



# CMS

## Monte Carlo production in the WLCG Computing Grid

José M. Hernández

CIEMAT, Madrid

P. Kreuzer

RWTH, Aachen

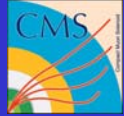
On behalf of the CMS production group

**XVI International Conference on Computing in  
High Energy and Nuclear Physics**



# Outline

- Introduction
- Experience with former production system
- The new production system ProdAgent
- Production operations and performance
- New production components and prospects
- Summary and conclusions



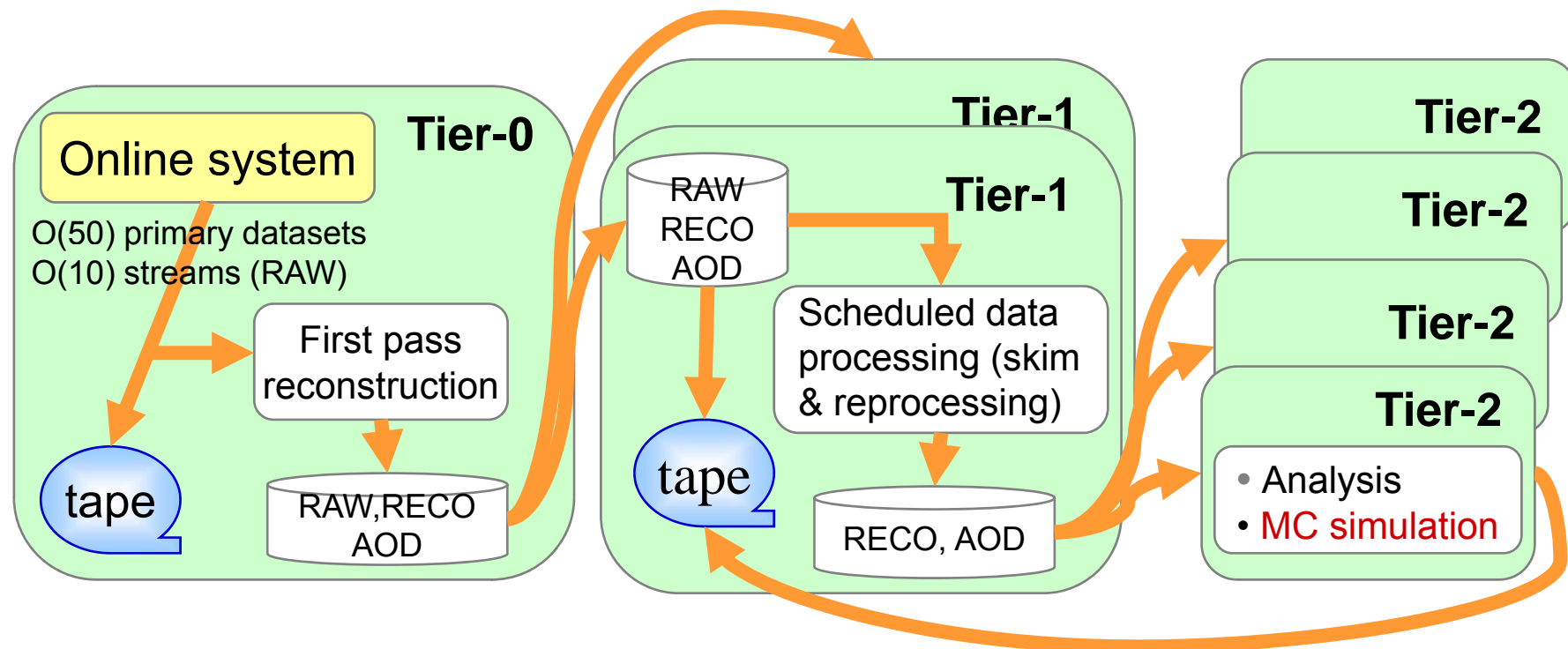
# Introduction

- Efficient and performant MC production system crucial for delivering large simulated data samples for detector studies and physics analyses
- WLCG provides a large amount of distributed computing, storage and network resources for data processing
- Reliability and stability in Grid services and at the sites are key issues given the complexity of services and heterogeneity of resources
- A robust, scalable, automated and easy to maintain MC production system that can make an efficient use of the LCG resources is mandatory
- Major boost in performance and scale in CMS MC production since last CHEP conference
- Production system reengineering to incorporate experience gained in running the previous system. Production organization tighten
- Integration aspects, operational experience and production performance will be presented



