

Dirac integration with a general purpose bookkeeping DB: a complete general suite for distributed resources exploitation

F Bianchi, *INFN*, M Chrzaszcz, *IFJ PAN*, V Ciaschini, *INFN*, M Corvo, *INFN*, C De Santis, *INFN*, D Del Prete, *INFN*, A Di Simone, *INFN*, G Donvito, *INFN*, A Fella, *INFN*, P Franchini, *INFN*, F Giacomini, *INFN*, A Gianelle, *INFN*, R Grzymkowski, *IFJ PAN*, S Longo, *INFN*, S Luitz, *SLAC*, E Luppi, *INFN*, M Manzali, *INFN*, M Rama, *INFN*, G Russo, *INFN*, S Pardi, *INFN*, L Perez, *INFN*, B Santeramo, *INFN*, R Stroili, *INFN*, L Tomassetti, *INFN*, and M Zdybal, *IFJ PAN*

Abstract—In the context of High Energy Physics computing field the R&D studies aimed to the definition of the data and workload models have been carried on and completed by the SuperB community beyond the experiment life itself. The work resulted of great interest for a generic mid- and small size VO to fulfill Grid exploiting requirements involving CPU-intensive tasks. We present the R&D line achievements in the design, developments and test of a distributed resource exploitation suite based on DIRAC. The main components of such a suite are the information system, the job wrapper and the new generation DIRAC framework. The DB schema and the SQL logic have been designed to be able to be adaptive with respect to the VO requirements in terms of physics application, job environment and bookkeeping parameters. A deep and flexible integration with DIRAC features has been obtained using SQLAlchemy

technology allowing mapping and interaction with the information system. A new DIRAC extension has been developed to include this functionality along with a new set of DIRAC portal interfaces aimed to the job, distributed resources, and metadata management. The results of the first functionality and efficiency tests will be reported.

I. INTRODUCTION

IN HEP as well in other fields, a typical issue is the necessity to manage and analyze huge amount of data or simulate large amount of events. Today several solutions are yet available for this purpose, but lot of them are developed for particular needs of a specific customer. During SuperB R&D activities, a solution useful for SuperB was deployed which can be easily adopted by a generic small and mid-size VO.

II. SUITE DESCRIPTION

In order to develop a simple, standard and long term solution, suite components (adopted or developed) should be flexible enough to be adapted in order to fulfill needs of a generic VO.

Suite components include DIRAC[1], the Information System and the Job Wrapper (see figure 1).

- DIRAC is a well known and widely adopted framework to manage Grid resources, job submission, workflow definition, user authentication, authorization and accounting.
- The Information System holds metadata related to simulations and information about dataset structures and data placement on Grid resources. Information System relies on a PostgreSQL[2] database. DIRAC interacts with Information System via SQLAlchemy[3].
- The Job Wrapper, executed in bundle with jobs, updates Information System about simulations status and data placement using a REST interface. Job Wrapper is a python script.

A. The Dirac extension

DIRAC can be extended to add specific functionalities required by the user. For example in SuperB[5] needs was an interface with a PostgreSQL BookKeeping database and a

F. Bianchi is with University of Torino, Turin, Italy and INFN - Sezione di Torino, Turin, Italy

M Chrzaszcz is with Physik-Institut, Universitat Zurich, Zurich, Switzerland and Henryk Niewodniczanski Institute of Nuclear Physics Polish Academy of Sciences, Krakow, Poland

V Ciaschini is with INFN - CNAF, Bologna, Italy

M Corvo is with INFN - Sezione di Padova, Padua, Italy

C De Santis is with INFN - Sezione di Roma Tor Vergata, Rome, Italy and Department of Physics, University of Rome Tor Vergata, Rome, Italy

D Del Prete is with INFN - Sezione di Napoli, Naples, Italy

A Di Simone is with INFN - Sezione di Roma Tor Vergata, Rome, Italy and Department of Physics, University of Rome Tor Vergata, Rome, Italy

G Donvito is with INFN - Sezione di Bari, Bari, Italy

A Fella is with INFN - Sezione di Pisa, Pisa, Italy and Department of Mathematics and Computer Science, University of Ferrara, Ferrara, Italy

