

# DIRAC-Belle & Amazon EC<sup>2</sup>



R. G., A. Casajus & T. Fifield



- Motivation
- Strategy
- New DIRAC Developments
- Results
- Summary

# Motivation

- Belle (and Belle II in the future) has large peaks of CPU need for Monte Carlo simulation.
- It has been shown that Amazon CE<sup>2</sup> can be used to provide extra CPU resources (\*).
- Need to make a large scale exercise to determine efficiency and cost.
- Use a well established tool easily integrate with existing resources and to reduce development time.

(\*) "Belle Monte-Carlo Production on the Amazon EC2 Cloud" CHEP 2009.

# Strategy



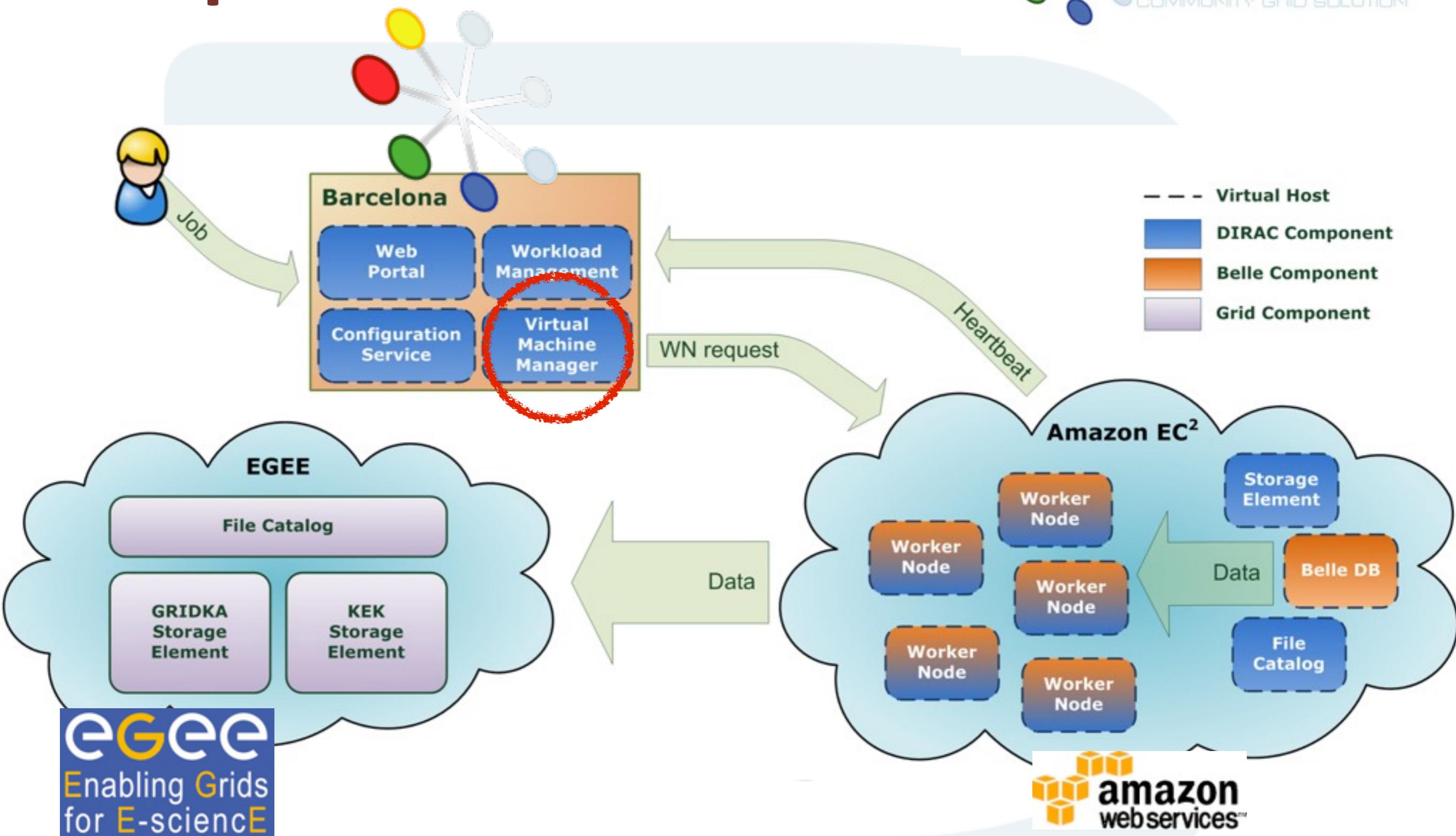
- Use DIRAC framework to control execution and integrate cloud resources with grid and local clusters:
  - Grid: DIRAC Pilot Jobs
  - Cloud: DIRAC Virtual Machines
- Divide exercise into different phases:
  - Phase I:
    - Test cloud
  - Phase II:
    - Incorporate local resources
  - Phase III:
    - Incorporate grid resources
- Make extensive use of virtualization for future portability and scalability.

# DIRAC Design choices



- **VO Centric:**
  - Gives to the community, the VO, a central role in the relation of its users with their computing resources.
- **Modularity:**
  - To achieve optimal scalability and flexibility, a highly modular design was decided.
- **Pull Scheduling:**
  - Implements pull scheduling with late binding of payload to resource to extract optimal performance out of the ever changing underlying resources.

# Proposed DIRAC solution



# Central Components (@Barcelona)



- Configuration Server:
  - Describes resources and the way to access them: Sites, CE's, SE's, FC's,...
  - Users and Policies.
  - Location of DIRAC components.
  - Configurable parameters.
- Web Portal:
  - Monitoring
  - Accounting
  - Interactivity
- Workload Management:
  - Job repository and agents to access the resources
- Virtual Machine:
  - Scheduler for submission
  - Manager for control and monitoring



# Remote Component (@ Amazon)



- Storage Virtual Machine:
  - DIRAC Storage Element/File Catalog
    - cache input data
  - PostgreSQL
    - Belle condition data
- Execution Virtual Machine:
  - Job Agent
  - Virtual Machine Monitor

# Other Components



- User Interface
  - job preparation
  - job submission
- Grid Resources
  - Interfaced by DIRAC
  - Storage Elements
    - KEK, GRIDKA, CYFRONET,...
  - LFC
    - KEK

# Pilots -> Virtual Machines (I)



## DIRAC Execution Virtual Machine:

- Replaces the pilot.
- Can be installed and configured as necessary.
  - Full root access at boot time.
  - Limit user access at execution time.
- Executes:
  - DIRAC Configuration (private) Slave:
    - reduce bootstrap time for any DIRAC component.
  - DIRAC Job Agent:
    - matches and executes payload.
  - **DIRAC VirtualMachineMonitorAgent:**
    - Report status and usage of VM.
    - Asynchronously uploads output data.
    - Halts the VM if requested or if idle.

# Pilots -> Virtual Machines (II)



New DIRAC WMS extensions:

- **Virtual Machine Manager (+DB):**

- Persists information relative to Virtual Machines
- Receives heart beats from VM's (Monitor Agent)
- Monitors status

- **Virtual Machine Scheduler:**

- Checks:
  - TaskQueues (pending Jobs)
  - Running/Submitted VM's
  - Configuration Parameters
- Decide whether new VM's have to be requested:
  - requires appropriated credentials for provider (Amazon)
- Registers in the DB the submitted VM's.
- Can be executed remotely (for absolute privacy of credentials).