# The CMS Computing Model

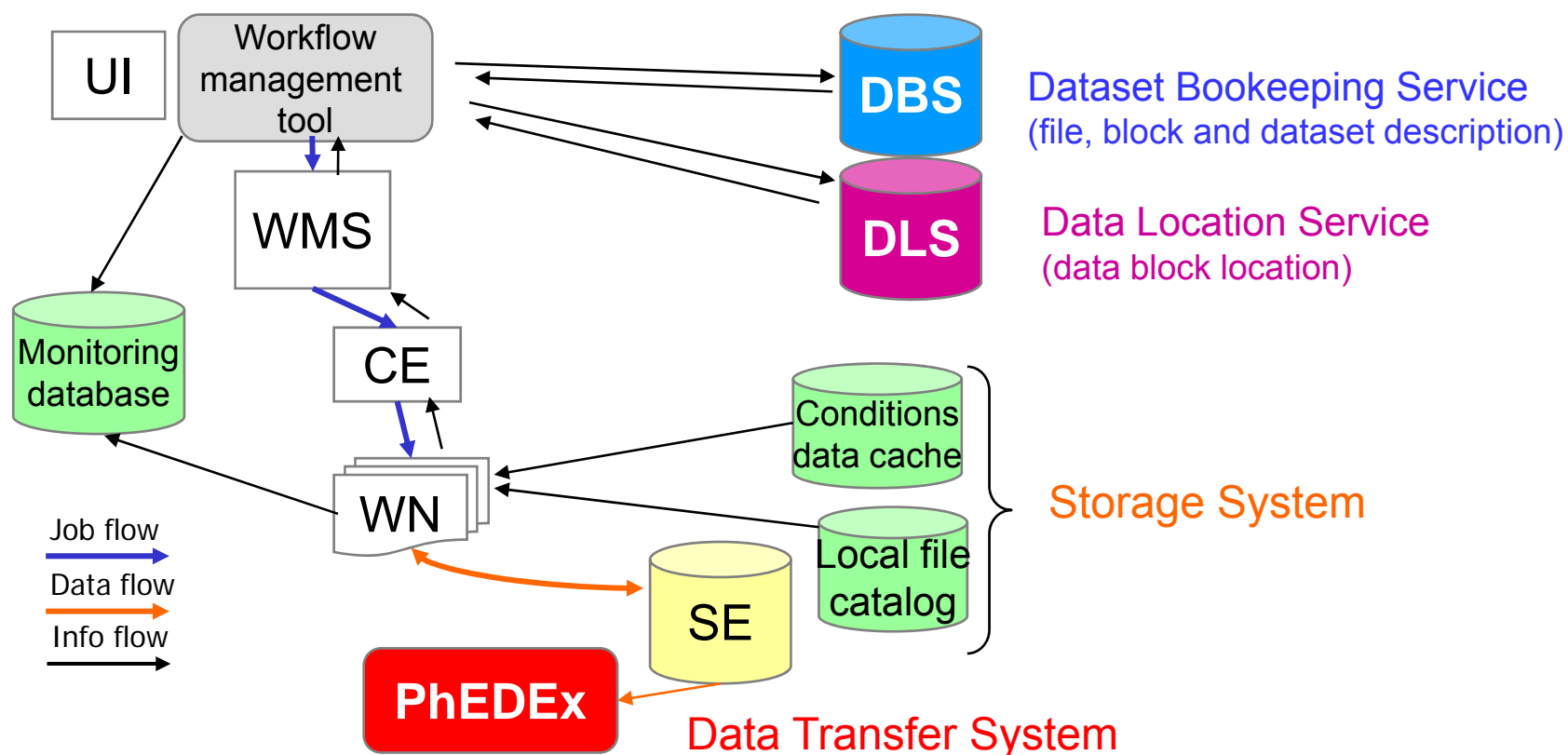
- Distributed computing model for data storage, processing and analysis
- Grid technologies (Worldwide LHC Computing Grid, WLCG)
- Tiered architecture of computing resources
- Several Petabytes of data (real and simulated) every year