P Franchini is with INFN - CNAF, Bologna, Italy

F Giacomini is with INFN - CNAF, Bologna, Italy

A Gianelle is with INFN - Sezione di Padova, Padua, Italy

R Grzymkowski is with Henryk Niewodniczanski Institute of Nuclear Physics Polish Academy of Sciences, Krakow, Poland

S Longo is with INFN - Sezione di Padova, Padua, Italy

S Luitz is with SLAC, USA

E Luppi is with Department of Physics, University of Ferrara, Ferrara, Italy and INFN - Sezione di Ferrara, Ferrara, Italy

M Manzali is with Department of Physics, University of Ferrara, Ferrara, Italy and INFN - Sezione di Ferrara, Ferrara, Italy

M Rama is with INFN - Sezione di Padova, Padua, Italy

G Russo is with INFN - Sezione di Napoli, Naples, Italy

S Pardi is with INFN - Sezione di Napoli, Naples, Italy

L Perez is with INFN - Sezione di Pisa, Pisa, Italy

B Santeramo is with INFN - Sezione di Bari, Bari, Italy and Department of Physics, University and Polytechnic of Bari, Bari, Italy

R Stroili is with INFN - Sezione di Padova, Padua, Italy

L Tomassetti is with Department of Mathematics and Computer Science, University of Ferrara, Ferrara, Italy and INFN - Sezione di Ferrara, Ferrara, Italy

M Zdybal is with Henryk Niewodniczanski Institute of Nuclear Physics Polish Academy of Sciences, Krakow, Poland

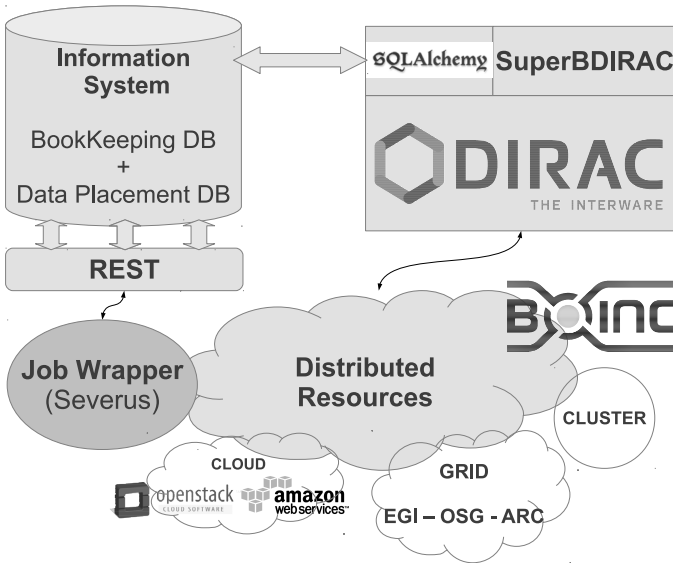


Fig. 1. Project bird's eye view.

webportal able to display monitoring data from this database where required. This particular DIRAC extension was named SuperBDIRAC. In SuperBDIRAC a new service, named SBKService, provides an SQLAlchemy layer able to connect DIRAC to a generic SQL DBMS. SBKService maps the bookkeeping database using Object Relational Mapping (see section IV). Webportal extension is designed to create and manage simulations, monitoring related jobs and sites. New functionalities integration into webportal permits to use only one interface to manage and monitor the entire stack of simulation related tasks.

III. BOOKKEEPING DB

MonteCarlo events production needs a method to identify data files and Storage Elements that holds them. A BookKeeping database, named SBK (SuperB BookKeeping), have been developed to manage metadata associated with data files. The same database stores information about simulations, executed jobs and output data, site availability in terms of installed and supported software. SBK is used also to schedule jobs submission in order to complete simulations. Database is modeled to fulfill general requirements of a typical simulation production.

Entities in SBK are Session, Production and Request (see figure 2).

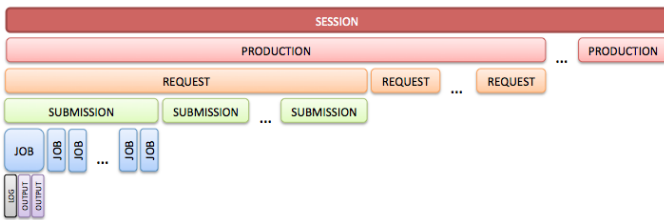


Fig. 2. Entities in BookKeeping database.