# And there we go !



- Tasks taken from Belle official MC requests.
- Preparations:
  - “Execution VM” image with preinstalled Belle SW and **DIRAC**.
  - “SE VM” image with **DIRAC FC/SE** and postgresql DB.
- Input at KEK:
  - collection of bash scripts
  - event description and backgrounds (1 GB/ 1 M evt)
- Output sent back to grid SE’s (**GRIDKA & KEK**).
  - 22 GB / 1 M evt

# Starting up (April 14th)



File Edit View History Bookmarks Tools Help

https://belle01.ecm.ub.es/DIRAC/Belle-Production/dirac\_admin/jobs/JobMonit amazon Ec2 cost

Most Visited Getting Started Latest Headlines LHCb Guía TV - Programa...

Manage ... Jobs ... Data Op... Virtual M... Elasticfox Product... WMS his... Job plots ... Pilot plot... Ama > +

Systems Jobs Production Data Web Tools Virtual machines Help Selected setup: Belle-Production LHCb WMS

**JobMonitoring**

JobId Status MinorStatus ApplicationStatus Site JobName LastUpdate [UTC] LastSignOfLife [UTC]

JobId	Status	MinorStatus	ApplicationStatus	Site	JobName	LastUpdate [UTC]	LastSignOfLife [UTC]
670	Running	Job Initialization	Unknown	DIRAC.Amazon.us	e000049r000702	2010-04-14 17:27	2010-04-14 17:2
385	Running	Job Initialization	Unknown	DIRAC.Amazon.us	e000049r000120	2010-04-14 17:23	2010-04-14 17:2
1030	Waiting	Pilot Agent Submis	Unknown	DIRAC.Amazon.us	e000045r000448	2010-04-14 14:42	2010-04-14 14:4
1031	Waiting	Pilot Agent Submis	Unknown	DIRAC.Amazon.us	e000045r000449	2010-04-14 14:42	2010-04-14 14:4
1032	Waiting	Pilot Agent Submis	Unknown	DIRAC.Amazon.us	e000045r000450	2010-04-14 14:42	2010-04-14 14:4
1022	Waiting	Pilot Agent Submis	Unknown	DIRAC.Amazon.us	e000045r000435	2010-04-14 14:42	2010-04-14 14:4
1023	Waiting	Pilot Agent Submis	Unknown	DIRAC.Amazon.us	e000045r000436	2010-04-14 14:42	2010-04-14 14:4
1021	Waiting	Pilot Agent Submis	Unknown	DIRAC.Amazon.us	e000045r000429	2010-04-14 14:42	2010-04-14 14:4
1019	Waiting	Pilot Agent Submis	Unknown	DIRAC.Amazon.us	e000045r000372	2010-04-14 14:42	2010-04-14 14:4
1020	Waiting	Pilot Agent Submis	Unknown	DIRAC.Amazon.us	e000045r000428	2010-04-14 14:42	2010-04-14 14:4
1017	Waiting	Pilot Agent Submis	Unknown	DIRAC.Amazon.us	e000045r000369	2010-04-14 14:42	2010-04-14 14:4
1018	Waiting	Pilot Agent Submis	Unknown	DIRAC.Amazon.us	e000045r000371	2010-04-14 14:42	2010-04-14 14:4
1015	Waiting	Pilot Agent Submis	Unknown	DIRAC.Amazon.us	e000045r000364	2010-04-14 14:42	2010-04-14 14:4
1016	Waiting	Pilot Agent Submis	Unknown	DIRAC.Amazon.us	e000045r000367	2010-04-14 14:42	2010-04-14 14:4
1014	Waiting	Pilot Agent Submis	Unknown	DIRAC.Amazon.us	e000045r000363	2010-04-14 14:42	2010-04-14 14:4

Submit Reset Global Sort Current Statistics Global Statistics

Page 1 of 31 Items displaying per page: 25 Displaying 1 - 25 of 752

jobs > Job monitor ricardo@ dirac\_admin (/DC=es/DC=irisgrid/O=ecm-ub/CN=Ricardo-Graciani-Diaz)

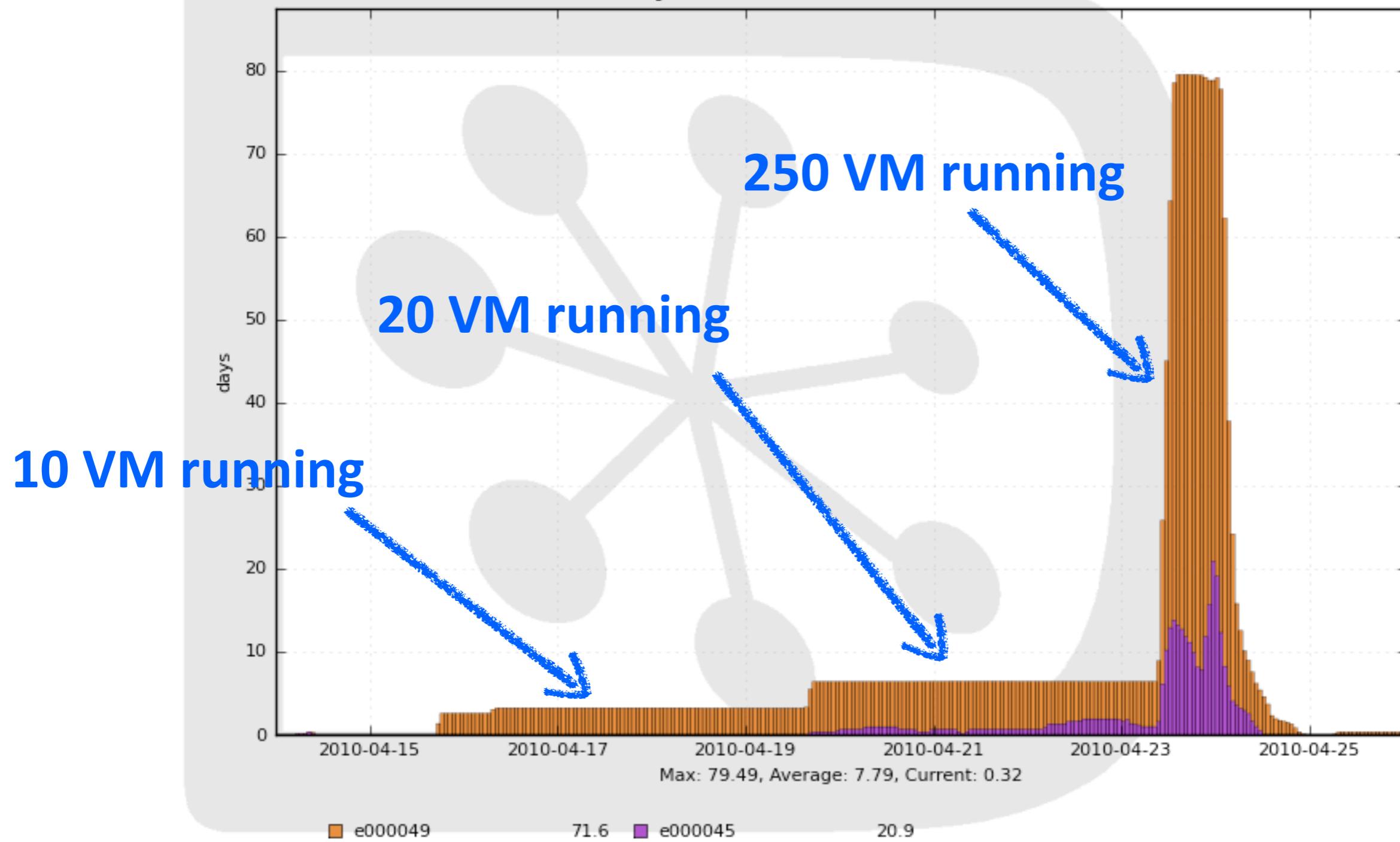
https://belle01.ecm.ub.es/DIRAC/Belle-Production/dirac\_admin/jobs/JobMonitor/display#

