

Analysis Tools and Support

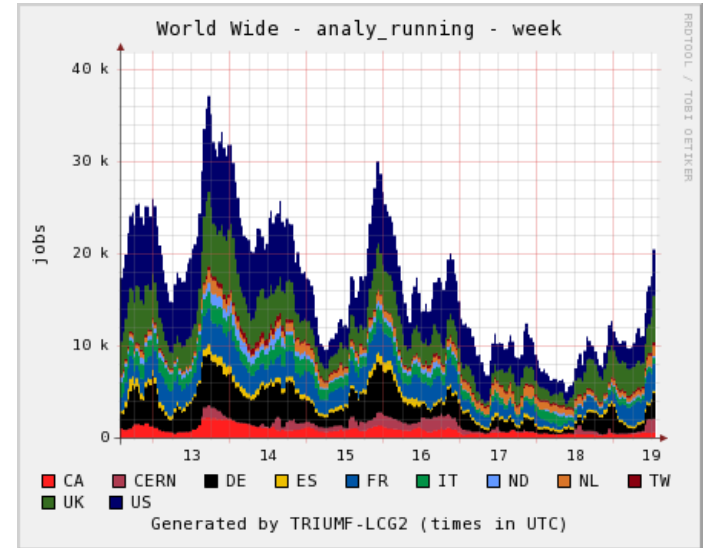
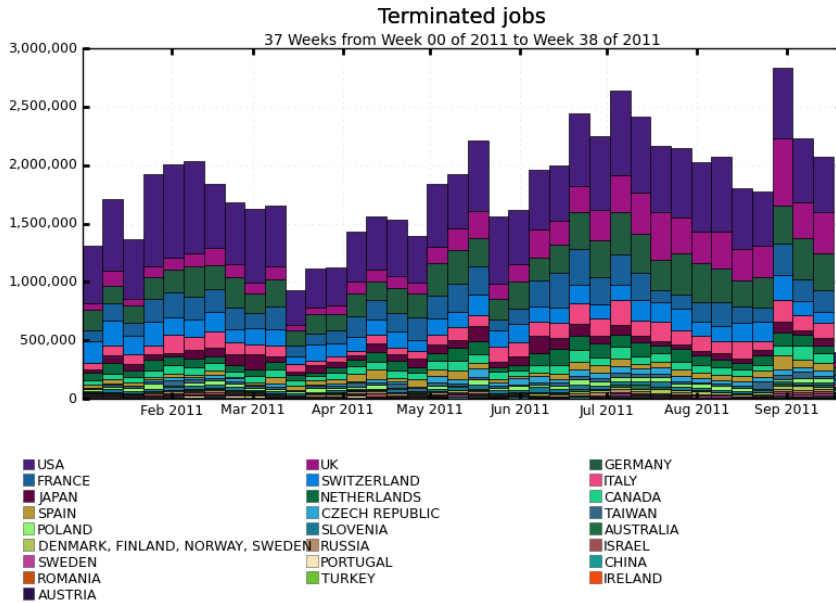
Dan van der Ster – CERN IT-ES
EGI TF 2011, Lyon, France



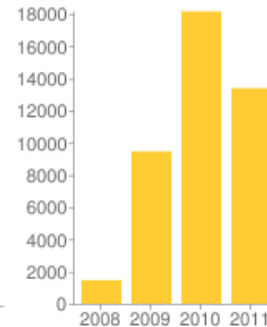
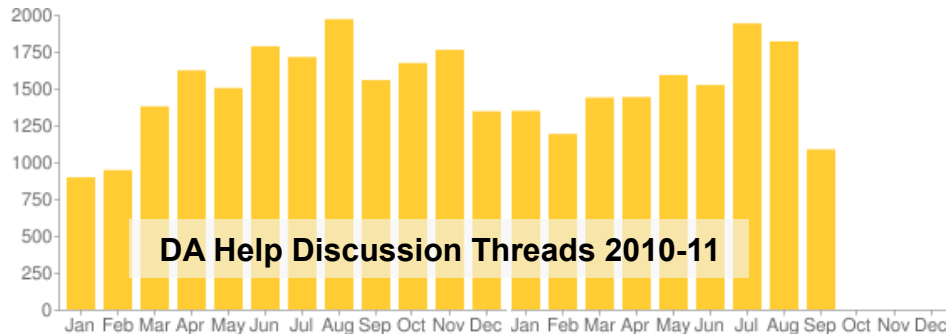
- Overview of Analysis Tools and Support
- Review of 2011
 - User activity
 - Development achievements
 - Operations experience
- Future Plans

