# A Multidimensional Approach to the Analysis of Grid Monitoring Data*

S. Andreozzi, G. L. Rubini, INFN-CNAF[†], Bologna, Italy
Sergio Fantinel, Lab. Naz. di Legnaro[‡], Legnaro, Italy
N. De Bortoli, G. Tortone, INFN-NA[§], Napoli, Italy

## Abstract

Analyzing Grid monitoring data requires the capability of dealing with multidimensional concepts intrinsic to Grid systems. Typically, monitoring tools store data into databases that rely on the relational model. This model is not ideal for the efficient support of data analysis that requires quick navigation through the different data dimensions. In this paper, we discuss the application of On-Line Analytical Processing (OLAP), an approach to the fast analysis of shared multidimensional information. Our proposal is exemplified over monitoring data collected from a production-quality Grid system by means of the GridICE monitoring tool.

## INTRODUCTION

Grid computing is concerned with the virtualization, integration and management of services and resources in a distributed, heterogeneous environment that supports collections of users and resources across traditional administrative and organizational domains [1].

One aspect of particular importance is Grid monitoring, that is the activity of measuring significant Grid resource-related parameters to analyze usage, behavior and performance of a Grid system, and to detect and notify fault situations, contract violations and user defined events.

The monitoring data of Grid systems are intrinsically multidimensional, therefore they require to be aggregated along different dimensions. The meaningful dimensions that we have identified are the physical dimension referring to geographical location of resources, the Virtual Organization (VO) dimension, the time dimension and the resource identifier dimension.

We propose the application of On-Line Analytical Processing (OLAP), a meaningful approach to the analysis of shared multidimensional information. It is exemplified over monitoring data collected from a production-quality Grid system by means of the GridICE [4, 5] monitoring tool.

## GRID MONITORING WITH GRIDICE

GridICE is a monitoring tool designed specifically for Grid systems. It integrates with the Grid Information Service

(GIS) and extends the GLUE schema [2] to deal with detailed host and job monitoring metrics. Historical data are kept to enable retrospective analysis useful in a number of different situations, e.g., usage metering or Service Level Agreement (SLA) monitoring of resource providers. Three are the identified viewpoints that require a different view of monitoring data: the VO, the site and the operations domain. GridICE offers a web interface that enables the browsing of the data collection according to these different viewpoints (see Figure 1). GridICE is currently integrated in the LCG [6] middleware. In Figure 2, we present an example of deployment scenario. Two sites are involved: INFN-LNL in Legnaro and INFN-CNAF in Bologna. The central database of the GridICE server that maintains the raw monitoring data is optimized to reduce the redundancy of data.

## ENABLING MULTIDIMENSIONAL ANALYSIS

In this section, we introduce the main concepts of multidimensional analysis of data as regards OLAP. We present the typical steps to transform data relying on a relational schema to data in the form of structures called 'OLAP cubes'. These are created by a reorganization of data contained inside a relational database, thus transforming operational data into dimensional data. Subsequently, we propose a design suitable for data stored in the GridICE database.

### Background on OLTP and OLAP

OLTP is an acronym for On-Line Transaction Processing. This is a class of programs that facilitates and manages transaction-oriented applications, typically for data entry and retrieval in contexts such industries, banking, airlines, mail order, supermarkets, and manufacturing. Generally, the OLTP concept is associated to the concept of Relational Data Base System (RDBMS). As regards user and system orientation, it is 'customer oriented' and it is used for query and transaction processing. It manages current data that are typically too detailed to be easily used for decision making. The database design relies on the entity-relationship (ER) data model. In this context, GridICE is a typical example of an OLTP application, as it continuously stores monitoring data in a PostgreSQL RDBMS.