# Data processing workflow

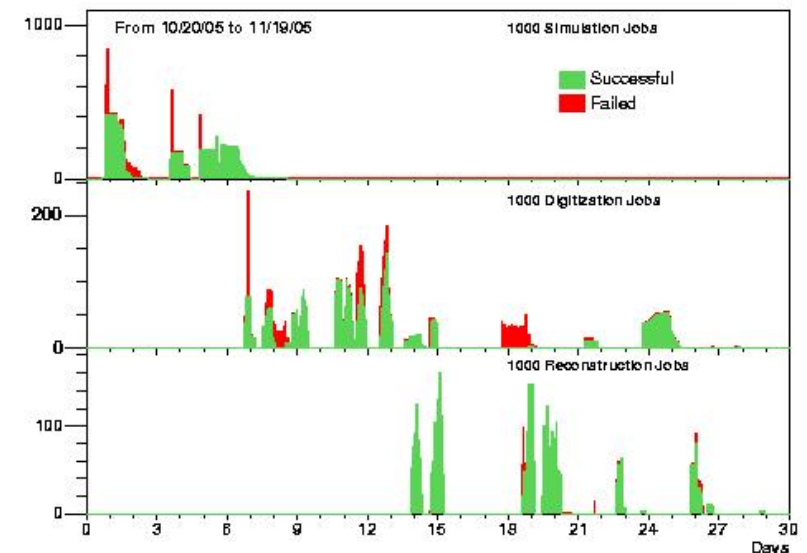
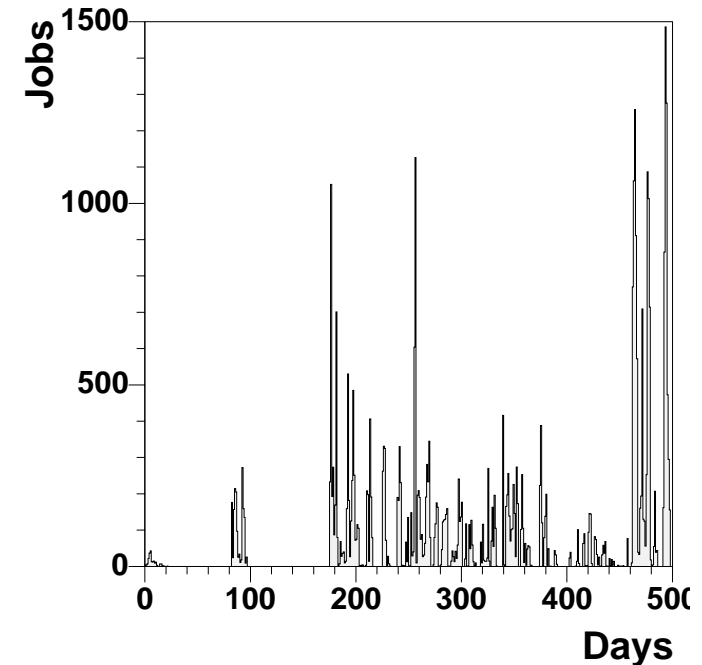
- Data Management System to discover, access and transfer event data in a distributed computing environment
  - Fileblock-based data replication and tracking
  - Local (trivial) file catalogue
  - Merging processing step
  - Data location based processing





# Production with old system (McRunjob)

- Porting of McRunjob production system to the LCG Grid presented in past CHEP
- Limited scale (~few hundred jobs/day, ~3Mevt/month)
  - Poor automation in job and error handling, limited monitoring, poor coupling to DMS
  - Lack of robustness in handling Grid and site inefficiencies
  - Complex processing and publication workflows (no processing step chaining, metadata creation in old CMS event data model)
  - Not efficient use of available resources
  - Only a fraction of computing resources available via LCG (significant production in local farms and OSG US Grid)