- Support data analysis on the grid for the Heavy User Communities
 - Typical workflow is **distributed data mining**: users want to **process large amounts of input data** to produce summary statistics.
- Tools:
 - **Ganga**: end-user tool for batch and grid job submission and management
 - **DIANE**: lightweight pilot job framework for improved grid efficiency
 - **CRAB**: service for grid job management, with end-user client tailored to the CMS experiment
 - **HammerCloud**: grid site testing service for commissioning, validation, and ongoing functional testing

Approaching 3 million jobs per week with peaks of ~40,000 running concurrently

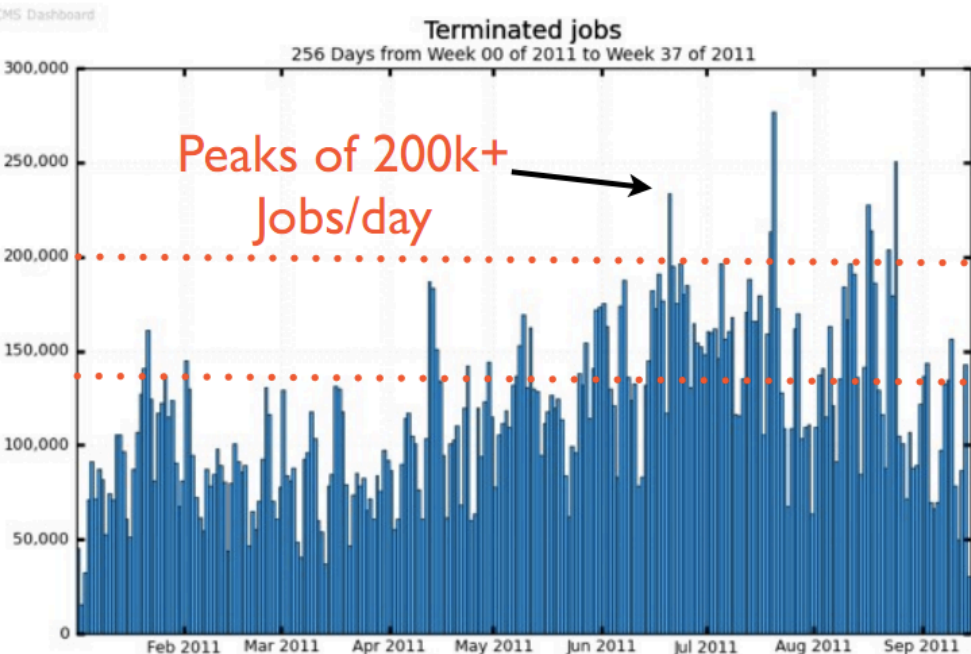
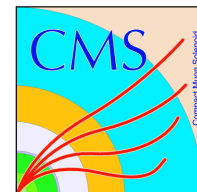