OLAP is an acronym for On-Line Analytical Processing. It is an approach to the quick provision of answers to complex database queries. It is used in business reporting

† e-mail address: sergio.andreozzi,gianluca.rubini@cnaf.infn.it
‡ e-mail address: sergio.fantinel@lnl.infn.it
§ e-mail address: natascia.debortoli,gennaro.tortone@na.infn.it

| | | | | | Computing Resources | | | | | | | Storage Resources | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Site | Q# | Slot# | SlotFree | SlotLoad | RunJob | WaitJob | JobLoad | Power | WN# | CPU# | CPULoad | Available | Total | % |
| ba.infn.it | 3 | 186 | 105 | 8% | 25 | 0 | 19% | 205K | 21 | 42 | 9% | 376.0 Gb | 1.8 Tb | 80% |
| bo.infn.it | 6 | 90 | 90 | 0% | 0 | 0 | 0% | 67K | 15 | 30 | 0% | 1.7 Tb | 1.8 Tb | 5% |
| bo.ingv.it | 2 | 0 | 0 | - | 0 | 6 | 0% | - | - | - | - | 23.9 Gb | 34.3 Gb | 30% |
| ca.infn.it | 3 | 42 | 9 | 78% | 11 | 0 | 55% | - | - | - | - | 139.2 Gb | 142.1 Gb | 2% |
| cern.ch | 4 | 986 | 95 | 90% | 427 | 55 | 81% | - | - | - | - | 610.0 Gb | 688.2 Gb | 1% |
| cnaf.infn.it | 10 | 36 | 36 | 0% | 0 | 0 | 0% | 57K | 6 | 12 | 0% | 178.7 Gb | 234.9 Gb | 8% |
| cr.cnaf.infn.it | 7 | 1192 | 1178 | 1% | 0 | 0 | 0% | 3M | 462 | 924 | 29% | 1.5 Gb | 999.7 Gb | 100% |
| ct.infn.it | 3 | 174 | 57 | 67% | 0 | 0 | 0% | 267K | 28 | 46 | 29% | 1.7 Tb | 2.0 Tb | 13% |
| fe.infn.it | 3 | 42 | 12 | 71% | 10 | 0 | 5% | 53K | 7 | 14 | 68% | 20.5 Gb | 25.6 Gb | 20% |
| lnf.infn.it | 3 | 18 | 0 | 100% | 6 | 33 | 100% | 16K | 3 | 3 | 0% | 1.0 Tb | 1.1 Tb | 11% |
| lnl.infn.it | 5 | 730 | 40 | 94% | 109 | 0 | 73% | 466K | 75 | 150 | 1% | 178.3 Gb | 1.3 Tb | 87% |
| mi.infn.it | 3 | 177 | 0 | 100% | 59 | 8 | 100% | 269K | 31 | 62 | 54% | 1.2 Tb | 2.2 Tb | 14% |
| na.infn.it | 9 | 108 | 36 | 66% | 24 | 10 | 0% | 182K | 18 | 36 | 61% | 581.2 Gb | 1.7 Tb | 17% |
| oat.ts.astro.it | 3 | 6 | 6 | 0% | 0 | 0 | 0% | 1K | 1 | 2 | 0% | 29.5 Gb | 34.3 Gb | 14% |
| pd.infn.it | 3 | 330 | 225 | 31% | 35 | 0 | 13% | 510K | 56 | 112 | 31% | 259.8 Gb | 1.2 Tb | 28% |
| pg.infn.it | 3 | 12 | 12 | 0% | 0 | 0 | 0% | 3K | 2 | 4 | 0% | 214.9 Gb | 216.5 Gb | 1% |
| pi.infn.it | 3 | 51 | 48 | 5% | 0 | 0 | 0% | 52K | 8 | 16 | 0% | 4.3 Gb | 13.4 Gb | 68% |
| pv.infn.it | 3 | 18 | 18 | 0% | 0 | 0 | 0% | 7K | 3 | 6 | 0% | - | - | - |
| roma1.infn.it | 6 | 216 | 30 | 86% | 62 | 15 | 67% | 238K | 31 | 62 | 61% | 46.0 Gb | 83.7 Gb | 4% |
| roma2.infn.it | 3 | 0 | 0 | - | 0 | 0 | 0% | 19K | 4 | 4 | 25% | 9.2 Gb | 9.7 Gb | 5% |
| to.infn.it | 7 | 96 | 0 | 100% | 24 | 0 | 100% | 114K | 12 | 24 | 53% | 421.2 Gb | 1.9 Tb | 79% |
| ts.infn.it | 3 | 6 | 6 | 0% | 0 | 0 | 0% | 4K | 1 | 2 | 0% | 29.1 Gb | 35.1 Gb | 17% |
| TOTAL | 95 | 4516 | 2003 | 4% | 792 | 127 | 32% | 5M | 784 | 1569 | 26% | 8.7 Tb | 17.7 Tb | 29% |

