

# *LHCb use of batch systems*

*A. Tsaregorodtsev,  
CPPM, Marseille*



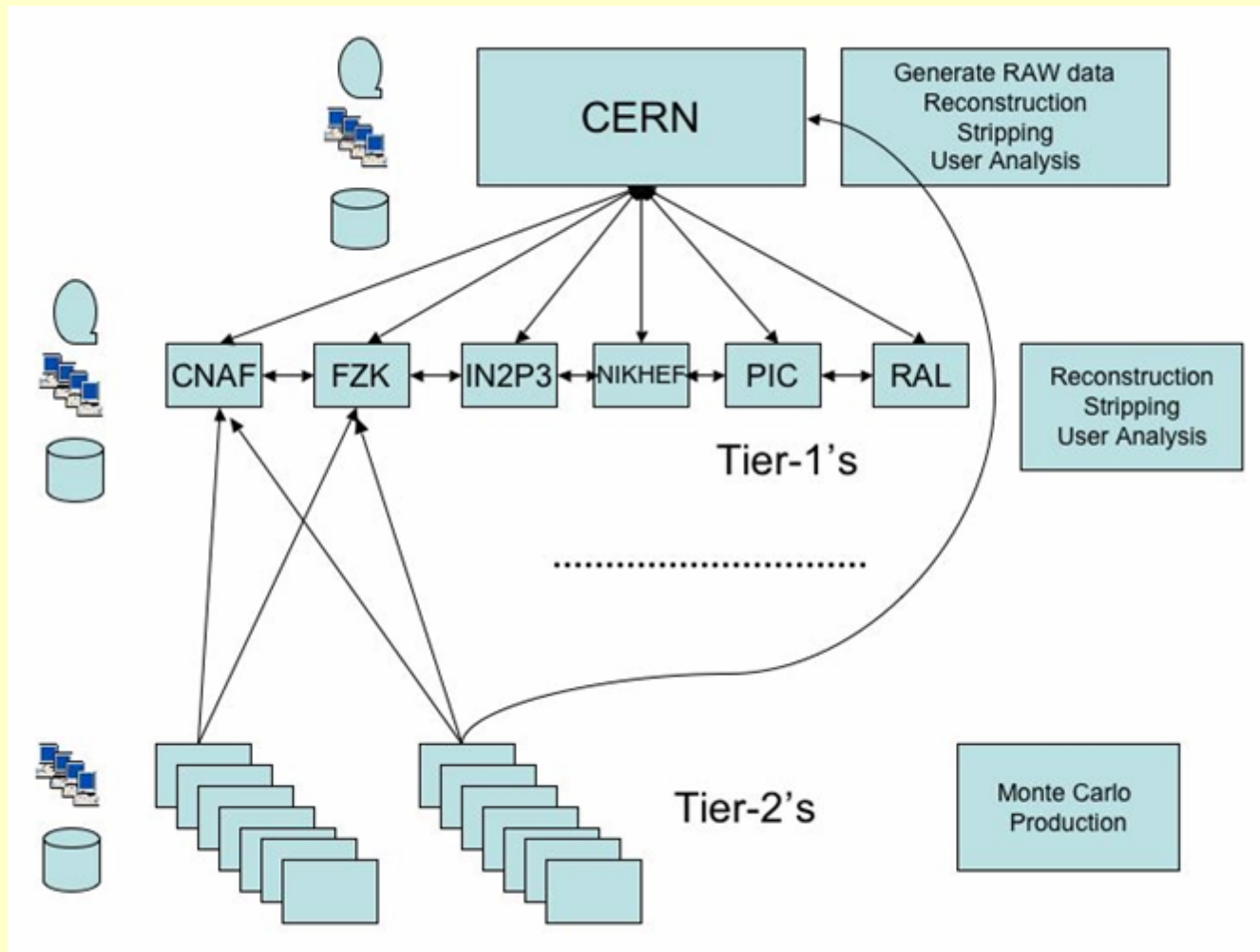
*HEPiX 2006 , 4 April 2006, Rome*

# Outline

---

- ◆ LHCb Computing Model
- ◆ DIRAC production and analysis system
- ◆ Pilot agent paradigm
- ◆ Application to the user analysis
- ◆ Conclusion

