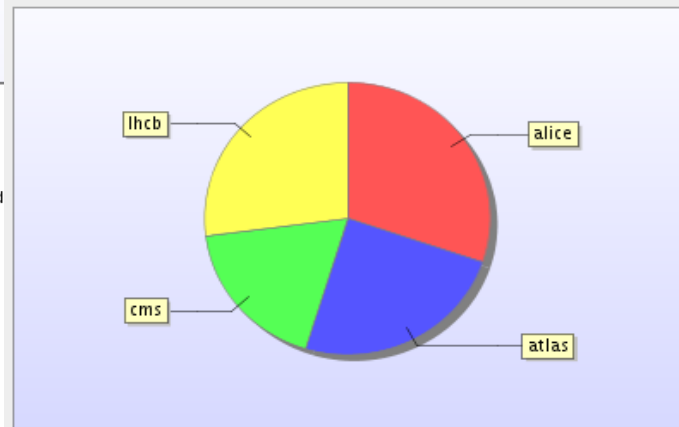
Database statistics

Total Number of Records in Database: 12502797

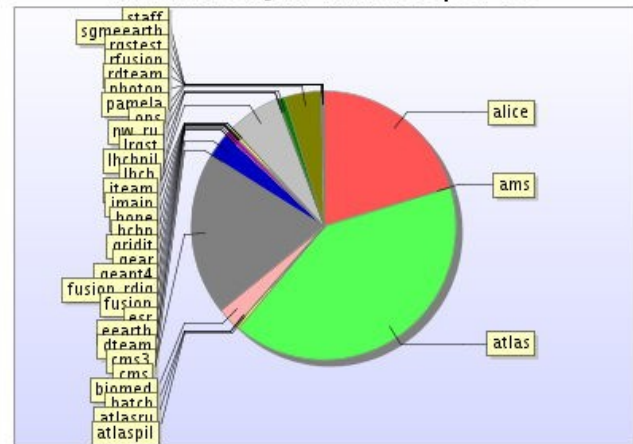
Number of Job Records per Site



Normalised CPU time (SpectInt2000*hour = 1000) per VO



Number of Job Records per VO



Some statistics from RDIG accounting

- **12.5 millions of job records since 2004**
- **44.5 millions kSI2k*hours of CPU time**
(21.6 millions plain hours = 2.5 thousands years)
- **Jobs from 33 VOs ever appeared in RDIG**
 - **alice**, ams , **atlas**, atlaspil, atlasru, batch, biomed, **cms**, cms3, dteam, eearth, esr, fusion, fusion_rdig, geant4, gear, gridit, hcbp, hone, imain, iteam, **lhcb**, lhcbpil, lrgst, nw_ru, ops, pamela, photon, rdteam, rfusion, rgstest, sgmeearth,staff
- **Normalized CPU time last year – 27.2 millions kSI2k*hours**
 - ALICE - 29%
 - ATLAS - 23%
 - CMS - 17%
 - LHCb - 25%
 - non-LHC VOs - 6%



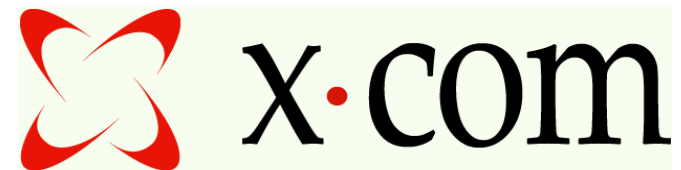
Why making our own accounting

- Representation of accounting data in any view necessary for Russian VO managers
- Having two sources is more reliable than only one
- “Backup” repository allows to perform crosschecks with the main EGEE/CESGA accounting portal
- Keeping detailed information on each job is useful to reveal errors
 - Example
 - double job counting was discovered after one of APEL updates:
 - *Format of a job key record was changed*
 - *If log data on a site are reprocessed with new APEL, some jobs are to be counted twice*
 - *about 200'000 excessive jobs were found for one year*

Collaboration with Romanian Tier2 Federation

- **Grid monitoring for Romanian sites**
 - CPUs
 - Disks
 - Site computational power (kSI2000)
 - VO jobs
- **Middleware on sites**
 - gLite
 - monitoring using some RDIG developments
 - Alien
 - special modules for collection data from ALICE
MonALISA were created

- **SKIF-GRID is a federation of Russian and Belorussian supercomputer centers**
- **Goal:**
 - HPC infrastructure for parallel computations
 - Break-through in development of domestic HPC technologies
- **Software:**
 - X-COM
 - UNICORE
 - AntMon, Ganglia
- **Monitoring**
 - Centralized monitoring site
 - Supercomputer resources
 - User jobs