# The execution: Phase I

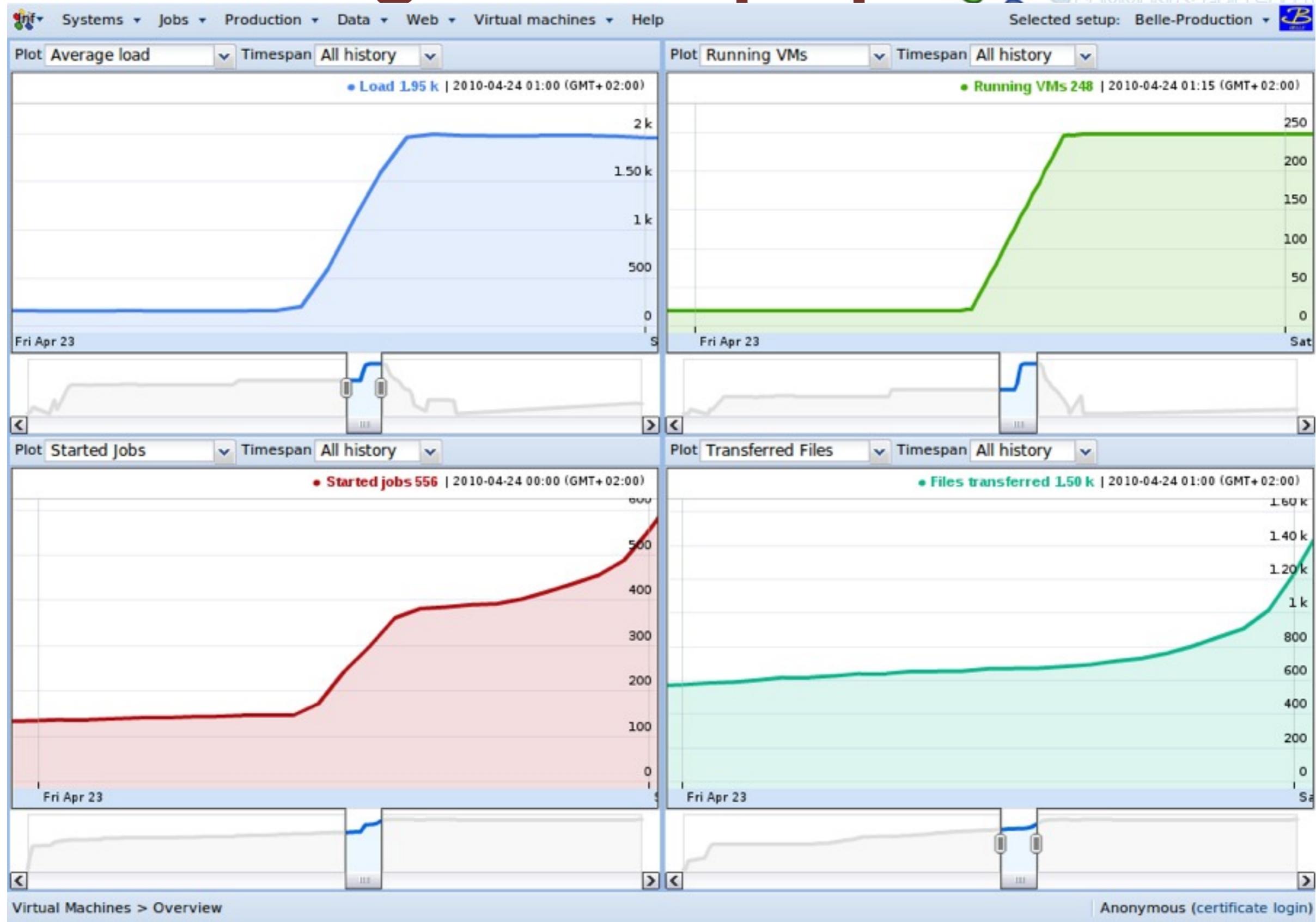


CPU days consumed by simulation Experiment / hour

12 Days from 2010-04-13 to 2010-04-25



# Monitoring the ramp up

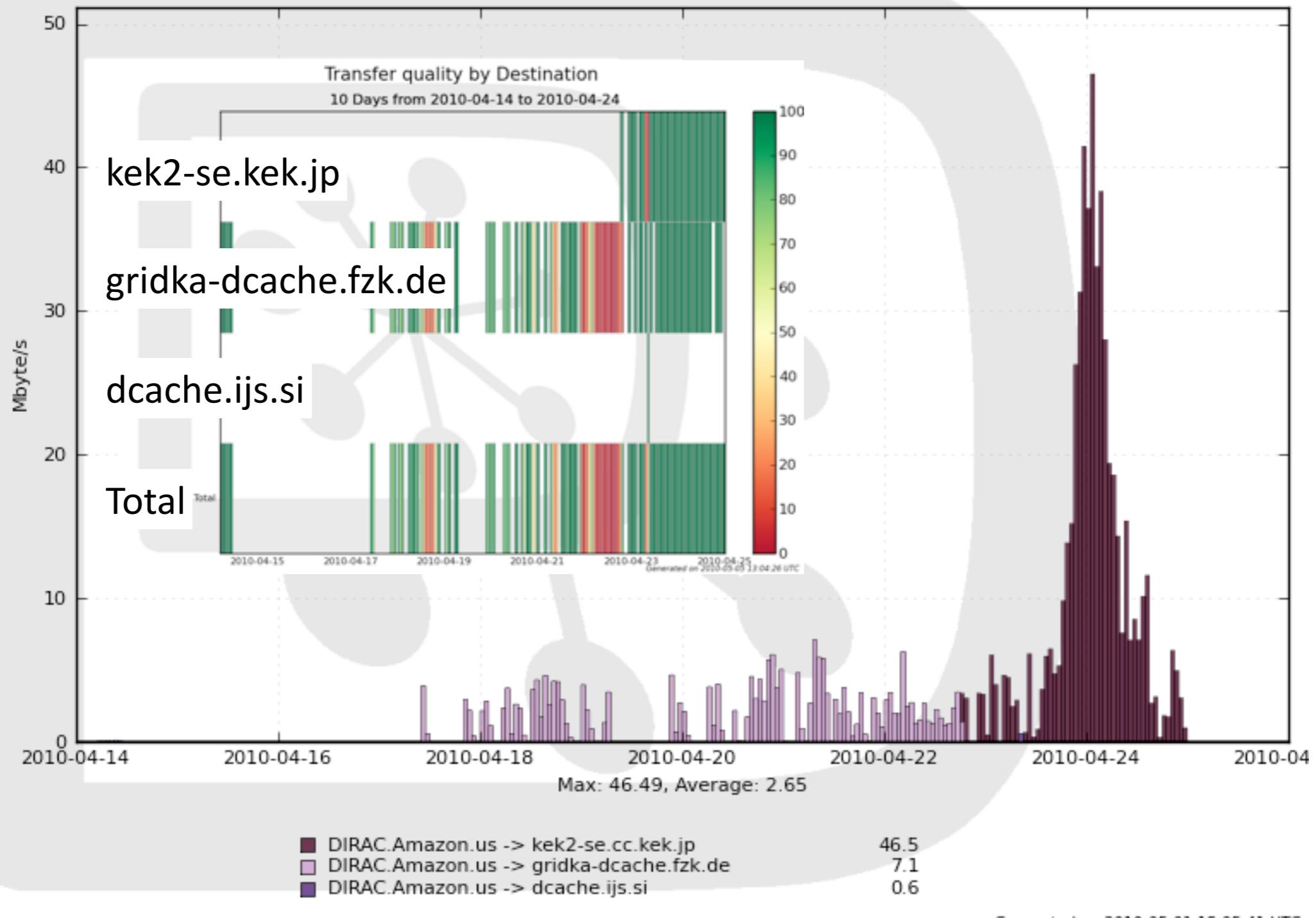


# And the data back to grid



Transferred data by Channel

11 Days from 2010-04-13 to 2010-04-25



# Results (I)



- Phase I (cloud test):

production ready:

- 5% of Belle production in 10 days
- 120 M evt (~2.7 TB)
- 2250 CPU days used

proven stability and scalability:

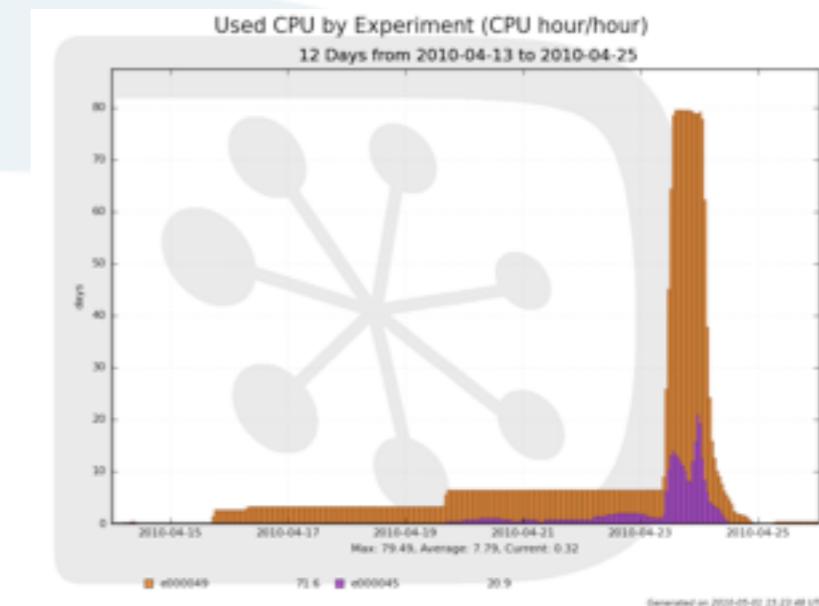
- 2000 CPUs peak achieved in < 4 hours
- > 90 % efficiency in CPU usage

– first cost estimation:

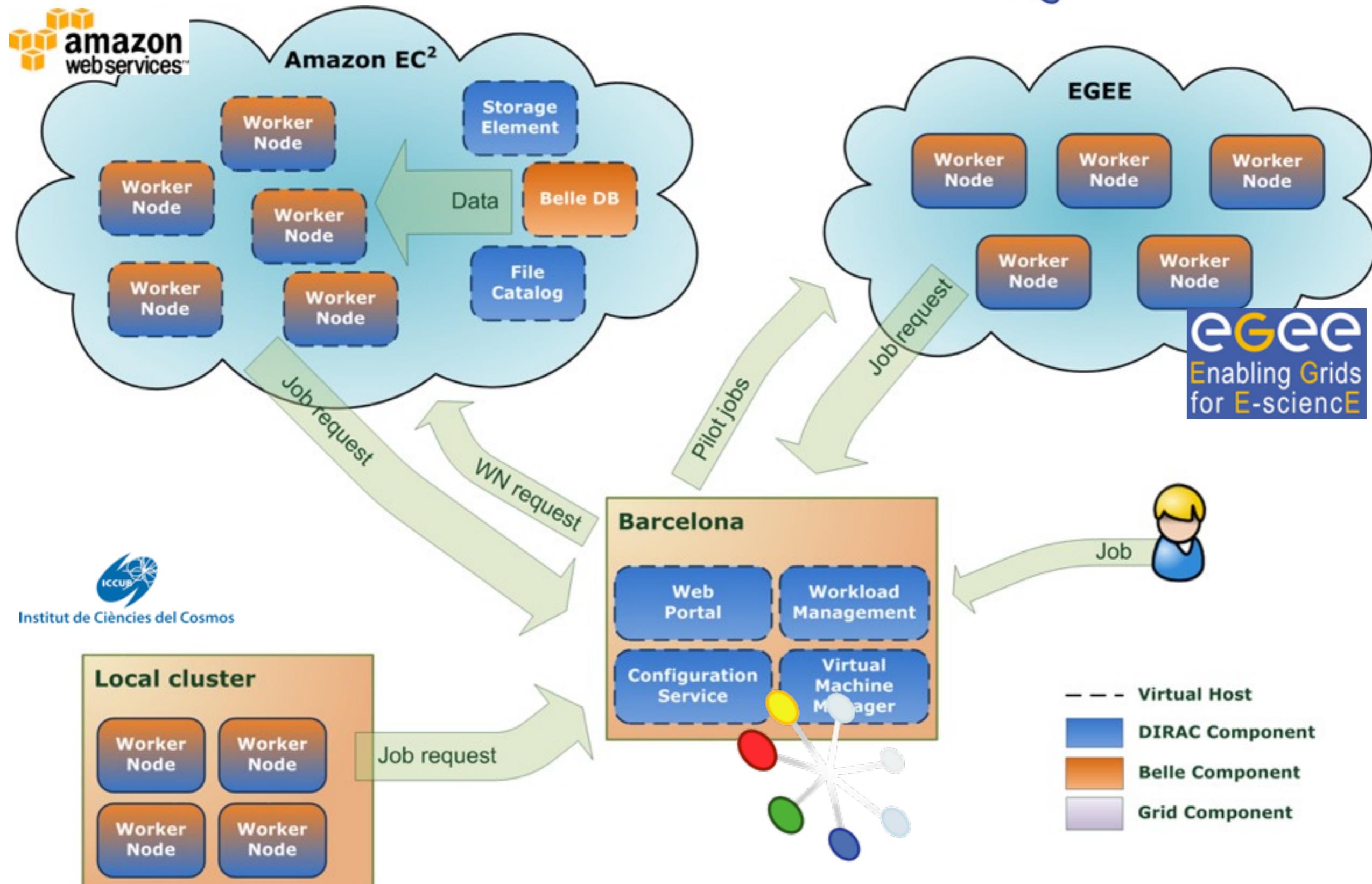
- 0.46 USD/10k evt

– input data pre-uploaded to Amazon SE VM.