# LHCb Computing Model



# DIRAC overview

*DIRAC* – *D*istributed *I*nfrastructure with *R*emote *A*gent *C*ontrol

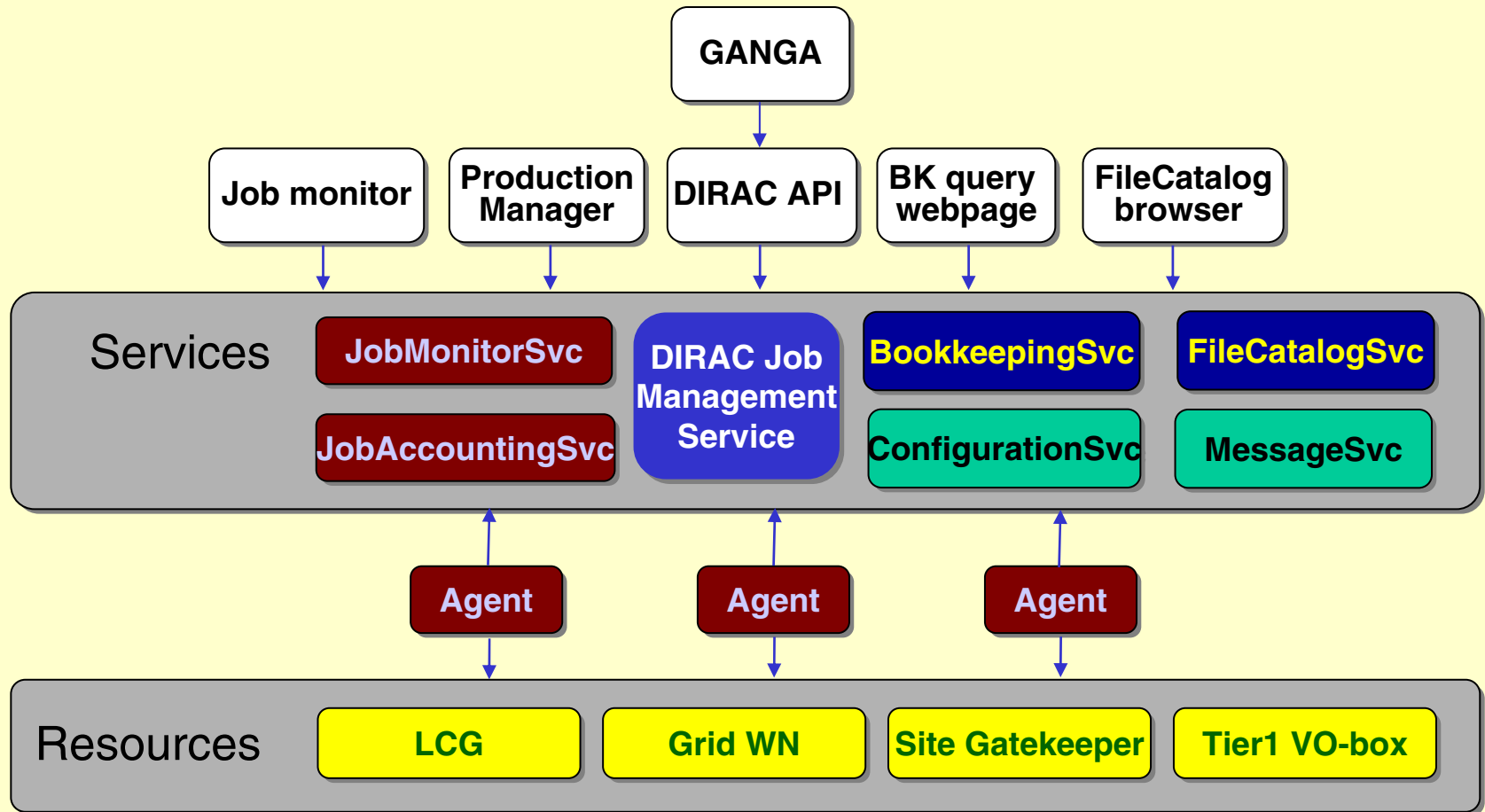
- ◆ LHCb grid system for the Monte-Carlo simulation data production and analysis
- ◆ Integrates computing resources available at LHCb production sites as well as on the LCG grid
- ◆ Composed of a set of light-weight services and a network of distributed agents to deliver workload to computing resources
- ◆ Runs autonomously once installed and configured on production sites
- ◆ Implemented in Python, using XML-RPC service access protocol



# DIRAC design goals

- ◆ Light implementation
  - ✦ Must be easy to deploy on various platforms
  - ✦ Non-intrusive
    - No root privileges, no dedicated machines on sites
  - ✦ Must be easy to configure, maintain and operate
- ◆ Using standard components and third party developments as much as possible
- ◆ High level of adaptability
  - ✦ There will be always resources outside LCGn domain
    - Sites that can not afford LCG, desktops, ...
  - ✦ We have to use them all in a consistent way
- ◆ Modular design at each level
  - ✦ Adding easily new functionality

# DIRAC Services, Agents and Resources



# DIRAC Services

- ◆ DIRAC Services are permanent processes deployed centrally or running at the VO-boxes and accepting incoming connections from clients (UI, jobs, agents)
- ◆ Reliable and redundant deployment
  - ✦ Running with watchdog process for automatic restart on failure or reboot
  - ✦ Critical services have mirrors for extra redundancy and load balancing
- ◆ Secure service framework:
  - ✦ XML-RPC protocol for client/service communication with GSI authentication and fine grained authorization based on user identity, groups and roles





# WMS Service

- ◆ DIRAC Workload Management System is itself composed of a set of central services, pilot agents and job wrappers
- ◆ The central Task Queue allows to apply easily the VO policies by prioritization of the user jobs
  - ✦ Using the accounting information and user identities, groups and roles (VOMS)
- ◆ The job scheduling happens in the last moment
  - ✦ With Pilot agents the job goes to a resource for immediate execution
- ◆ Sites are not required to manage user shares/priorities
  - ✦ Single long queue with guaranteed LHCb site quota will be enough