GridNNN monitoring and accounting

- **Grid support for national nanotechnology network of Russia**



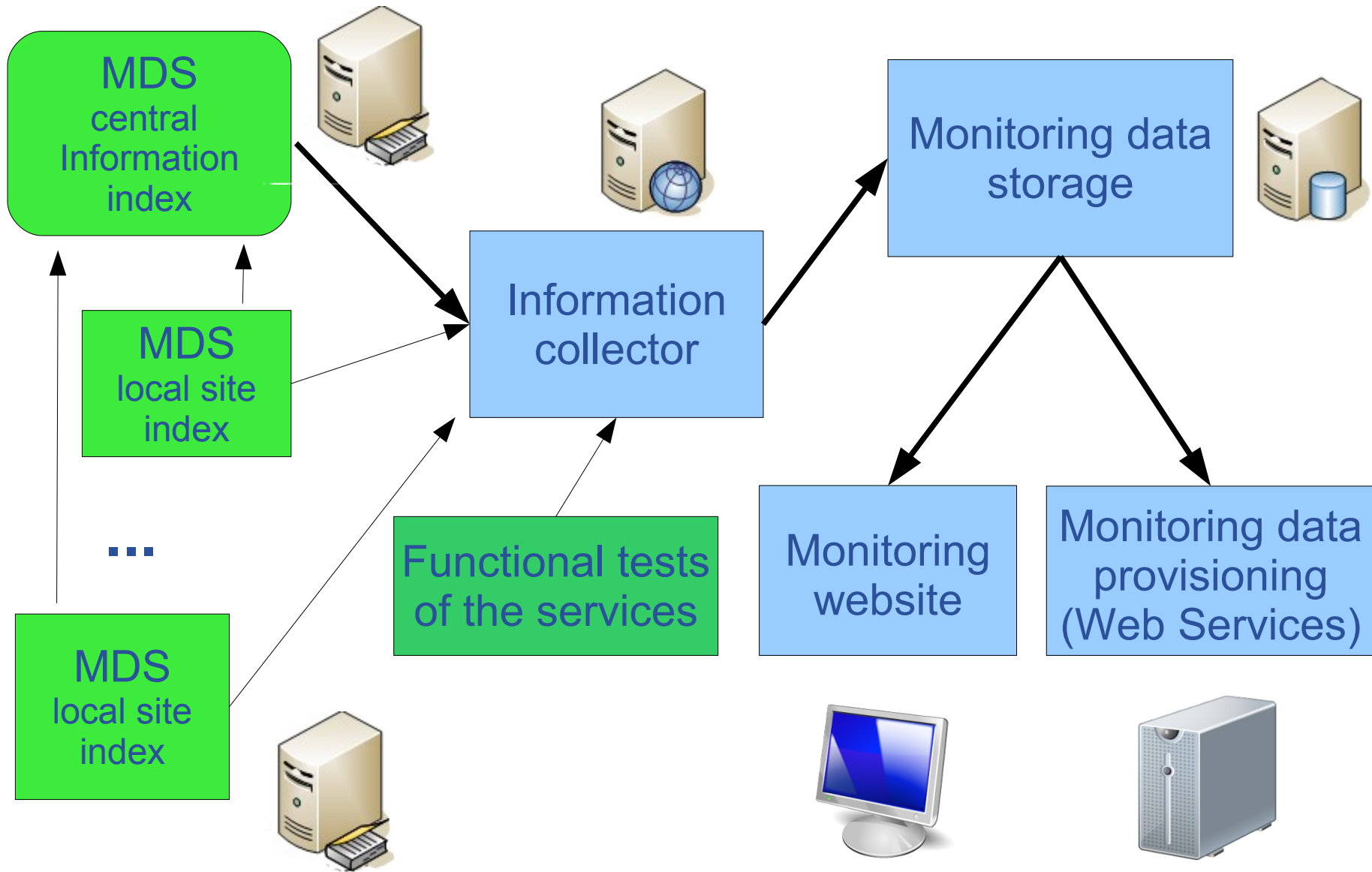
- To provide for science and industry an effective access to the distributed computational, informational and networking facilities
- Expecting breakthrough in nanotechnologies
- Supported by the special federal program

- **Main technical points**

- based on a network of supercomputers (about 15-30)
- has two grid operations centers (main and backup)
- is a set of grid services with unified interface
- partially based on Globus Toolkit 4



GridNNN monitoring - structure



- **There are a lot of cases when already developed systems are not enough**
- **We have positive experience in the area of creating multipurpose monitoring and accounting systems for Grid systems**
- **There are both developed and just planned project concerning Grid monitoring and accounting in Russia**
- **We plan to make our developments well-documented and public available for further external use**

Thank you for your attention!