# Experience from old production system

- Very valuable experience gained implementing and running McRunjob on the Grid input for the design of new production system
- Proper automation, monitoring and error handling necessary to deal with processing in a distributed environment with inherent instability and unreliability. Robustness is key
- Good support and responsiveness from sites is critical. White list strategy
- Need tools for services and infrastructure availability monitoring
- Robustness and tuning in data I/O is very important
- Small files bad for storage, transfer and cataloguing
- Proper job prioritization (production/analysis) required
- Need Grid bulk operations (submission, tracking)
- Central software (pre-)installation in shared file system OK
- A global replica catalogue did not scale
- Need to simplify processing and publication workflows
- Production on the Grid is a manpower-intensive task



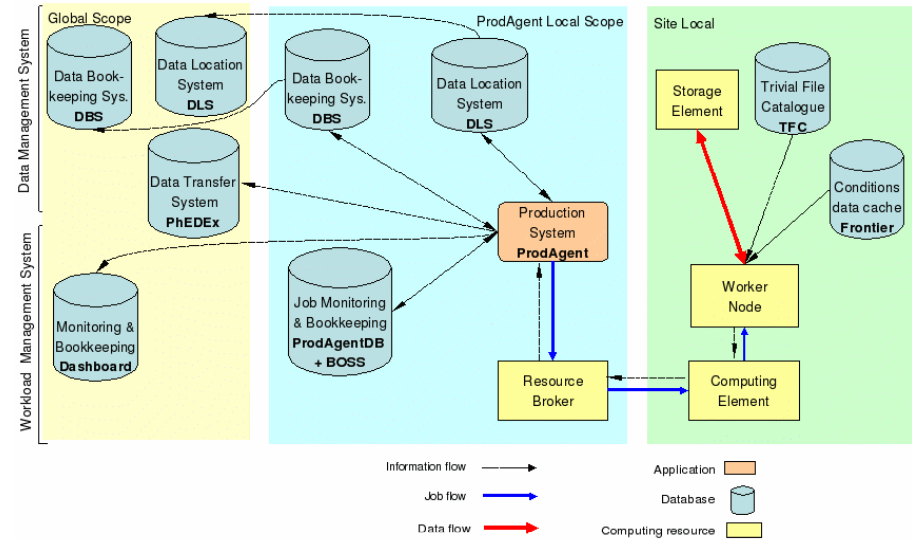
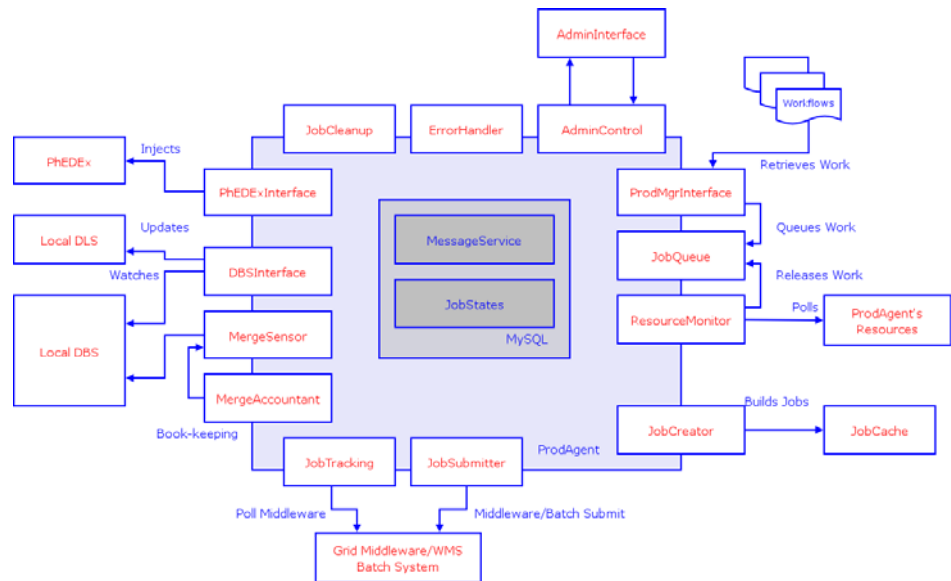
# The new production system: ProdAgent