# DIRAC Agents

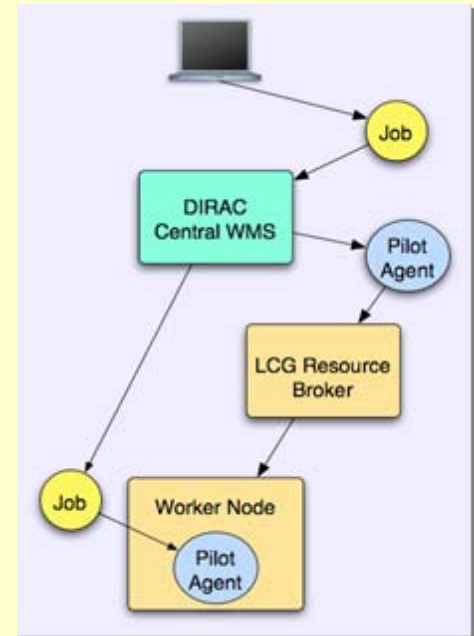
- ◆ Light easy to deploy software components running close to a computing resource to accomplish specific tasks
  - ✦ Written in Python, need only the interpreter for deployment
  - ✦ Modular easily configurable for specific needs
  - ✦ Running in user space
  - ✦ Using only outbound connections
- ◆ Agents based on the same software framework are used in different contexts
  - ✦ Agents for centralized operations at CERN
    - E.g. Transfer Agents used in the SC3 Data Transfer phase
    - Production system agents
  - ✦ Agents at the LHCb VO-boxes
  - ✦ Pilot Agents deployed as LCG jobs

# Pilot agents

- ◆ Pilot agents are deployed on the Worker Nodes as regular jobs using the standard LCG scheduling mechanism
  - ✦ Form a distributed Workload Management system
- ◆ Once started on the WN, the pilot agent performs some checks of the environment
  - ✦ Measures the CPU benchmark, disk and memory space
  - ✦ Installs the application software
- ◆ If the WN is OK the user job is retrieved from the central DIRAC Task Queue and executed
- ◆ In the end of execution some operations can be requested to be done asynchronously on the VO-box to accomplish the job

# Distributed Analysis

- ◆ The Pilot Agent paradigm was extended recently to the Distributed Analysis activity
- ◆ The advantages of this approach for users are:
  - ✦ Inefficiencies of the LCG grid are completely hidden from the users
  - ✦ Fine optimizations of the job turnaround
    - It also reduces the load on the LCG WMS
- ◆ The system was demonstrated to serve dozens of simultaneous users with about 2Hz submission rate
  - ✦ The limitation is mainly in the capacity of LCG RB to schedule this number of jobs