Maximum: 2,830,330 , Minimum: 0.00 , Average: 1,757,865 , Current: 395,105

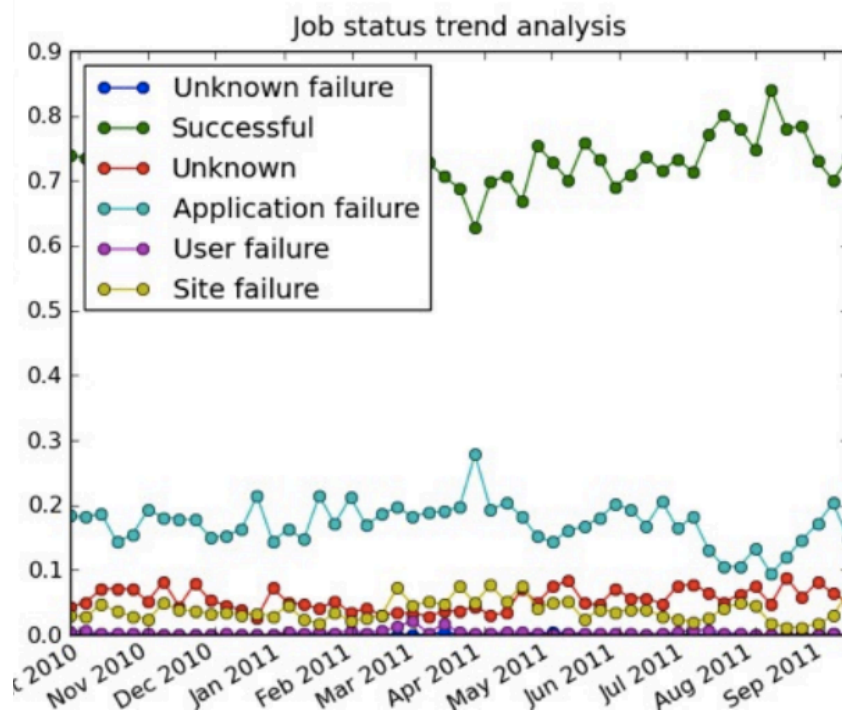
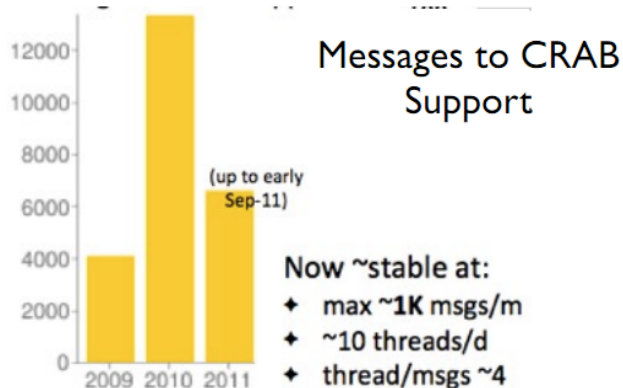


1452 distinct users
in the past 180 days





Job success rate ~ 75%



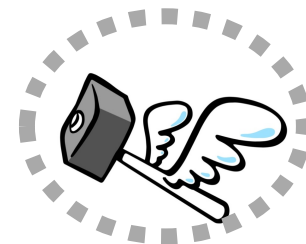
- Ganga is an job submission and management tool for local, batch and grid
 - <http://ganga.web.cern.ch/ganga/>



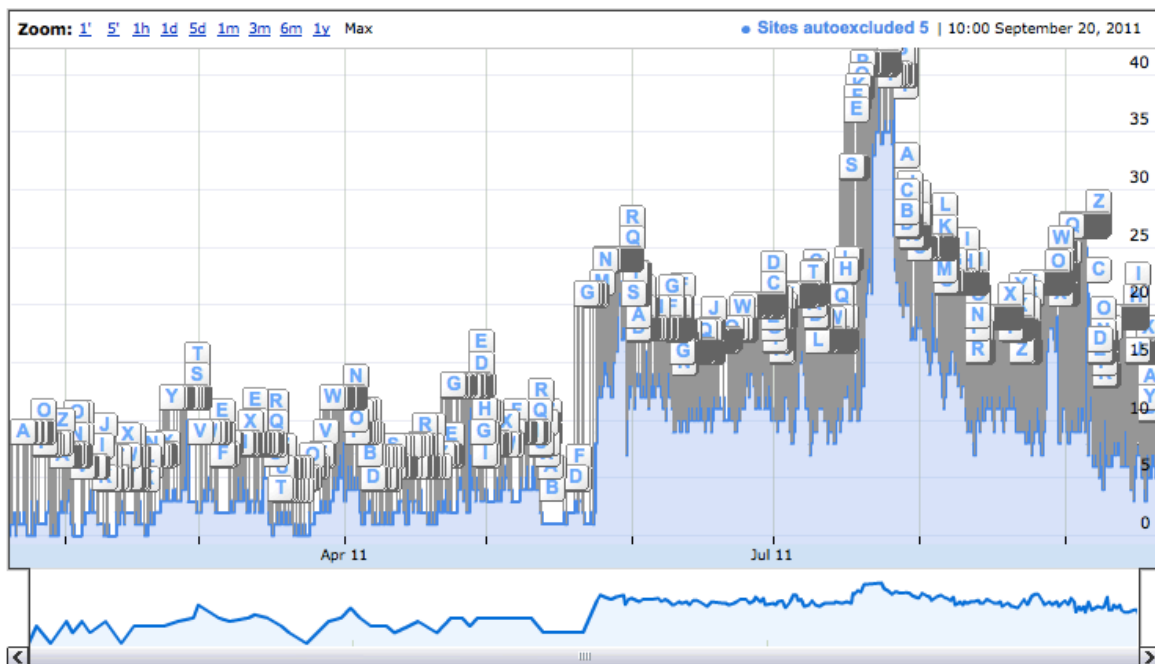
- Core development in 2011:
 - Releases 5.5.20->5.5.27 and 5.6.0->5.6.10 (mostly VO changes)
 - NEW: **Task management** (i.e. staged execution of a large job set) (as of 5.6.0)
 - NEW: Automatic **resubmission of failed subjobs**
- Ongoing maintenance of the VO plugins, for example:
 - ATLAS support for output data “**merging**” jobs
 - ATLAS support for **PanDA production** jobs
 - ATLAS **Job Execution Monitor** – live job peeking
 - LHCb support for **LHCbDirac v6**
 - LHCb added **Task management** features
- New VO: SuperB experiment application plugins in development
- Ganga 5.7 will bring a new job state: ***prepared***
 - `j=Job()` -> configure the `j.application` -> call `j.application.prepare()`
 - Persists the application state for use in the future, i.e. take an **app snapshot**
 - To re-execute the exact same code on different datasets
 - To inspect a historical application configuration weeks/months in the future
 - 5.7 beta is being tested now – targeted for fall 2011 release.