- Design aiming at automation, ease of maintenance, scalability, avoid single points of failure, modular architecture with support to multiple Grids
- Automation in job handling, data discovery and registration, proper monitoring, accounting and error handling
- Built-in robustness in interactions with Grid and site services
- Integration with new CMS Event Data Model, Data Management System and event processing framework
  - Simplified processing and publication workflows (chaining of processing steps possible, no event metadata anymore)
  - Fileblock-based data location and replication, merging step of output files in workflows, use of local trivial file catalogue (simple LFN to local PFN rules)
- Data location based processing
  - Send processing jobs to the data prelocated by transfer system
  - Read files directly from storage via posix I/O local access protocols





# ProdAgent architecture and workflow



- Work split in atomic loosely coupled components
- Data processing, bookkeeping, tracking and monitoring occurs in local-scope
- Data bookkeeping and location information promoted to global-scope databases and data transfer system after successful processing
- Scaling achieved by running any number of concurrent ProdAgent instances
- Many more details in *CMS MC Production System Development & Design* talk



# Production operation

- A good organization is key for an efficient production
- Production manager coordinates work from physics groups (requestors), production teams, developers of processing framework and production
- Before production, software releases, configurations and production system are validated via production of large test samples
- Physics groups enter production requests via web interface with production parameters → Production manager approves and assign them
- Requests are injected into ProdAgent instances. Production sites distributed between a number of production teams (white lists)
- Typically production run either in single step (generation-simulation-digitization-reconstruction) or split into two steps (gen-sim + digi-reco)
- Different levels of production monitoring, bookkeeping and accounting:
  - Local ProdAgent instance monitor (production team)
  - Global production accounting database (production manager)
  - Global CMS Dashboard (sites)
  - Global Dataset Bookkeeping Service (requestors and analyzers)



# Production performance

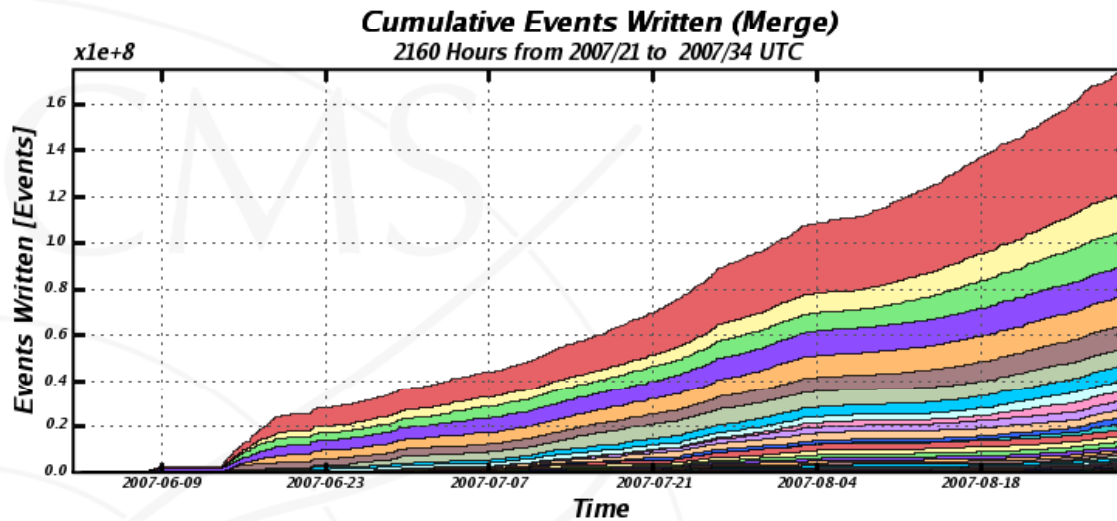
- Overall MC production @ CMS in 2007

Production Period	Number of MC events	Production Rate	Production Teams / PAs
Winter07	90 M	35M/month	6 / 10
Spring07	67 M	46M/month	6 / 9
Summer07	145 M	64M/month	5 / 7