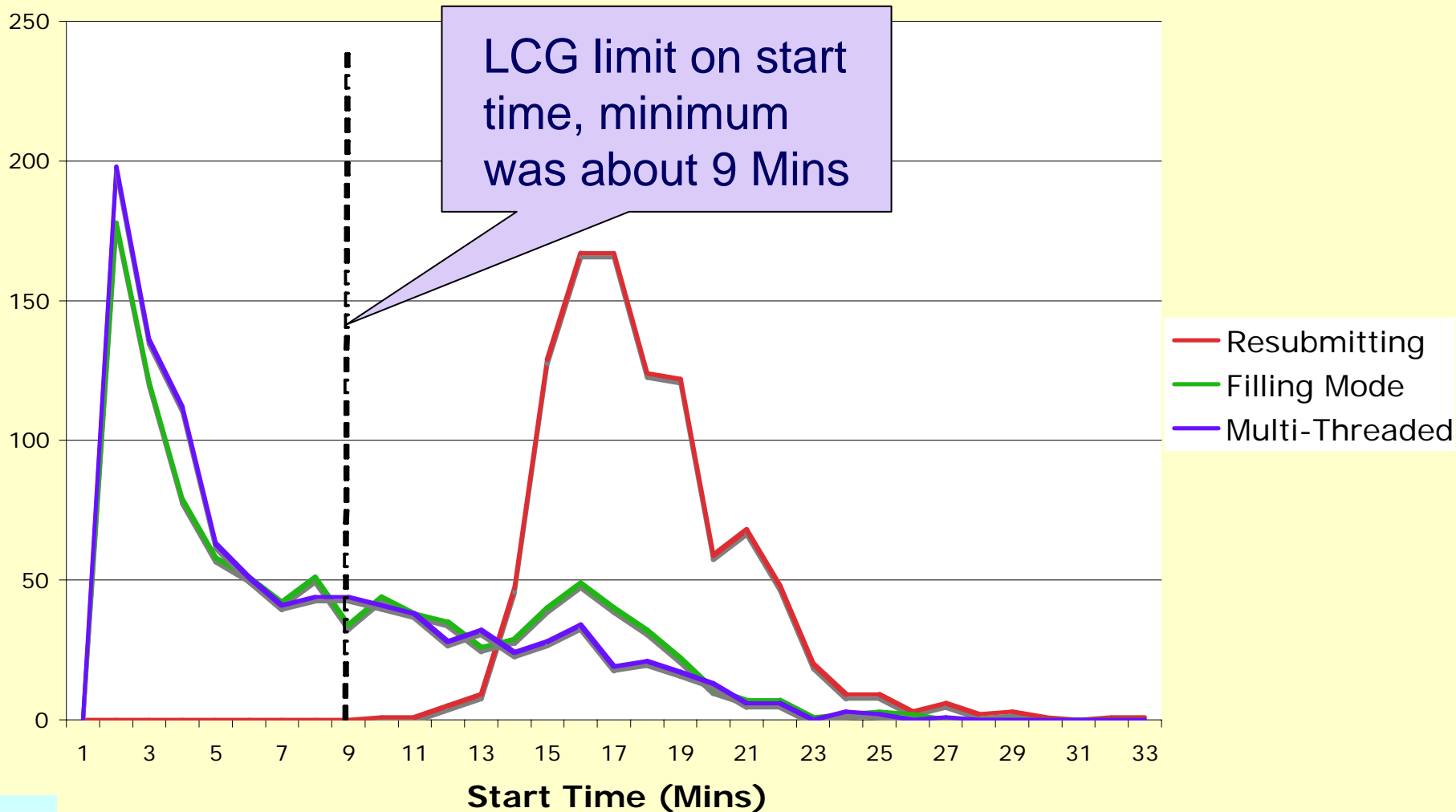


# DIRAC WMS Pilot Agent Strategies

- ◆ The combination of pilot agents running right on the WNs with the central Task Queue allows fine optimization of the workload on the VO level
  - ✦ The WN reserved by the pilot agent is a first class resource - there is no more uncertainty due to delays in the local batch queue
  
- ◆ DIRAC Modes of submission
  - ✦ 'Resubmission'
    - Pilot Agent submission to LCG with monitoring
    - Multiple Pilot Agents may be sent in case of LCG failures
  - ✦ 'Filling Mode'
    - Pilot Agents may request several jobs from the same user, one after the other
  - ✦ 'Multi-Threaded'
    - Same as 'Filling' Mode above except two jobs can be run in parallel on the Worker Node

# Start Times for 10 Experiments, 30 Users

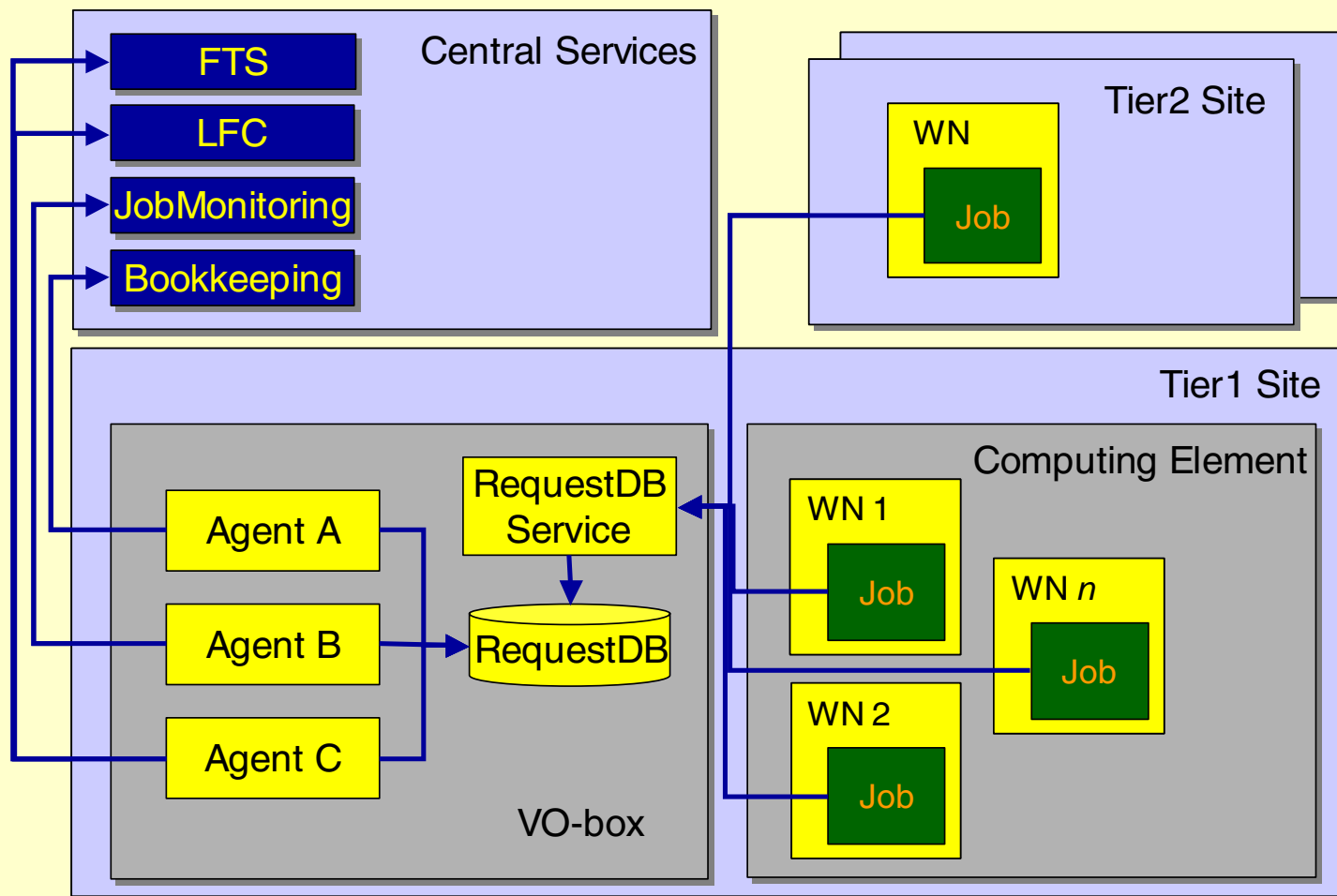
Start Times for 30 Users  
3000 Jobs, 1.5 Million Events