Session defines a simulation (eg. FastSim and FullSim): parameters and software for simulation are defined by VO

managers. Production is a Session subset that produce all the needed to simulate a particular scenario (eg. background in detector). Request is a Production subset. The required number of events to complete a Request is defined during its creation. Request completion is monitored via SBK, allowing job re-submission in order to complete it.

SBK design uses the relational model, and the current implementation relies on PostgreSQL (version 9.1) RDBMS which is SQL compliant and, exploiting its hstore datatype, allows to solve some major architectural issues concerning the dataset management of physical parameters. Hstore fields play a major role in the database architecture because storing sets of key/value pairs within a single PostgreSQL value is useful in various scenarios, such as rows with many attributes that are rarely examined, or semistructured data. Beside hstore, some other powerful PostgreSQL features have been exploited: its procedural language (PL/pgSQL) and schemas for a better management of user privileges. An extensive use of views and trigger procedures has been done too.

A. Normalization studies

During the development phase, the SBK database has been continuously analyzed in order to guarantee its normal form (NF) 1, 2 and 3 compliance. Four hierarchical levels (production, request, submission and ob+log+output+stat) have been identified for fastsim/fullsim to simplify table definitions and relations. In the production version the SBK database is NF1, NF2 and NF3 compliant with exception of hstore fields. A hstore is a string containing key- ζ value couples and for this reason, if splitted into each couple key- ζ value, it's not NF1 compliant. Since hstore fields permit to reduce database complexity and are rarely accessed (100 updates every 6 months), a trade-off has been accepted keeping hstore columns not normalized.

B. Stress tests

Extensive stress tests to check PostgreSQL and HTTP REST interface system robustness have been carried out by means of the Tsung tool (<http://tsung.erlang-projects.org/>). Tsung allows to create virtual machines for testing scalability and performance of IP based client/server applications in order to do load and stress testing of servers. It can be distributed on several client machines and is capable to simulate hundreds of thousands of virtual users concurrently. For the test phase the REST interface has been configured to establish permanent DB connections to save connection slots. During the stress tests up to $100 \text{ users} \cdot \text{s}^{-1}$ have been created. Each user has carried out a connection and 8 insert/update operations on a mock-up database which reproduced the real behavior of a production job. Stress test results were good, being the system capable to sustain 10000 DB transactions (1 transaction = 1 connection+8 insert/update) in 100s (900 operations $\cdot\text{sec}^{-1}$).

IV. BOOKKEEPING DB INTEGRATION

SBK integration in DIRAC is based on SQLAlchemy. SQLAlchemy uses Object Relational Mapper paradigm:

database entities are mapped as python objects. SQLAlchemy adoption simplify code writing, reading and documenting, provides an abstraction layer capable to manage in a transparent way a wide variety of database backends, giving freedom to change it without needs to re-write code. A new DIRAC service, named SBKService, integrates SQLAlchemy functionalities. SBKService interacts with other DIRAC components like any other service. SBKService maps SBK structure and expose methods to perform needed database operations.

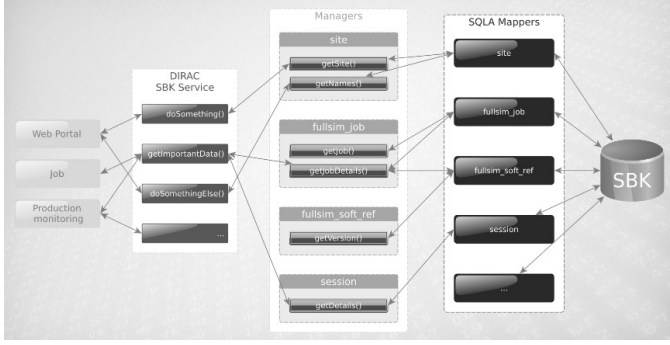


Fig. 3. SBKService schema.

V. JOB WRAPPER COMPONENT

Job wrapper named Severus takes care of main operations of a simulation job: copy software and input files to WN, copy output files in SE and register them in LFC, copy log in SE and bookkeeping DB, update job status in bookkeeping DB.