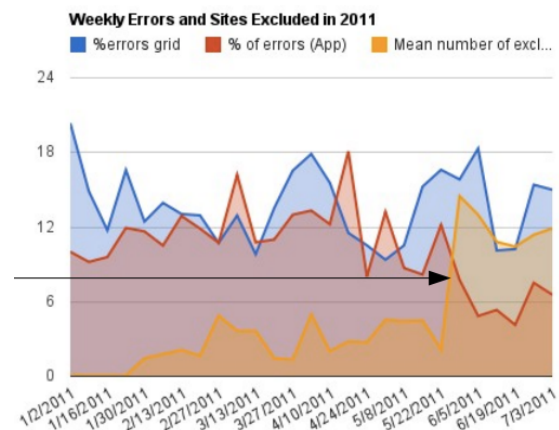
- **HammerCloud (HC)** is a grid site testing system serving two use-cases:
 - **Stress Testing**: on-demand large-scale stress tests to many sites simultaneously
 - **Functional Testing**: frequent short jobs to all sites to perform end-to-end validation
- Developed around **Ganga** using python backend and **Django** web frontend.
- **VO flexibility**: ATLAS, CMS, LHCb plugins in production



- **Auto-Exclusion** plugin used by ATLAS:
 - Use the stream of test jobs to **set PanDA queues offline and back online**
 - Relieves grid shifters and **increases grid reliability**