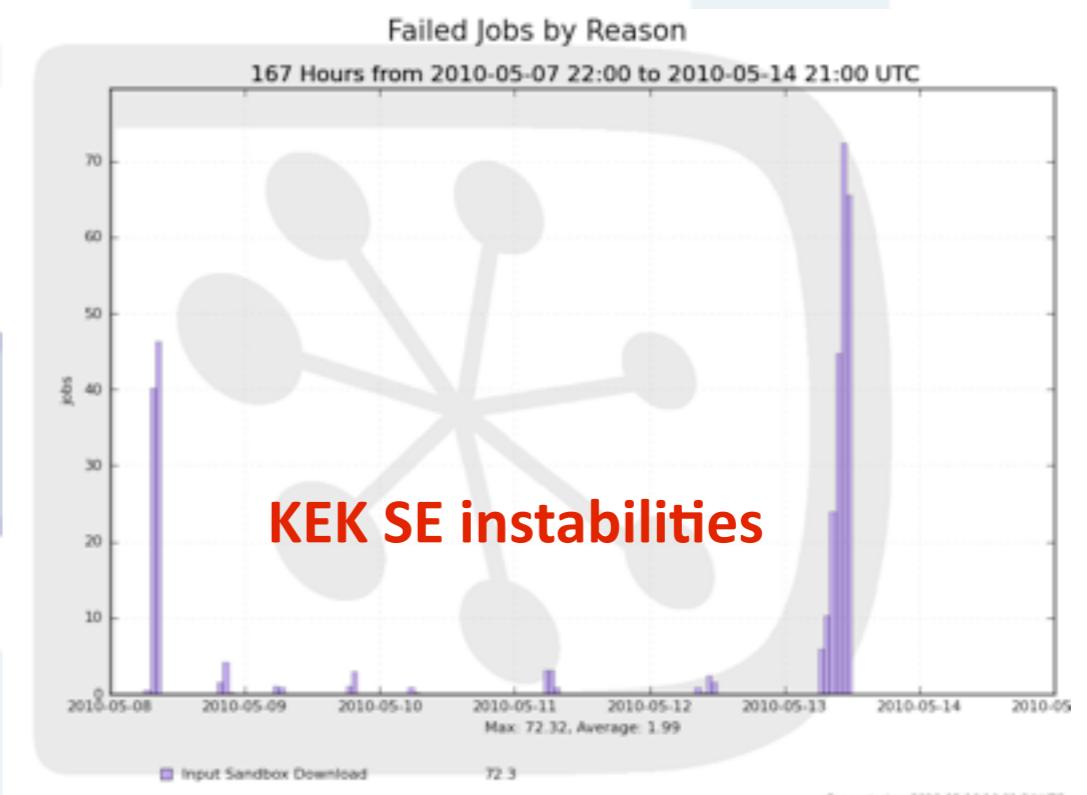
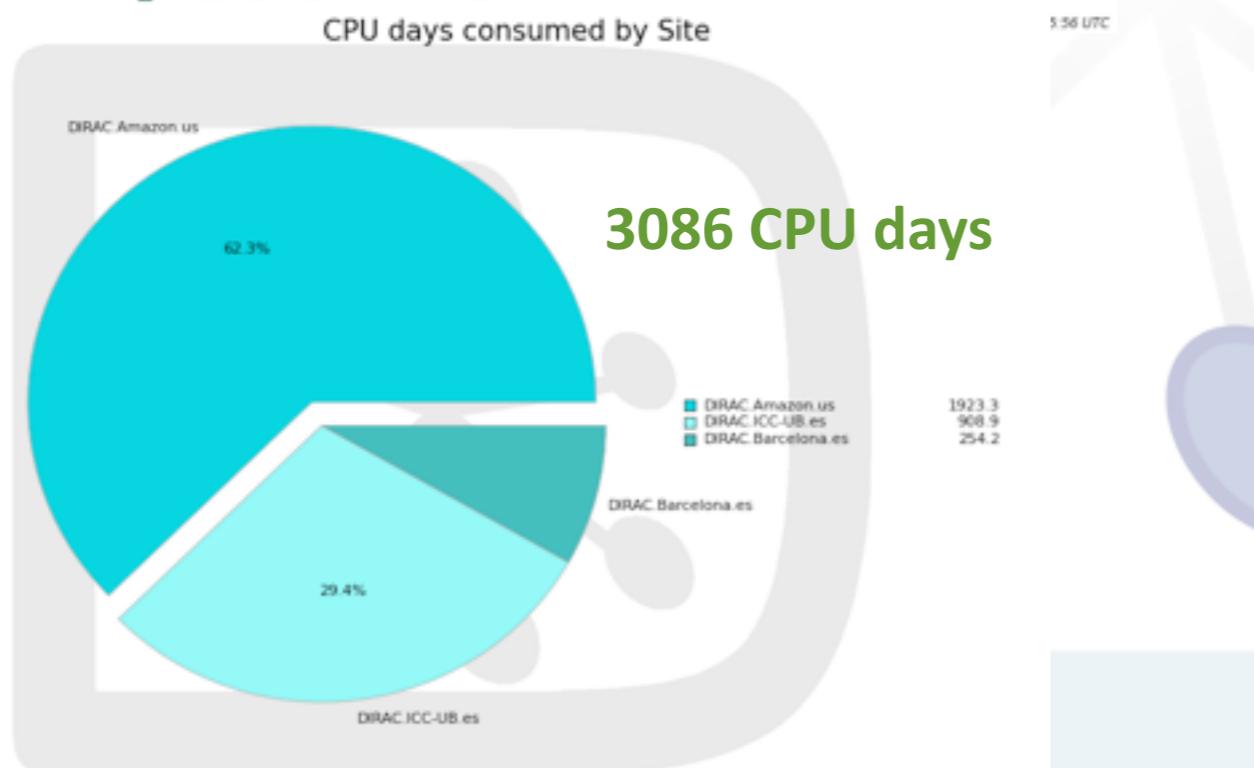
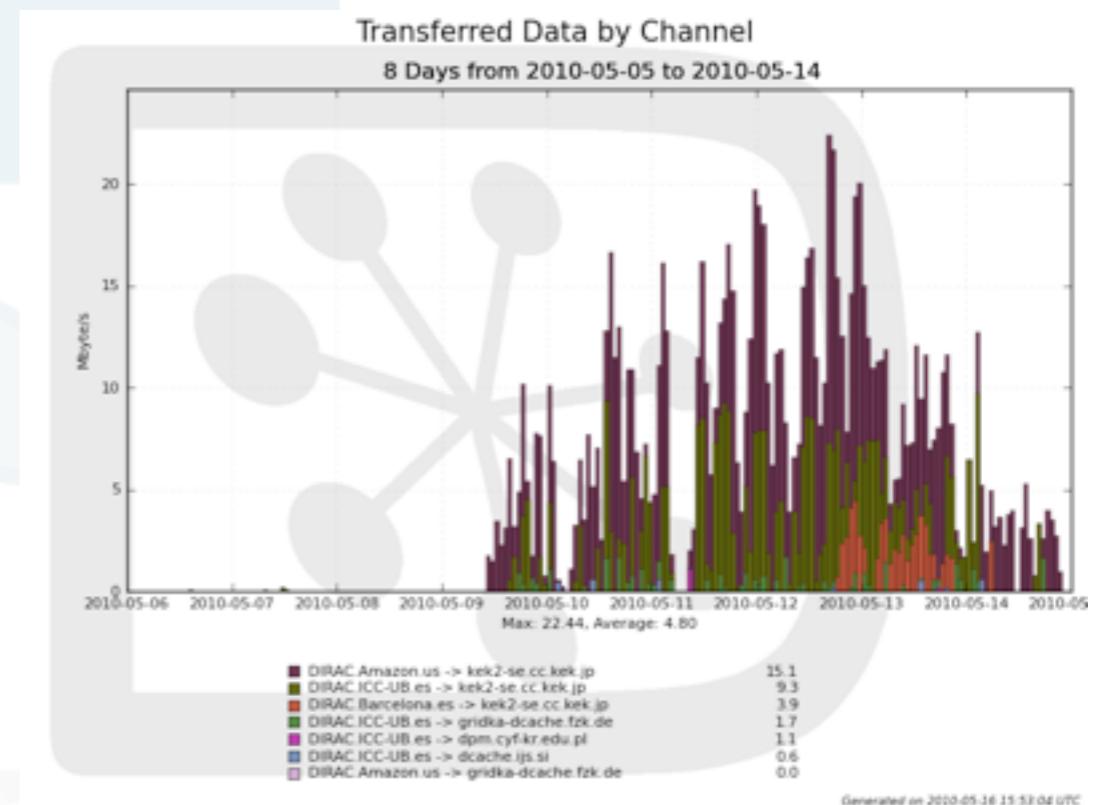
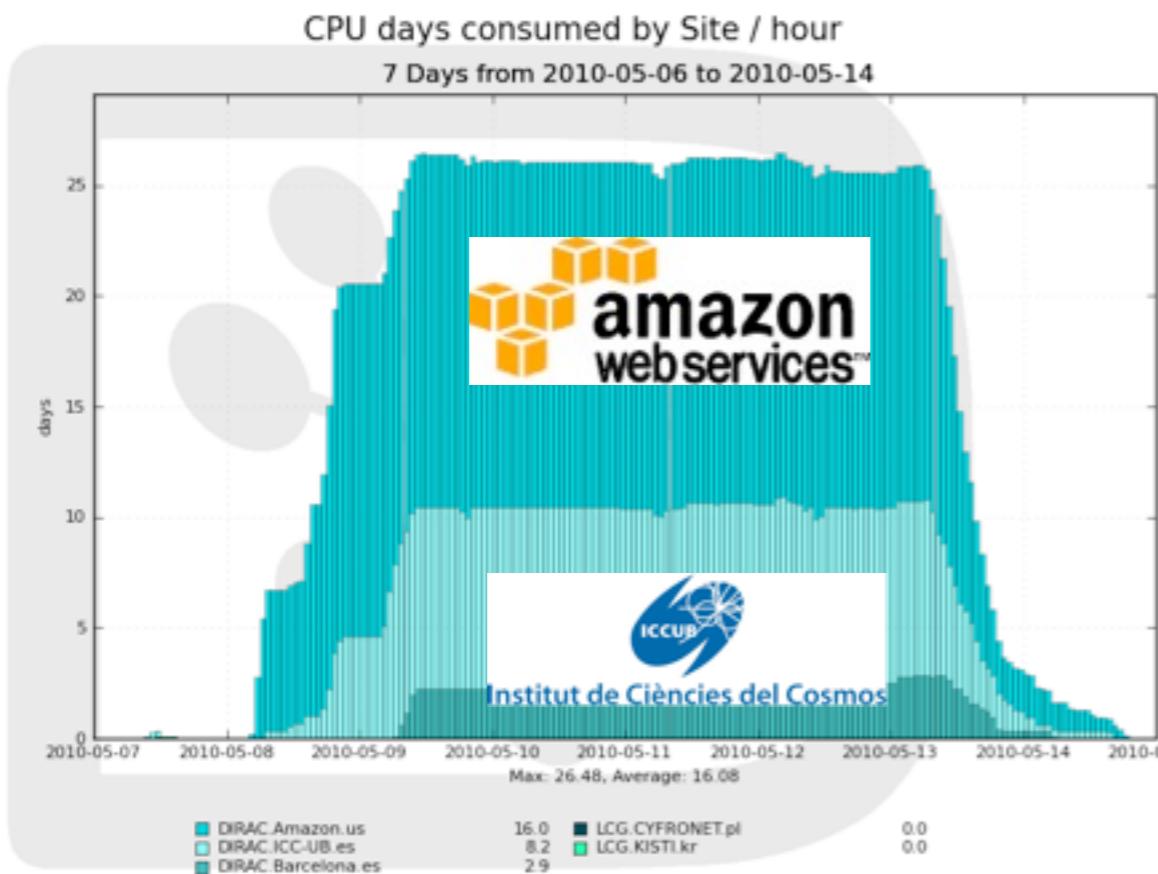
– few bug fixes