- Production resources
  - 1 T0 / 7 T1s / 34 T2s (active in production)
  - ~ 5k batch slots for T0/T1s
  - ~ 6k batch slots for T2s } up to 11K slots available for production
- MC production by default at T2s but T0/T1 resources also used if available
- Production teams: 1 OSG team, 4-5 LCG teams



# Production performance: yield

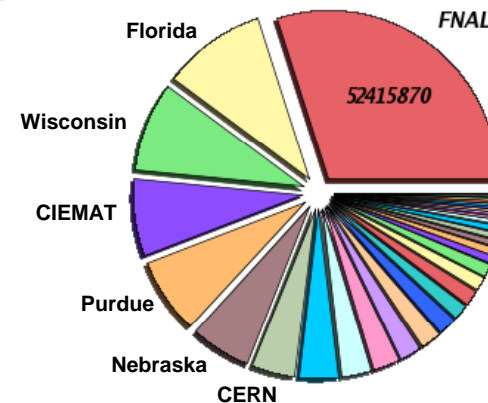


~175 M events in last 3 months  
 ~ 2/3 done at Tier-2 sites  
 ~ 50% done at OSG sites

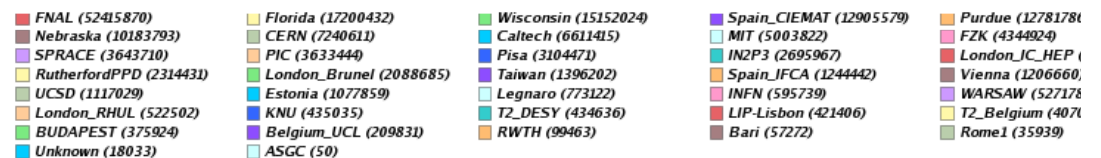


Total: 174953424.00 Events, Average Rate: 22.50 Events/s

**Events Written by Site (Sum: 174953424 Events)**  
12 Weeks from 2007/21 to 2007/34 UTC



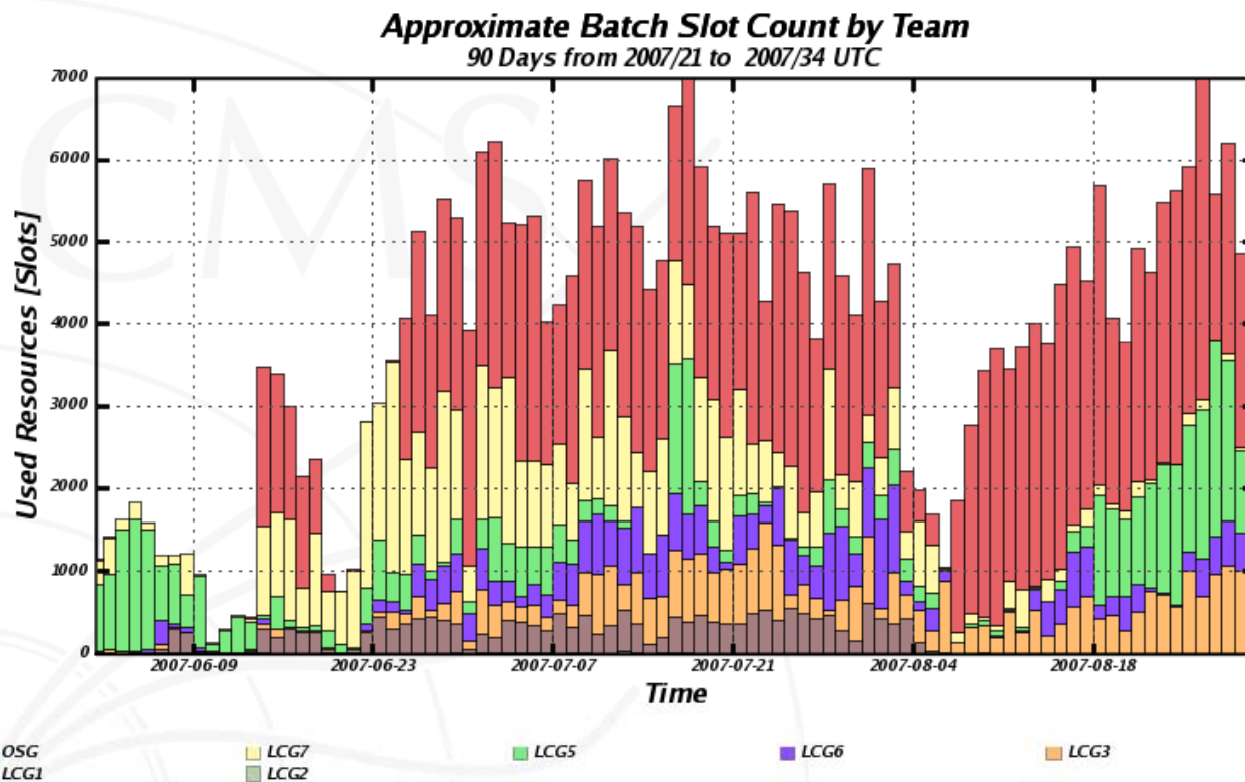
- Processing time:  
GEN-SIM ~40s/evt for Minbias to  
700s/evt for heavy channels
- Event sizes:  
GEN-SIM ~ 0.5 Mbytes  
DIGI-RECO ~ 1.5 MBytes