Strict exclusion in May decreased app errors by 50%

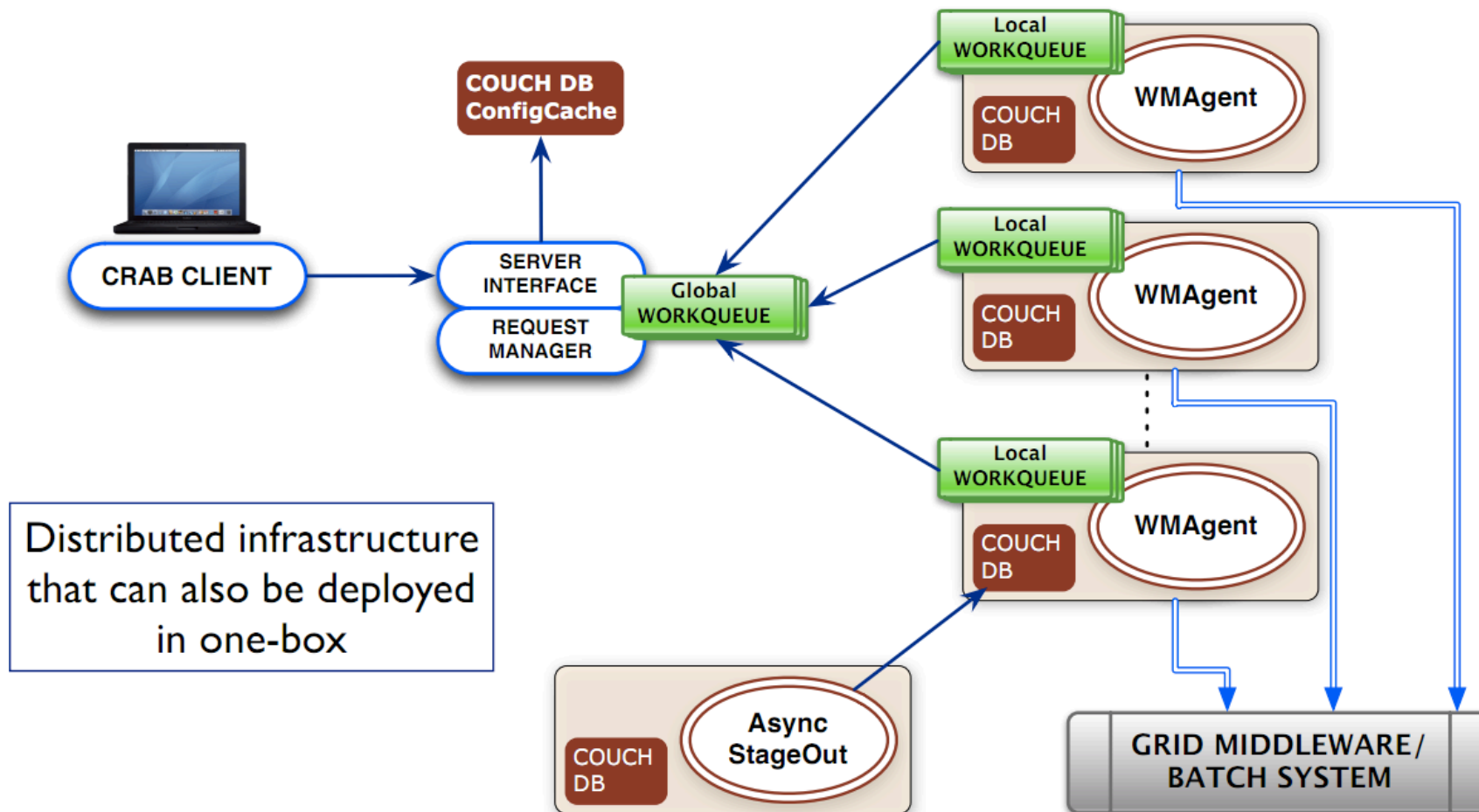


- **HammerCloud 4** introduced a VO plugin system:
 - VOs with Ganga support can add a **HammerCloud plugin** to submit jobs and retrieve their relevant statistics.
- Next steps for HammerCloud **5**:
 - **Oracle** Database support (to allow increased scalability over MySQL)
 - Highly-scalable **metric storage and display** system:
 - CMS wants to store *~200 performance metrics per job* and perform data mining on these metrics afterwards.
 - Investigating **NoSQL** databases for this.
 - General **performance improvements** to submit and monitor more and more test jobs.

- The **CMS Remote Analysis Builder** is the official CMS tool that allows to the end user to enable analysis job submissions through the grid
- By providing a simple CLI it allows to hide the underline complexities given from heterogeneous resources coming from different sources (middleware, specific experiment services, specific site implementations)
- **CRAB 2** is the currently in production system were a standalone system and client-server architecture are enabled
- CRAB is evolving into an architecture that will correspond to a new major release: **CRAB 3**