# Next steps



# Cloud + Local: Phase II



# Results (II)



- Phase II (local + cloud integration):

- production ready:

- 7% of Belle production in 6 days
    - 170 M evt (~3.6 TB)
    - 3100 CPU days

- proven interoperability:

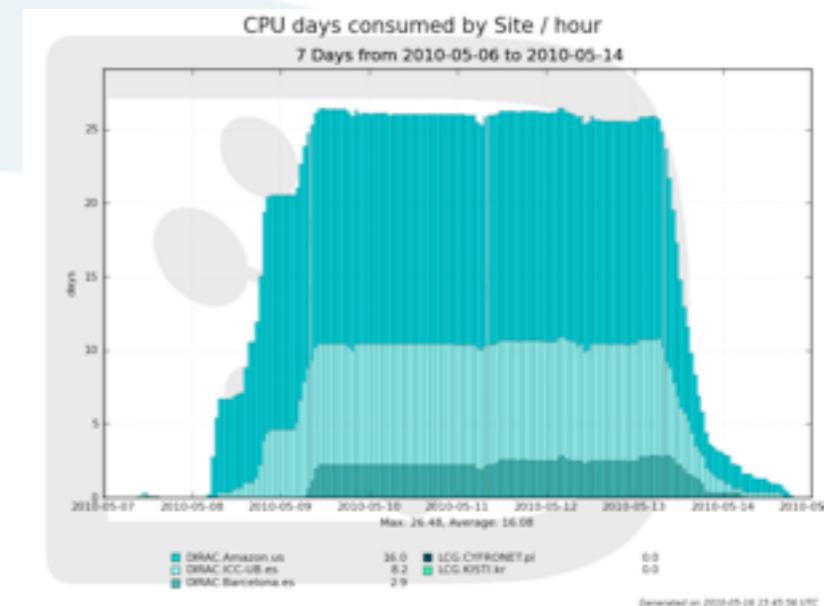
- 60 % cloud resources / 40 % local resources
    - transparent integration of heterogeneous resources
    - > 95 % efficiency in CPU usage

- improved cloud cost:

- 0.20 USD / 10k evt

- input data downloaded from KEK SE

- issues with stability of grid Storage



# On Amazon Cost



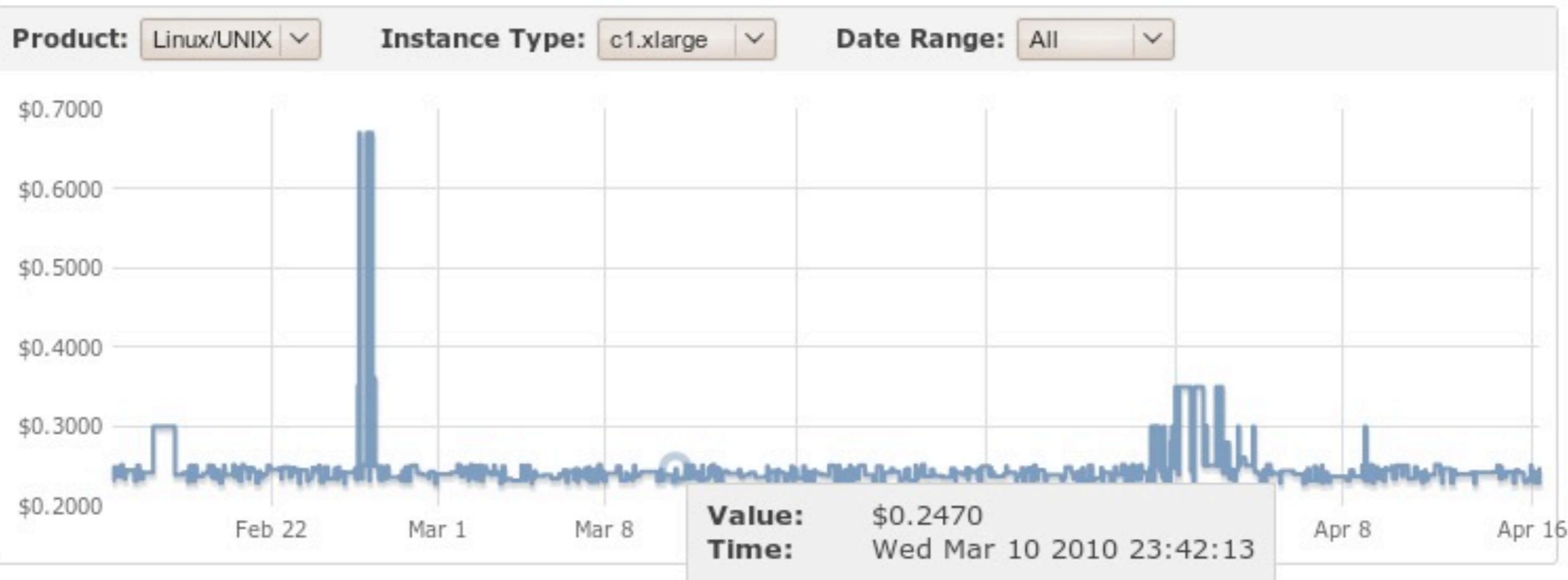
- In all cases we used the High-CPU Extra Large Virtual Machine, that offers the best CPU/cost rate:
  - 7 GB of RAM
  - 8 Cores (20 EC2 Units or 67 HEP Spec 2006)
  - 1.7 TB disk
- Purchase options:
  - On-demand (most expensive): **0.68 USD/hour**
    - fix price, request & pay
  - Reserved (cheapest): **0.24 USD/hour**
    - fix price (1-3 years), pay & use
  - Spot:
    - price fluctuates, request (with top price) & pay