# VO-box

- ◆ VO-boxes are dedicated hosts at the Tier1 centers running specific LHCb services for
  - ✦ Reliability due to retrying failed operations
  - ✦ Efficiency due to early release of WNs and delegating data moving operations from jobs to the VO-box agents
- ◆ Agents on VO-boxes execute requests for various operations from local jobs:
  - ✦ Data Transfer requests
  - ✦ Bookkeeping, Status message requests

# LHCb VO-box architecture



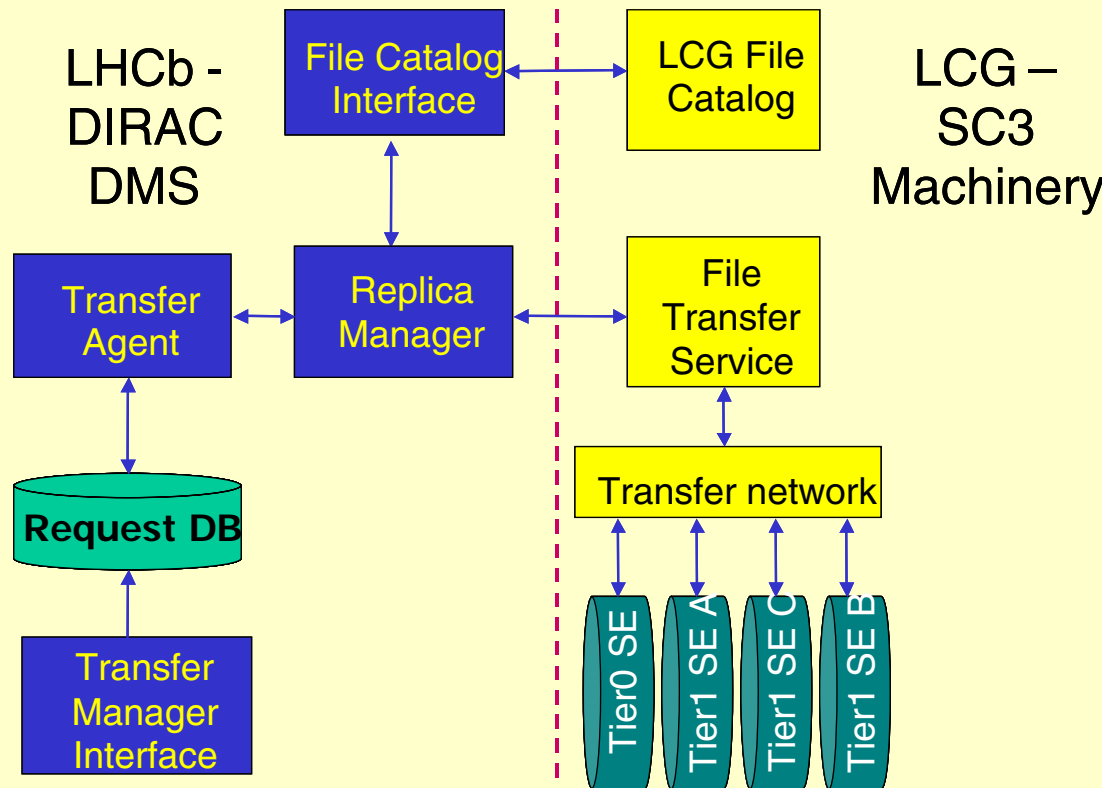


# Transfer Agent example

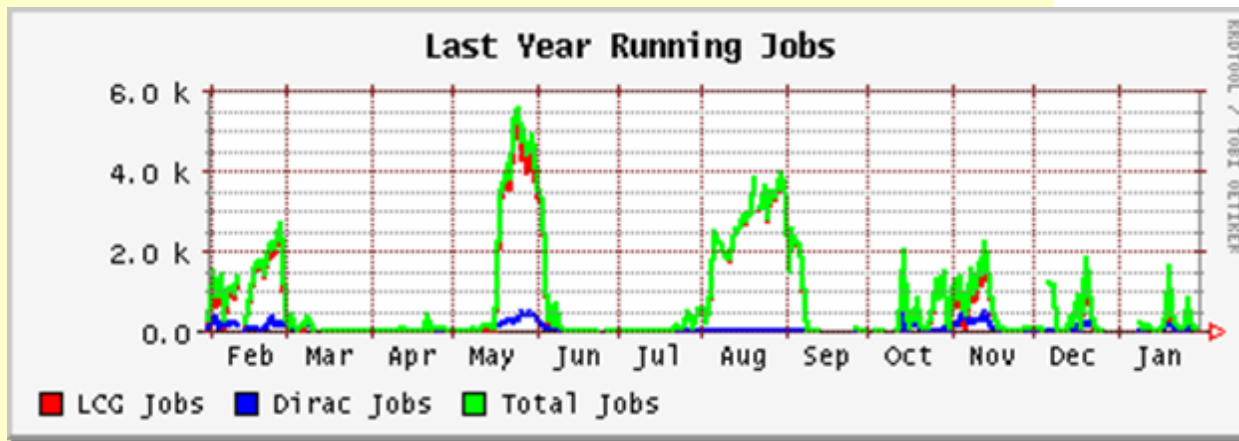
- ◆ Request DB is populated with data transfer/replication requests from Data Manager or jobs

- ◆ Transfer Agent

- ◆ checks the validity of request and passes to the FTS service
- ◆ uses third party transfer in case of FTS channel unavailability
- ◆ retries transfers in case of failures
- ◆ registers the new replicas in the catalog

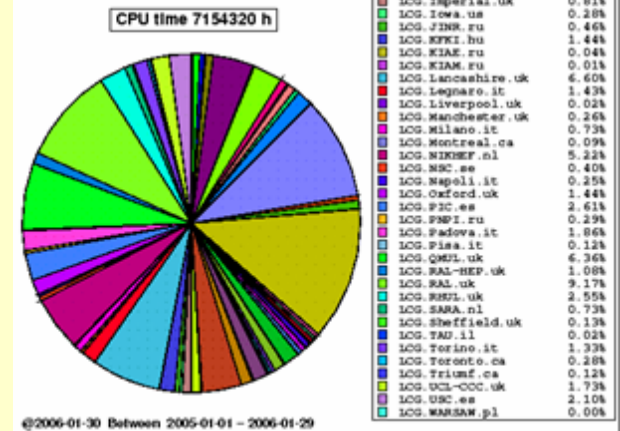


# DIRAC production performance



- ◆ Up to over 5000 simultaneous production jobs
  - ✦ The throughput is only limited by the capacity available on LCG
- ◆ ~80 distinct sites accessed through LCG or through DIRAC directly

DIRAC.Barcelona.es	0.19%
DIRAC.Bologna-12.it	0.63%
DIRAC.CERN.ch	0.41%
DIRAC.Cambridge.uk	0.00%
DIRAC.CracowAgd.pl	0.00%
DIRAC.IF-UFPR.br	0.14%
DIRAC.LINCOLNLINE.ch	0.71%
DIRAC.Lyon.fr	3.87%
DIRAC.FMPI.ru	0.00%
DIRAC.Santiago.es	0.13%