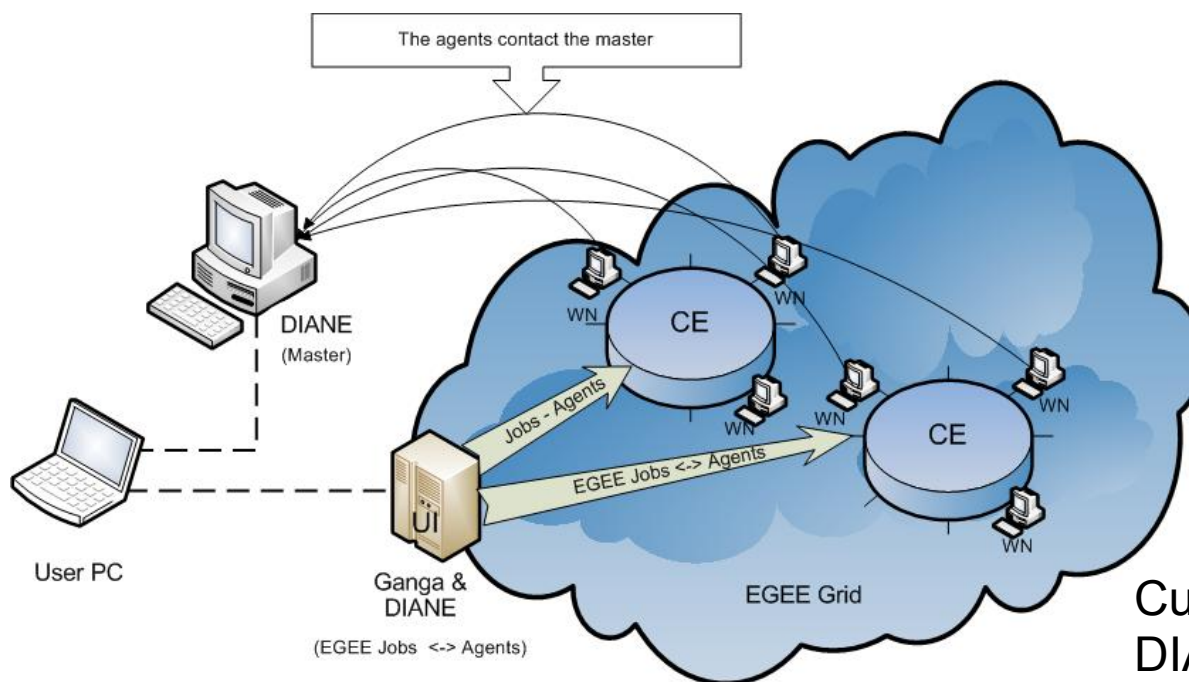
CRAB slides courtesy D.Spiga & M.Cinquilli (CERN IT-ES)

- The architecture which has been developed is getting more sophisticated c.w. CRAB2, being based on a client-server system implementing a **multi-tiers model**
 - Intermediate layer that receives all the client requests and stores them in a **global central queue**
 - **Local queues** that poll the global queue to retrieve new works; each local queue is associated to a server that processes the client request
- The key improvements in CRAB 3 are:
 - Thin stateless client and a stateful server: the workflow logic is entirely moved to the server level in order to encapsulate the end user in a protected environment
 - The server implements a RESTful based Web Service to enable the communication with the client
 - The monitoring is based on a **document oriented database to store job activities/ results** offering performance and simplicity advantages
 - The **AsynchronousStageOut** has been integrated to handle the user outputs
 - The deployment model has been improved with the aim to be quick and easy



- **CRAB 3** development activity is still ongoing but integration and commissioning tests started using a prototype version
 - Collecting **feedback** coming from **beta users** is the main objective
- The CRAB Team is working on the missing pieces to build the final deployment model
- The development effort is also focused on the physics domain aspects (eg: data publication, data merge, processed luminosity selections, ...)

- **DIANE** is a **pilot job framework** for users and small VO's to achieve more **efficient use** of grid resources
 - <http://it-proj-diane.web.cern.ch/it-proj-diane/>



Current version:
DIANE 2.4 (May 2011)

- Distributed analysis is very active! The HUCs have a few thousand users running millions of jobs per week.
- Tools:
 - **Ganga** used heavily in LHCb and ATLAS; SuperB in development.
 - **DIANE** is stable and ready for smaller VOs or individuals.
 - **HammerCloud** is a Ganga-app and with VO-plugins can be used for stress and functional testing.
 - **CRAB** is *the* CMS analysis system.
- Developments over the next year will continue to improve performance and reliability in grid analysis.
- Embracing new communities: contact the projects above if you are interested in using them in your VO.