# Spot Instance Price History



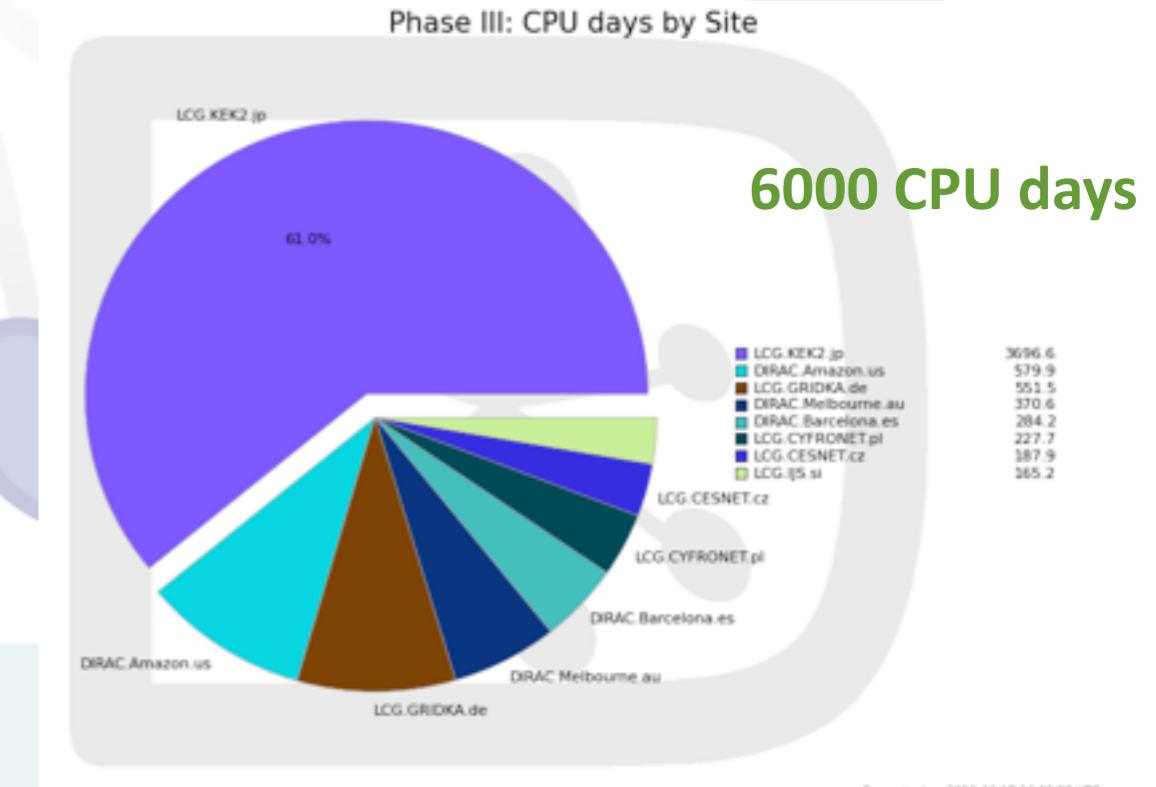
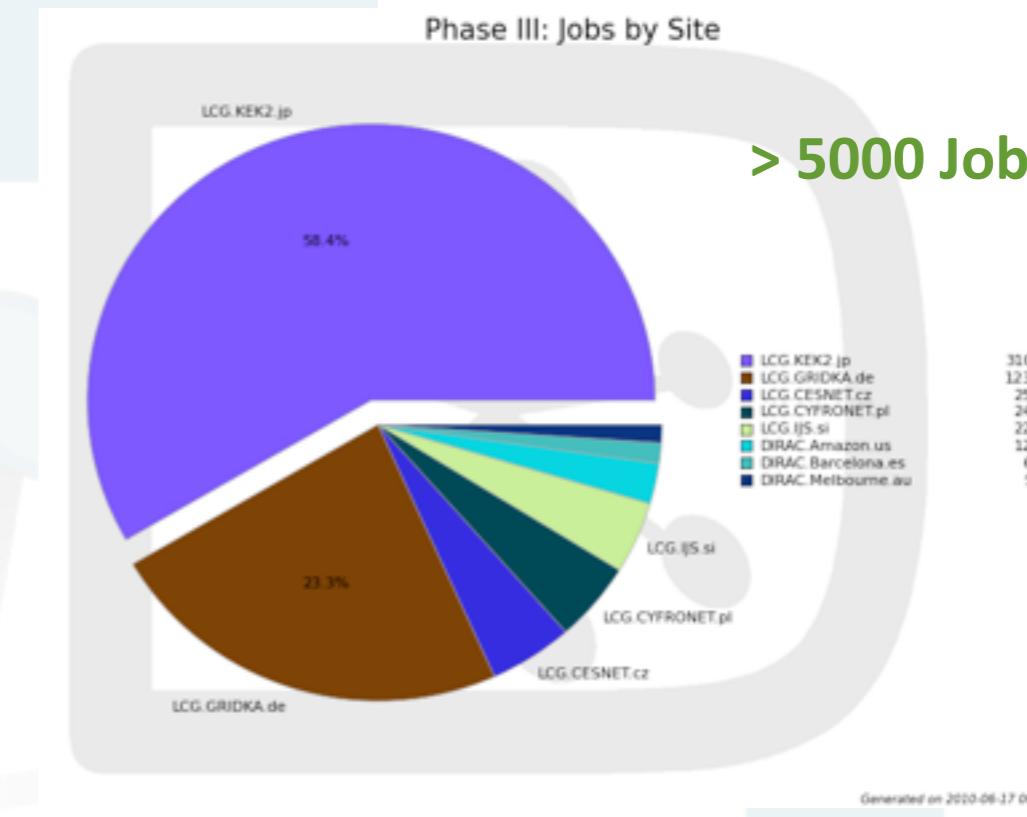
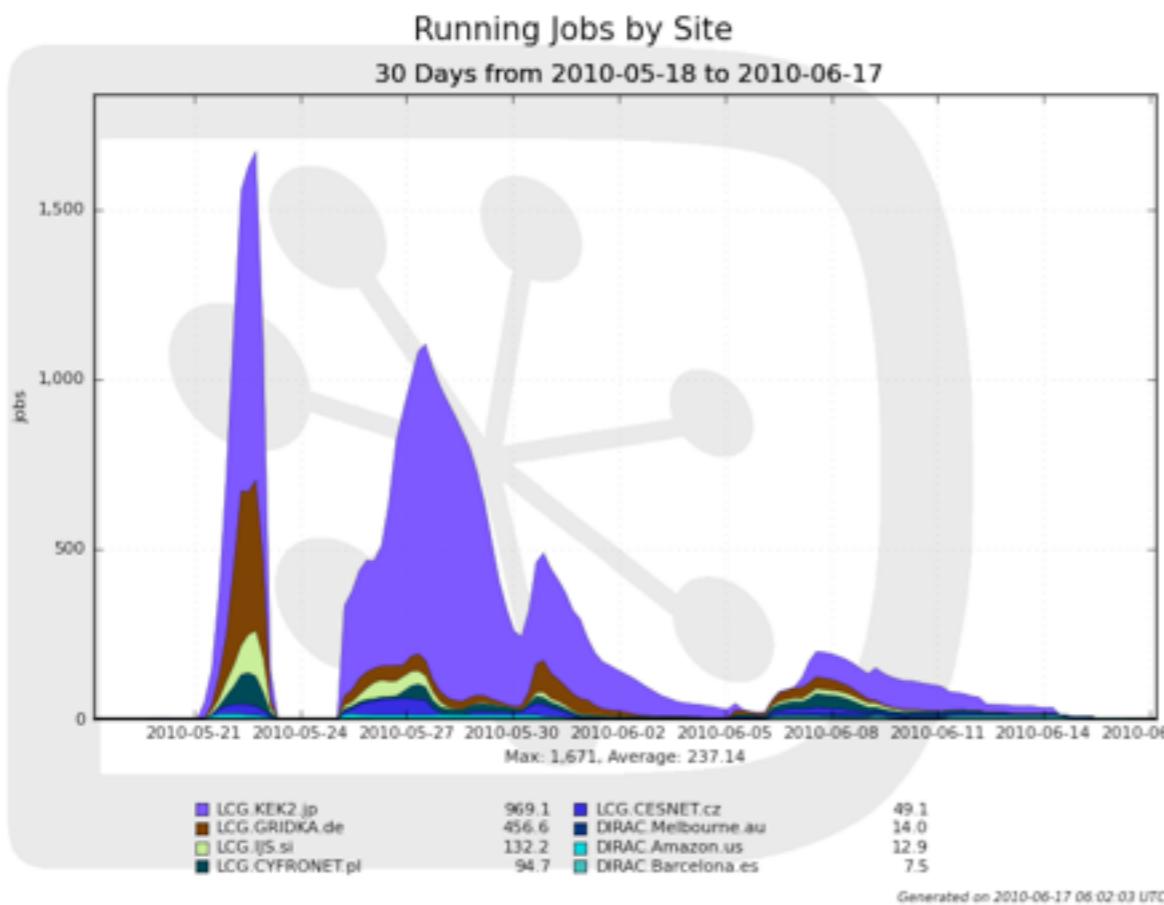
## Spot Instance Pricing History