Generated: Tue, 31 Aug 2004 11:52:44 +0200                    GridICE Homepage

Figure 1: INFN Grid: a GOC view

for sales, marketing, management reporting, data mining and similar areas. Such systems can organize and present data in various formats in order to accommodate the diverse needs of different users. As regards user and system orientation, it is 'market oriented' and used for data analysis. It manages large amounts of historical data, provides facilities for summarization and aggregation. The database design relies on star or snowflake schema model.

*From Relational Data to Analytical Data*

OLAP is based on a multidimensional data model where data are arranged in the form of cubes. It allows data to be modeled and viewed in multiple dimensions. The OLAP multidimensional data model is organized around one or more central themes called fact tables, where facts are numerical measures (see Figure 4).

Generally, OLAP dimensions are associated to perspectives that are interesting for organizations. Each dimension has a table associated to it (the dimension table). The meaningful dimensions identified for Grid systems are: time, Virtual Organization (VO), geography, Computing Element (CE) resource identity and Storage Element (SE) resource identity.

In general, data coming from OLTP systems are not directly usable for building an OLAP system. They need to be transformed from a relational form to an analytical form. The transformation process consists of four steps: (1) merging, a process needed when it is necessary to manage data related to multiple OLTP systems; (2) scrubbing, a process needed when data sources are multiple OLTP systems that have inconsistencies among them (e.g., different languages or different names for the same attribute); the goal of the scrubbing process is to delete these inconsistencies; (3) aggregation, needed because OLTP systems record all transaction details, while OLAP queries typically need summary data or data aggregated in some fashion; (4) organization of data in cubes, that refers to the construction of the fact table and the related dimensions.

In our case, merging and scrubbing processes are not necessaries because we are managing data coming from a unique OLTP system: the PostgreSQL database of GridICE. Conversely, the aggregation process is required because GridICE stores different attributes in many historical tables (see Figure 3). The logical schema of the database consists of many tables dedicated to the status of resources observed with configurable interval, e.g., the CE table that contains the current status of a set of CEs split into different historical tables. OLAP dimension tables can have concept hierarchies. A concept hierarchy defines a sequence of mapping from a set of low-level concepts to higher-level ones. Time and geography dimensions are order relations and form concept hierarchies: hour-day-week-month-year and site-country. In Figure 4, we propose
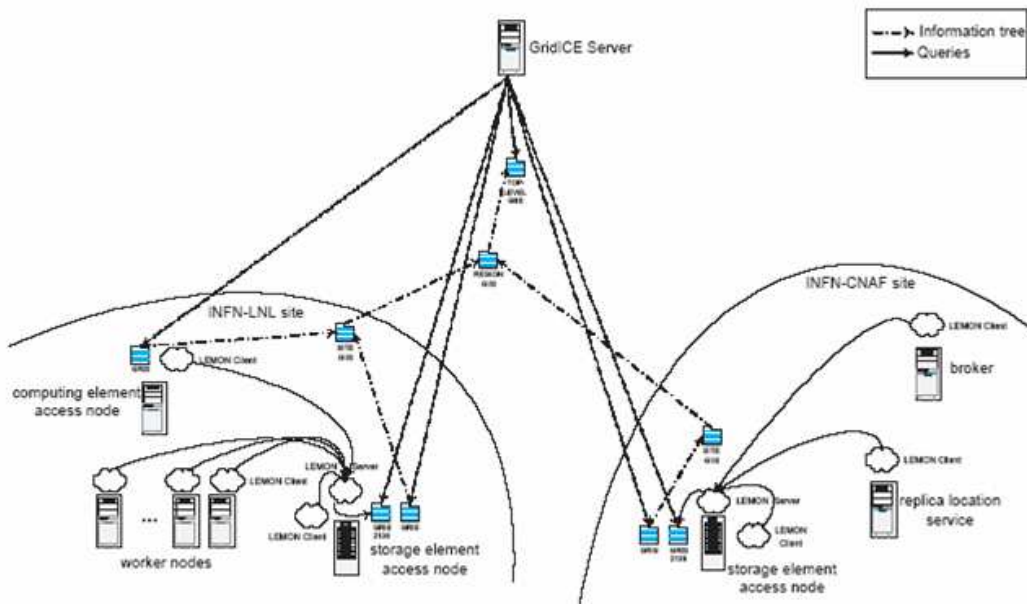
Figure 2: Deployment scenario

the GridICE fact table design.