Access type (lcg or direct) from the worker node to the site SE is automatically recognized and implemented using lcg utils. A module for each session takes care of properly setting environment variables according to simulation (aka session). A configuration file customizes its behavior at execution time. Configuration file has several sections:

- **OPTIONS:** general parameters
- **SOFTWARE:** info about the executable
- **REST:** info about the REST interface to be used for communications with DB
- **SITE:** info about the submission site where the job will be running. A module for each site is loaded
- **TARGETSITE:** info about site where replicas of output files must be written
- **INPUT:** info about the job input files
- **OUTPUT:** info for stage-out phase
- **EXPORT VARS:** list of environment vars to be exported on the WN
- **SESSION_NAME:** key-value pairs for simulation parameters, for the specified session

VI. SIMULATION PRODUCTION USE CASE: THE SUPERB EXPERIENCE

Simulation Production is designed to manage huge Monte-Carlo productions. A webportal, named WebUI[7], is available to manage the entire stack of operations related to this use case: user management, definition of new "Sessions", "Productions" and "Request", job submission and monitoring, sites

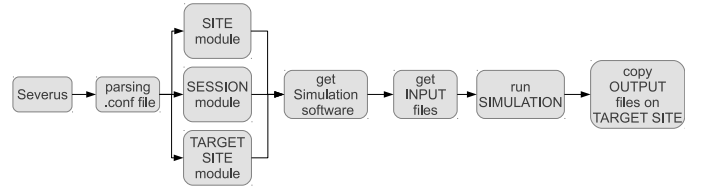


Fig. 4. Severus workflow.

management for productions. Production manager users can add and delete sites, add and delete CEs and SEs, set a site as enabled/disabled and supported/unsupported for a given session, manage "Sessions", "Productions" and "Request". Shifter users can submit and monitor jobs for a given Request. Job submission in WebUI is performed via Ganga, while submission configurations files are generated by WebUI itself taking data from SBK. Ganga submit jobs to grid via WMS using standard gLite commands. Job execution on WNs is driven by Severus (see section V), while job status updates in SBK is performed via REST interface. Stagein and Stageout as well as output files registration in LFC is performed even by Severus. SuperB DIRAC goal was the porting of the WebUI functionalities in DIRAC to manage jobs and their submission directly from DIRAC and exploit all the functionalities available in DIRAC: ie. Grid computing resources can be used via WMS or direct submission to CREAM CEs, computer clusters can be accessed via ssh connections, Cloud resources are available via VMDIRAC module, even desktop computers can be used via Boinc. User authentication via X509 certificates and authorization VOMS-role based are builtin in DIRAC. SuperB DIRAC enhance DIRAC integrating the generic bookkeeping database SBK via SQLAlchemy. DIRAC webportal is extended in order to provide an interface for Production manager actions as well as shifter tasks.