# Production performance: resource usage

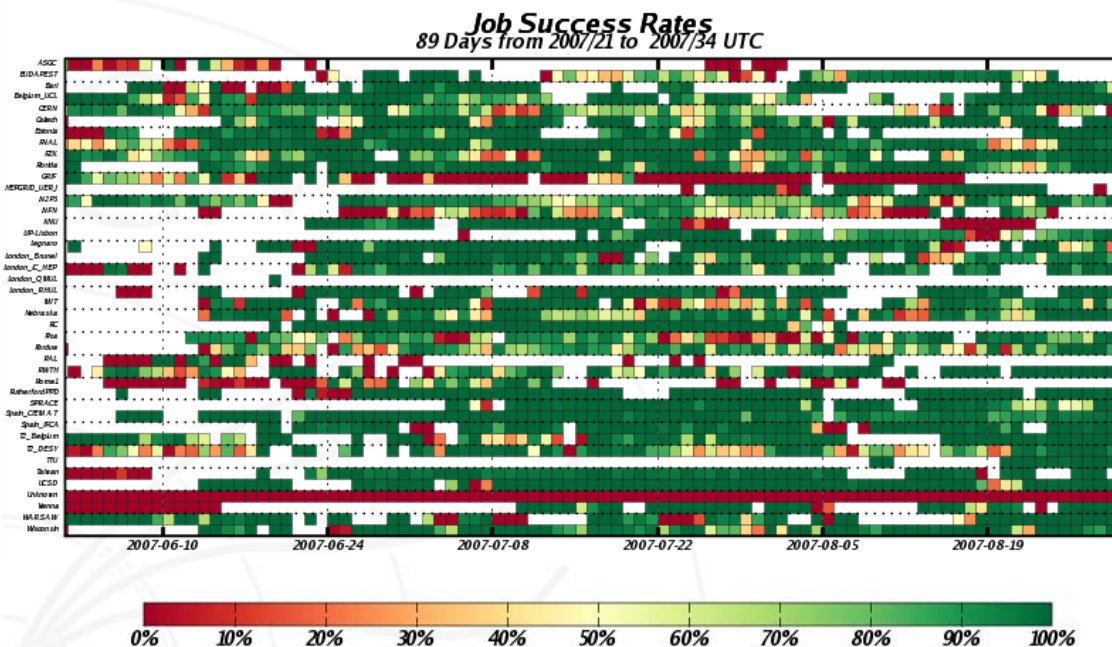
- Approximate slot usage (average number of slots continuously running production)
- Ramp up from June (~ 5000 slots used in average)
- In average ~50% resource occupation (production inefficiencies, no automatic resource management, manual job release by operators, many sites, many small production requests)