Cancel



[Close](#)

# Cloud + Local + Grid: Phase III



Different issues with SE at KEK:

- gridftp hosts restarting.
- limited number of gridftp servers.
- control channel timeout.
- SE getting full.

Random SW errors when executing on the grid

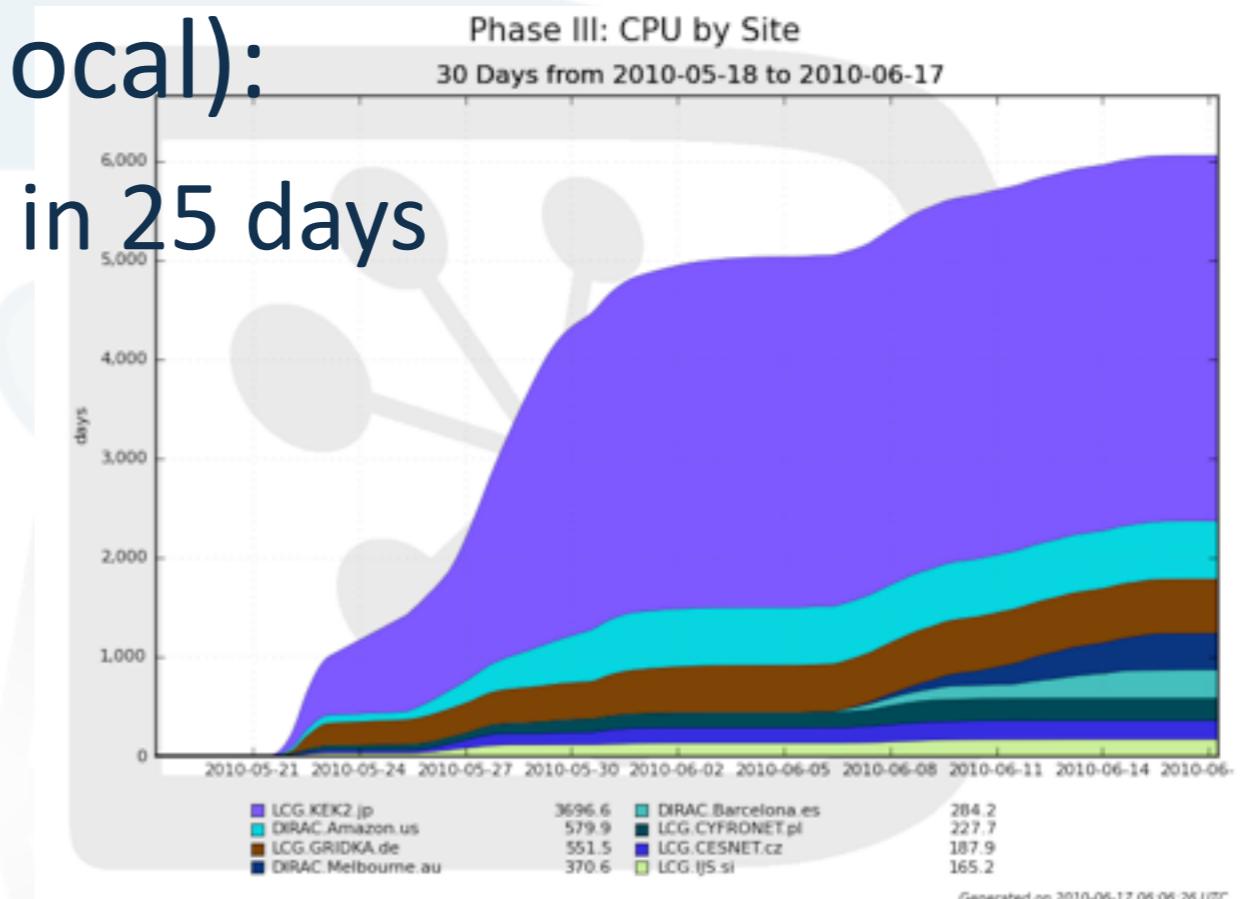
- Not seen in local or cloud (multi-process)

Issues with queue lengths

# Results (III)



- Phase III (grid + cloud + local):
  - ✓ 15% of Belle production in 25 days
    - Still running:
      - 6000 CPU days
    - ✓ Full interoperability:
      - cloud + grid + local
      - stability issues on grid SE.
      - main effort to make sure grid environment is appropriate
    - Other issues:
      - ✓ Jobs can be very long, 1-30 days, only the shorter ones are appropriate for the grid



# Summary



- Perfect behavior of DIRAC - Amazon integration:
  - >95% CPU usage efficiency.
  - No errors due to Amazon.
  - 0.20 USD / 10k evt (2009 estimate: 0.68 USD / 10k evt)
- Perfect integration of Local resources:
  - >95% CPU usage efficiency.
  - No errors due to Local clusters.
- Issues with grid resources:
  - Need to add extra redundancy and checks
    - Integrate execution in DIRAC framework
    - Add failover SE's
    - Need to test and validate resources (CE's and SE's)
- Main issue access to “large” input data files