Job Monitor							
Job Status	Submitter	Prod Series	Run num	Submission time	Last update	Status reason	
running	bsartermo	2010_July	33001413	2013-09-14 10:12:22	2013-09-14 10:29:46	FastSim executed in 665.1 seconds. Copying output...	
running	bsartermo	2010_July	33001355	2013-09-14 10:12:22	2013-09-14 10:29:46	FastSim executed in 676.4 seconds. Copying output...	
done	bsartermo	2010_July	33001367	2013-09-14 10:12:22	2013-09-14 10:29:46	Output files copied in 7.0 seconds. Job completed.	
running	bsartermo	2010_July	33001365	2013-09-14 10:12:22	2013-09-14 10:29:46	FastSim executed in 670.5 seconds. Copying output...	
done	bsartermo	2010_July	33001352	2013-09-14 10:12:22	2013-09-14 10:29:44	Output files copied in 7.0 seconds. Job completed.	
running	bsartermo	2010_July	33001450	2013-09-14 10:12:22	2013-09-14 10:29:44	FastSim executed in 658.9 seconds. Copying output...	
running	bsartermo	2010_July	33001368	2013-09-14 10:12:22	2013-09-14 10:29:44	FastSim executed in 724.2 seconds. Copying output...	
running	bsartermo	2010_July	33001415	2013-09-14 10:12:22	2013-09-14 10:29:43	FastSim executed in 672.8 seconds. Copying output...	
running	bsartermo	2010_July	33001365	2013-09-14 10:12:22	2013-09-14 10:29:43	FastSim executed in 654.1 seconds. Copying output...	
done	bsartermo	2010_July	33001334	2013-09-14 10:12:22	2013-09-14 10:29:43	Output files copied in 7.1 seconds. Job completed.	
done	bsartermo	2010_July	33001348	2013-09-14 10:12:22	2013-09-14 10:29:43	Output files copied in 7.1 seconds. Job completed.	
done	bsartermo	2010_July	33001325	2013-09-14 10:12:22	2013-09-14 10:29:43	Output files copied in 7.1 seconds. Job completed.	
running	bsartermo	2010_July	33001372	2013-09-14 10:12:22	2013-09-14 10:29:42	FastSim executed in 719.5 seconds. Copying output...	
running	bsartermo	2010_July	33001322	2013-09-14 10:12:22	2013-09-14 10:29:41	FastSim executed in 674.7 seconds. Copying output...	
done	bsartermo	2010_July	33001353	2013-09-14 10:12:22	2013-09-14 10:29:41	Output files copied in 6.9 seconds. Job completed.	
done	bsartermo	2010_July	33001391	2013-09-14 10:12:22	2013-09-14 10:29:40	Output files copied in 6.9 seconds. Job completed.	
running	bsartermo	2010_July	33001407	2013-09-14 10:12:22	2013-09-14 10:29:40	FastSim executed in 670.9 seconds. Copying output...	
done	bsartermo	2010_July	33001376	2013-09-14 10:12:22	2013-09-14 10:29:36	Output files copied in 7.0 seconds. Job completed.	
done	bsartermo	2010_July	33001423	2013-09-14 10:12:22	2013-09-14 10:29:34	Output files copied in 7.0 seconds. Job completed.	
done	bsartermo	2010_July	33001342	2013-09-14 10:12:22	2013-09-14 10:29:31	Output files copied in 7.0 seconds. Job completed.	
done	bsartermo	2010_July	33001328	2013-09-14 10:12:22	2013-09-14 10:29:31	Output files copied in 6.8 seconds. Job completed.	
done	bsartermo	2010_July	33001324	2013-09-14 10:12:22	2013-09-14 10:29:28	Output files copied in 6.8 seconds. Job completed.	
done	bsartermo	2010_July	33001327	2013-09-14 10:12:22	2013-09-14 10:29:28	Output files copied in 6.8 seconds. Job completed.	
done	bsartermo	2010_July	33001341	2013-09-14 10:12:22	2013-09-14 10:29:24	Output files copied in 7.0 seconds. Job completed.	
done	bsartermo	2010_July	33001359	2013-09-14 10:12:22	2013-09-14 10:29:20	Output files copied in 6.8 seconds. Job completed.	
done	bsartermo	2010_July	33001403	2013-09-14 10:12:22	2013-09-14 10:29:18	Output files copied in 7.1 seconds. Job completed.	

Fig. 5. Monitoring job data from SBK in DIRAC.

VII. FUNCTIONALITY TEST

A functionality test was performed to demonstrate the capability of SuperB DIRAC to integrate the monitoring of a bookkeeping database in DIRAC webportal. Functionality Test objective is to demonstrate that SuperB DIRAC is capable of substituting WebUI as monitoring tool for job submission and showing data

from a bookkeeping database. Test "Q-factor" is the exact correspondence of information as stored in SBK and displayed in SuperBDirac, like in WebUI monitoring page. Correct execution of simulation jobs is not important as long as the bookkeeping information is properly managed.

At present time, not all WebUI functionalities are implemented in SuperBDirac, in particular "Submission" job for a given "Request". The following procedure was used to perform this test. FastSim job submission is created via WebUI interface. Once submission is created, a set of scripts and configuration files are created. In particular the php submission script, generated by WebUI, is made of an array with all relevant parameters and UI commands needed to submit jobs in grid via standard glite commands. Since WebUI is still linked with SBK, its portal could be used as well as monitoring portal, useful for a check-cross between info displayed in SBK, WebUI and SuperBDirac.