## CHOOSING AN OLAP TOOL AND BUILDING THE OLAP DATABASE

Unfortunately, PostgreSQL does not provide a meaningful set of OLAP functions for our purposes. We investigated a number of open source software tools, but a worthy solution was not found. Therefore, we decided to resort to a commercial product, Oracle 9i, which provides an OLAP module [7].

The porting process required the following data transformations: the PostgreSQL integer datatype must be converted in Oracle numeric format; the PostgreSQL boolean datatype is not implemented in Oracle 9i, so it requires a booleanto-integer conversion; the PostgreSQL character datatype must be converted in an Oracle compatible format, we have chosen the varchar2 format.

Another aspect to consider is the VO dimension, that is a multi-value attribute for the CE entity (a CE can be authorized to many VOs). In order to avoid redundancy,

this would require the capability of defining a multi-value attribute in the CE fact table to which VO dimension is mapped. There is no native support for this capability in the selected tool, therefore as a first step we have decided to duplicate the tuples in the fact table, each of them referring to a single VO. Another source of redundancy is the need for reassembling the chunks of the historical tables of a certain entity into a single fact table.

## CONCLUSION

OLAP is an appealing solution for the analysis of multidimensional data in Grid systems. We have proposed a design of an OLAP hypercube that deals with structure of the data collected from GridICE. The main problems arise from the different involved database systems (PostgreSQL for OLTP and Oracle for OLAP), for the nature of monitoring data as regards the time dimension (static vs. dynamic attribute values) and for the dimensions involved in
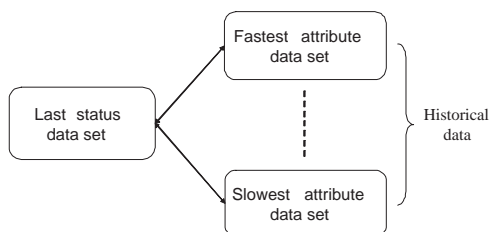


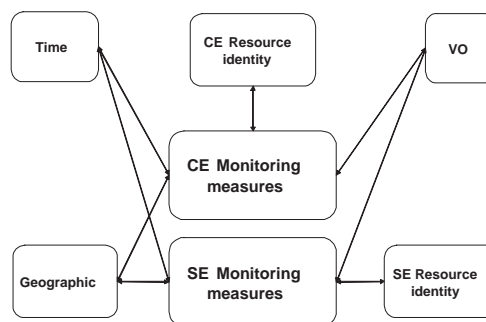Figure 3: Relational Schema - last status and historical tables splitting



Figure 4: Fact constellation schema for GridICE multidimensional data model

the Grid (i.e., VO dimension). Future work will be targeted at refining the logical OLAP schema in order to better deal with the identified problems. Further, an important goal is the integration of the OLAP functionalities in the GridICE monitoring tool.

# REFERENCES

[1] GGF Open Grid Services Architecture Working Group. https://forge.gridforum.org/projects/ogsa-wg

[2] GLUE Schema - Resources. http://www.cnaf.infn.it/∼sergio/glue

[3] S. Andreozzi. GLUE Schema Implementation for the LDAP Model. INFN Technical Report INFN/TC-04/16. 30 Sep 2004. http://www.lnf.infn.it/sis/preprint/pdf/INFN-TC-04-16.pdf

[4] S. Andreozzi, N. De Bortoni, S. Fantinel, A. Ghiselli, G.L. Rubini, G. Tortone, M.C. Vistoli. GridICE: a Monitoring Service for Grid Systems. Preprint copy, Aug 2004. To appear in Future Generation Computer Systems Journal, Elsevier.

[5] GridICE web site. http://grid.infn.it/gridice.

[6] LCG web site. http://lcg.web.cern.ch/LCG/

[7] Oracle OLAP resources. http://www.oracle.com/technology/products/bi/olap/olap.html