Maximum: 6998.60 Slots, Minimum: 131.27 Slots, Average: 3856.94 Slots, Current: 678.10 Slots

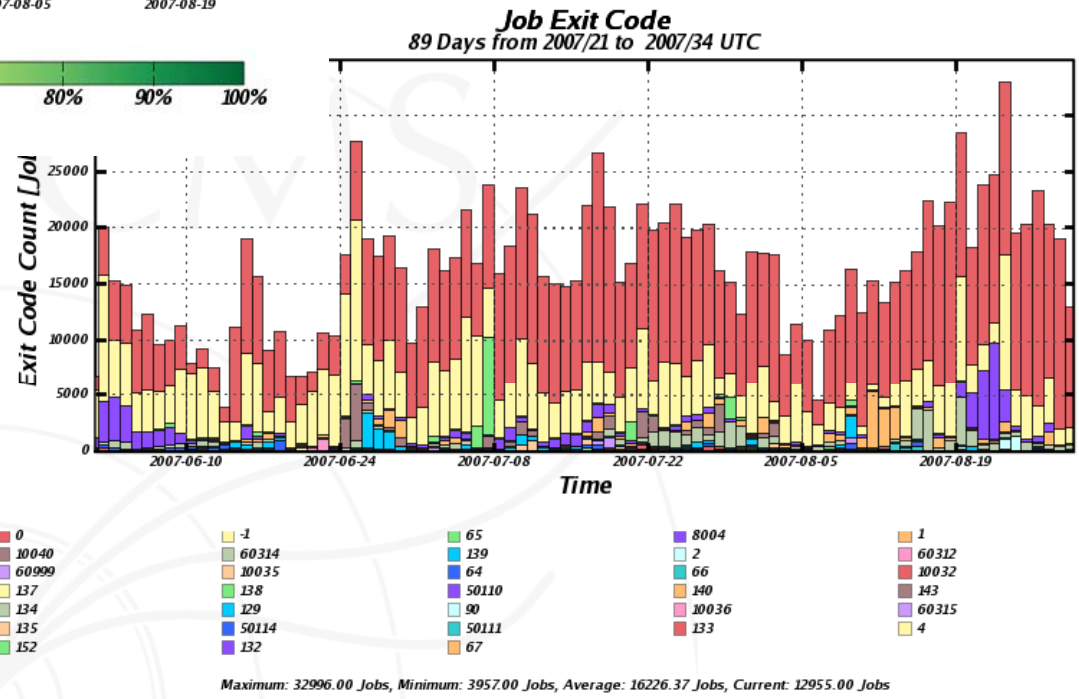


# Production performance: job efficiency



- 40+ sites used continuously for production with high job efficiency
- Average job efficiency ~ 75%, including application and Grid. (Failed jobs are automatically resubmitted and get done)

- > 20K/day production jobs
- Large fraction of aborted jobs by Grid





# New production components

- New production components under test, integration of deployment to further increase production scale and automation
- Workflow management automation with ProdManager to supply work to ProdAgent instances from requests entered in ProdRequest
- Better use of available resources with ResourceMonitor
- Internal JobQueue in ProdAgent to buffer jobs waiting for resources
- Bulk operations for job creation, submission and tracking
- Automatic multi-step processing automation with ProcSensor
- Better monitoring and accounting with ProcMon





# Summary and conclusions

- Major boost in performance and scale in CMS MC production since last CHEP conference
- Re-factored production system has brought automation, scalability, robustness and efficiency in handling the CMS distributed production system
- Much improved production organization, bringing together requestors, consumers, producers, developers and sites, has also contributed to the increase in scale
- Reached scale of more than 20K jobs and 65 Mevt/month with an average job efficiency of about 75% and resource occupation ~ 50%
- Production is still manpower intensive. New components being integrated to further improve automation, scale and efficient use of available resources while reducing required manpower to run the system