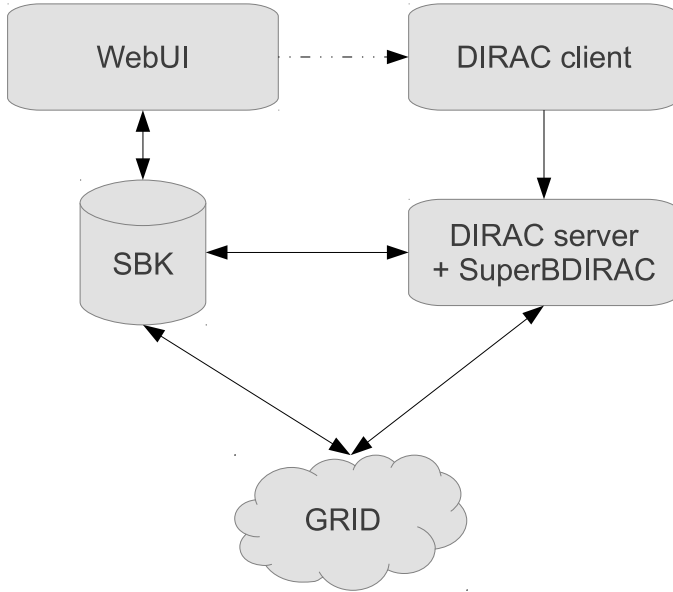


Fig. 6. Testbed schema.

The php submission script is parsed by a python script (mc_production.py), in particular the params array, in order to retrieve all needed parameters to properly submit jobs via DIRAC: production series, session name, min and max runnumber, configuration files location, physical parameters, events to simulate. Once taken all these parameters, mc_production.py uses DIRAC API to prepare and submit jobs via DIRAC client.

DIRAC server receive jobs from DIRAC client, than starts the normal workflow for job management in DIRAC: scheduling, pilot submission, payload retrieval and execution, stageout. Since DIRAC server is equipped with SuperBDirac, even bookkeeping monitoring is performed by this component. Every job simulated 3000 events: this value is set to have an execution time quite longer than 10 minutes. Physical parameters are the same of other official productions. 3 main bunch submission of 400 jobs were performed at INFN-T1 to obtain a total of 1200

simulation jobs. Status in SBK were "prepared", "running" and "done". In addition, 2 bunch submission of 10 jobs were performed, again at INFN-T1, to force some failure messages in SBK and catch it even in SuperBDirac monitoring. First failure sample was obtained setting a not-existing site as destination for stageout: error was detected during preliminary check, so status in SBK were "prepared" and "failed". In second failure sample, jobs were submitted using a proxy without Role=ProductionManager, so error occurred in stageout phase: status in SBK were "prepared", "running" and "failed". In summary, 1220 jobs were submitted for test.

Test	Jobs	prepared		running		done		failed		success rate
		A	B	A	B	A	B	A	B	
good-1	400	400	400	400	400	400	400	0	0	100%
good-2	400	400	400	400	400	400	400	0	0	100%
good-3	400	400	400	400	400	400	400	0	0	100%
fail-1	10	10	10	0	0	0	0	10	10	100%
fail-2	10	10	10	10	10	0	0	10	10	100%

TABLE I
TEST RESULTS - A) EXPECTED B) TEST RESULT

BookKeeping database was properly updated for every job in every test. All status change were promptly displayed as well in SuperBDirac as in WebUI, without any appreciable delay between two portals. SQLAlchemy didn't introduced any appreciable delay or information loss, at least in this functionality test. Table I reports, for every test, how many status were saved in SBK and displayed in WebUI and SuperBDirac. Success rate was established as ratio between status saved in SBK and status displayed in SuperBDirac: its value was 100% in all submissions. SuperBDirac could be considered good enough to integrate in DIRAC the capability of monitoring jobs metadata from a bookkeeping database.

VIII. CONCLUSIONS

DIRAC is a mature and stable framework to manage all grid-related tasks, easily adoptable by small as well large VOs. The Information System is designed to be adapted for a generic huge simulation production. The Job Wrapper acts as a bridge between simulation jobs and Information System. We propose SuperBDirac as a DIRAC extension capable to satisfy the needs of small and mid size VOs in terms of distributed resource exploitation.

REFERENCES

- [1] <http://diracgrid.org>
- [2] <http://www.postgresql.org/>
- [3] <http://www.sqlalchemy.org/>
- [4] <http://boinc.berkeley.edu/>
- [5] SuperB Technical Design Report, <http://arxiv.org/abs/1306.5655>
- [6] Fielding R T 2000 *Architectural Styles and The Design of Network-based Software Architectures*, PhD Thesis, University of California Irvine
- [7] A.Fella, E.Luppi, L.Tomassetti *A General Purpose Suite for Job Management, Bookkeeping and Grid Submission*. International Journal of Grid Computing & Applications (IJGCA) Vol.2, No.2, June 2011. DOI: 10.5121/ijgca.2011.2202.