DIRAC.ScotGrid.uk	2.81%
DIRAC.Suichb-eps.ch	0.00%
DIRAC.Eurich.ch	0.71%
ICG.ACAD.bg	0.15%
ICG.BHAM-HEP.uk	0.76%
ICG.Barcelona.es	0.44%
ICG.Saci.it	1.55%
ICG.Bologna.it	0.04%
ICG.CERN.ch	9.80%
ICG.CESGA.es	0.48%
ICG.COG.Fr	0.81%
ICG.CMNF-GRIDIT.it	0.03%
ICG.CMNF.it	12.74%
ICG.CNB.es	0.39%
ICG.CPPM.fr	0.30%
ICG.CSCS.ch	0.39%
ICG.CY1.cy	0.16%
ICG.Cagliari.it	0.47%
ICG.Cambridge.uk	0.03%
ICG.Catania.it	0.50%
ICG.Durban.uk	0.43%
ICG.Edinburgh.uk	0.03%
ICG.FZK.de	1.54%
ICG.Ferrara.it	0.07%
ICG.Firenze.it	1.03%
ICG.GR-01.gr	0.34%
ICG.GR-02.gr	0.26%
ICG.GR-03.gr	0.17%
ICG.GR-04.gr	0.06%
ICG.GRNET.gr	1.38%
ICG.HPC2N.se	0.00%
ICG.ICT.ro	0.13%
ICG.IFCA.es	0.02%
ICG.IHEP.eu	1.10%
ICG.IJGPP.fr	3.77%
ICG.IMTA.es	0.09%
ICG.IPP.bg	0.03%
ICG.IPSI-IPSP.fr	0.01%
ICG.ITEP.ru	0.92%
ICG.Imperial.uk	0.81%
ICG.Iowa.us	0.28%
ICG.JINR.ru	0.46%
ICG.KFI.hu	1.44%
ICG.KIAC.ru	0.04%
ICG.KIAM.ru	0.01%
ICG.Lancashire.uk	6.60%
ICG.Legnaro.it	1.43%
ICG.Liverpool.uk	0.02%
ICG.Manchester.uk	0.26%
ICG.Milano.it	0.73%
ICG.Montreal.ca	0.09%
ICG.NIKHEF.nl	5.22%
ICG.NIC.se	0.40%
ICG.Napoli.it	0.25%
ICG.Oxford.uk	1.44%
ICG.PIC.es	2.61%
ICG.FMPI.ru	0.29%
ICG.Padova.it	1.86%
ICG.Pisa.it	0.12%
ICG.QMUL.uk	6.36%
ICG.RAL-HEP.uk	1.08%
ICG.RAL.uk	9.17%
ICG.RHUL.uk	2.55%
ICG.SARA.nl	0.73%
ICG.Sheffield.uk	0.13%
ICG.TAD.il	0.02%
ICG.Torino.it	1.33%
ICG.Toronto.ca	0.28%
ICG.Triest.ca	0.12%
ICG.UCL-CCC.uk	1.73%
ICG.USC.es	2.10%
ICG.WARSAW.pl	0.00%

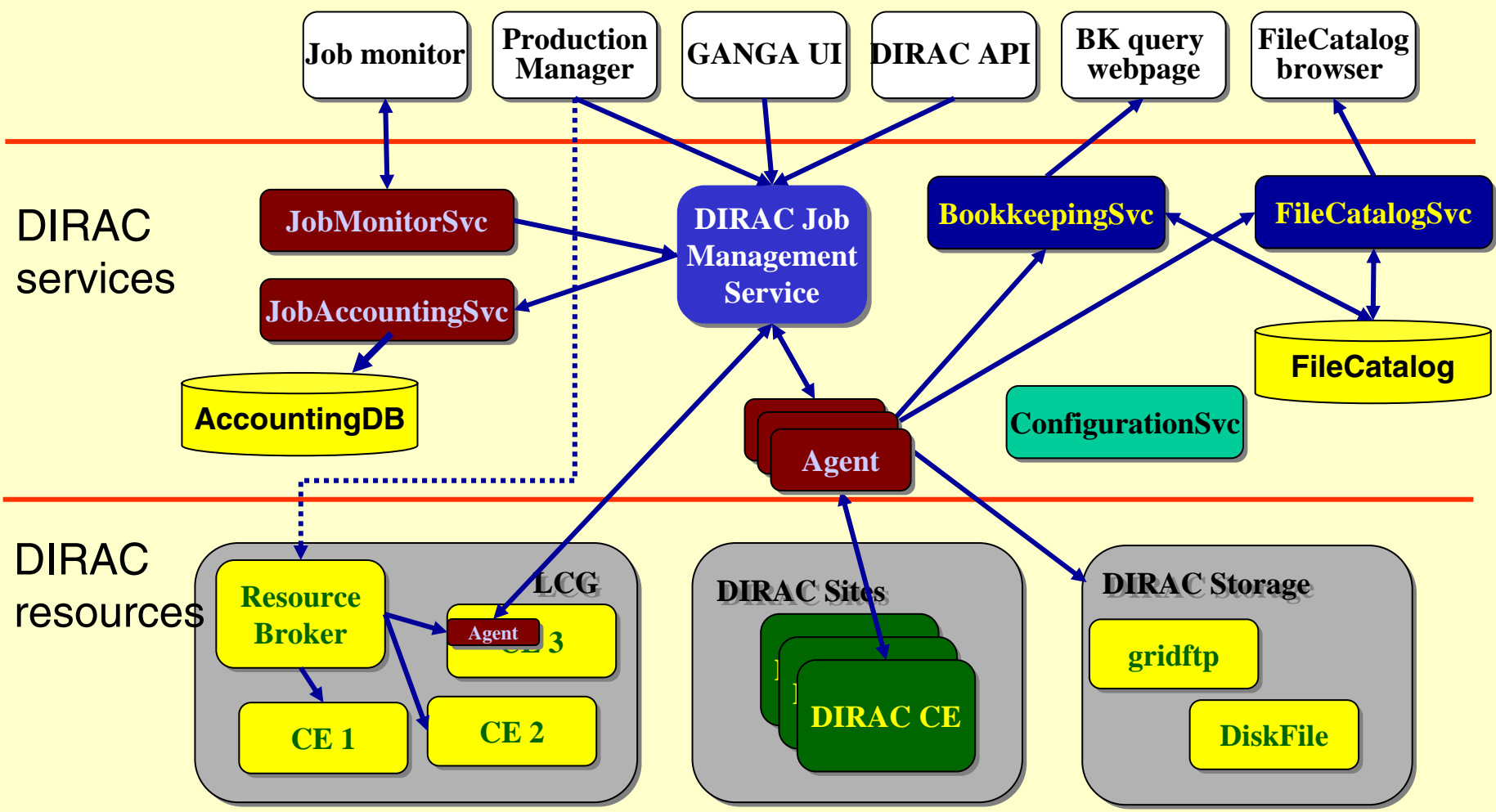


# Conclusions

- ◆ The Overlay Network paradigm employed by the DIRAC system proved to be efficient in integrating heterogeneous resources in a single reliable system for simulation data production
- ◆ The system is now extended to deal with the Distributed Analysis tasks
  - ✦ Workload management on the user level is effective
  - ✦ Real users (~30) are starting to use the system
- ◆ The LHCb Data Challenge 2006 in June
  - ✦ Test LHCb Computing Model before data taking
  - ✦ An ultimate test of the DIRAC system

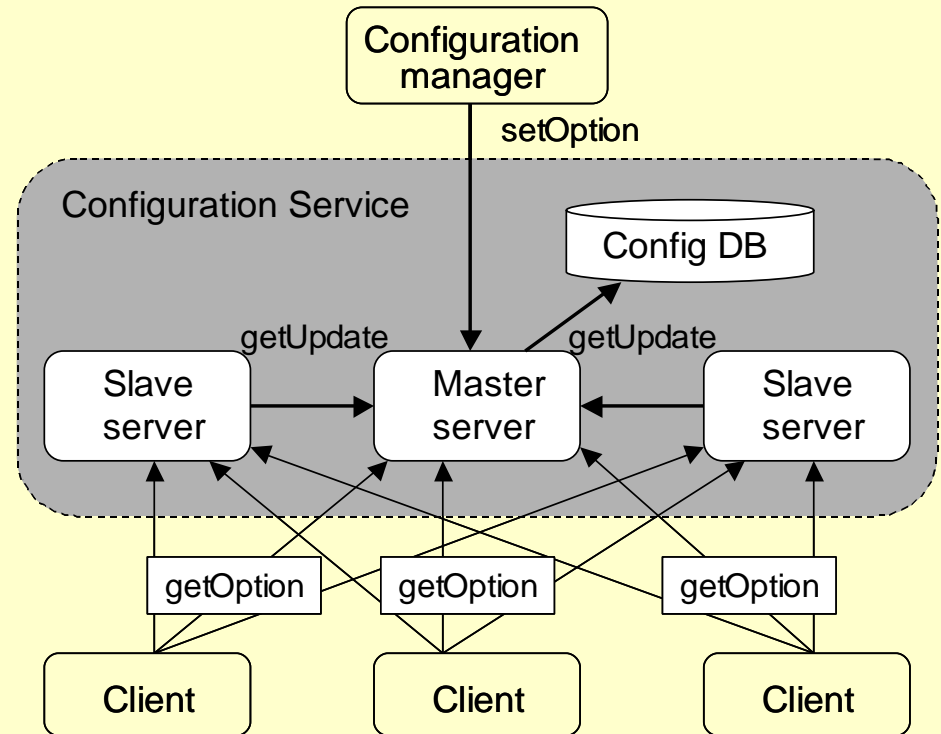
Back-up slides

# DIRAC Services and Resources



# Configuration service

- ◆ Master server at CERN is the only one allowing write access
- ◆ Redundant system with multiple read-only slave servers running at sites on VO-boxes for load balancing and high availability
- ◆ Automatic slave updates from the master information
- ◆ Watchdog to restart the server in case of failures



# Other Services

- ◆ Job monitoring service
  - ✦ Getting job heartbeats and status reports
  - ✦ Service the job status to clients ( users )
    - Web and scripting interfaces
- ◆ Bookkeeping service
  - ✦ Receiving, storing and serving job provenance information
- ◆ Accounting service
  - ✦ Receives accounting information for each job
  - ✦ Generates reports per time period, specific productions or user groups
  - ✦ Provides information for taking policy decisions

- ◆ DIRAC is a distributed data production and analysis system for the LHCb experiment
  - ✦ Includes workload and data management components
  - ✦ Was developed originally for the MC data production tasks
  - ✦ The goal was:
    - integrate all the heterogeneous computing resources available to LHCb
    - Minimize human intervention at LHCb sites
  - ✦ The resulting design led to an architecture based on a set of services and a network of light distributed agents



# File Catalog Service

- ◆ LFC is the main File Catalog
  - ✦ Chosen after trying out several options
  - ✦ Good performance after optimization done
  - ✦ One global catalog with several read-only mirrors for redundancy and load balancing
- ◆ Similar client API as for other DIRAC “File Catalog” services
  - ✦ Seamless file registration in several catalogs
  - ✦ E.g. Processing DB receiving data to be processed automatically

# DIRAC performance

## ◆ Performance in the 2005 RTTC production

- ◆ Over 5000 simultaneous jobs
  - Limited by the available resources
- ◆ Far from the critical load on the DIRAC servers

