

INFN-CNAF Monitor and Control System

**Stefano Antonelli, Donato De Girolamo, Luca dell'Agnello,
Daniele Gregori, Guido Guizzunti, Pier Paolo Ricci, Felice Rosso,
Vladimir Sapunenko, Riccardo Veraldi, Paolo Veronesi,
Cristina Vistoli, Giulia Vita Finzi, Stefano Zani**

INFN-CNAF, viale Berti-Pichat 6/2, 40127 Bologna

E-mail: daniele.gregori@bo.infn.it

Abstract.

CNAF is the national center of National Institute of Nuclear Physics (INFN) for R&D in the field of Information Technologies applied to High Energy Physics (HEP) experiments. It is involved in the management and development of the most important information and data handling services in behalf of the INFN. In 2005, the Italian Tier-1 for Large Hadron Collider (LHC) experiments has been inaugurated at INFN-CNAF.

Due to the huge complexity of Tier-1 center, the use of control systems is fundamental for management and operation of the center. At INFN-CNAF, several solutions have been adopted, from commercial to open source products up to entirely home-made systems.

Adopted open source solutions have been strongly adapted to specific needs; a wide set of customized sensors has been developed for various divisions like Network, Storage, Farming, Grid operation and National Services.

Finally, a dashboard has been developed, to which described control systems send critical alarms (sent via sms to an operator as well). The dashboard can be exploited to get an historical view of the Tier-1 and national services' state and to allow a quick web control.

In this article, the whole system, adopted customizations in monitoring and control as well as their integrations with the dashboard will be described.

1. Introduction

The Italian Tier1 center for LHC computing is located in Bologna at INFN-CNAF [1], the main INFN computing facility. The Tier1 provides storage and computing resources to LHC and other HEP experiments. Its complexity requires control systems for management and operation of the center. Historically since the start of the Tier-1 project many monitoring systems have been adopted. In fact each of the Tier-1 division had their own specific requirements and the monitoring solution was chosen by the group system administrators so as to fulfill the division needs. So this lead to a situation where different heterogeneous system coexists in our site, including commercial, open source software and home made ones. Besides that all the different monitoring systems are integrated in the Dashboard (Fig. 1) that shows the overall status of the whole Tier-1 center.

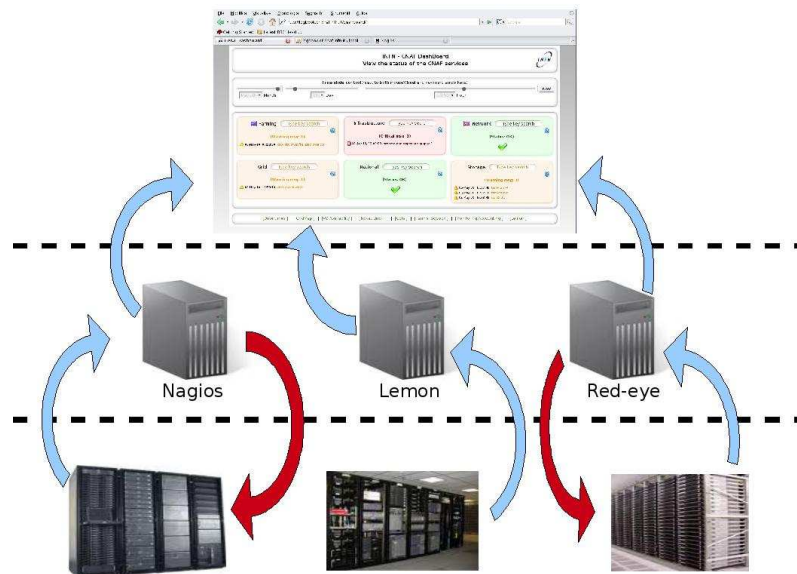


Figure 1. An overall view of the monitor system at INFN-CNAF. The servers of each INFN-CNAF division are the low level in this schema and they form the clients for the monitor system. These clients run a set of scripts that move information (blue arrow) towards the middle level that is composed of control server like Nagios, Lemon or Red-eye. The Middle level servers collect information, evaluate the results and eventually send alarms via email and sms. The Nagios and Red-eye Server could execute some feedback actions in order to normalize the status of the critical clients (red arrow). Finally, the monitor server send the status information of the most important services to the top level Dashboard that shows in a web page the overall status.

2. Dashboard

The Dashboard¹ is a tool developed by the INFN-CNAF staff which offers the opportunity, through a simple and intuitive web interface, to check the status of the fundamental services at the INFN-CNAF. The access to the Dashboard is not public but it is available only with Virtual Private Network (VPN) through “https” protocol. Information about status of all the divisions (Infrastructure, Network, Farming, Storage, Grid and National Services), are displayed in a concise manner which gives the possibilities to check, in a single web page, the status of the center’s services and infrastructure.

Each department have decided which services and server machines are considered critical for the functioning of the entire computer center. A MySQL database (DB) keeps track of states of all monitored services. Each department has a dedicated table in the database. The content of the DB is then parsed by appropriate scripts² and sent in html format to a web page which auto updates every five minutes allowing, through a convenient scroll bar (activated with the mouse or the arrows keyboard), to view the status of the center and to access to the alert history. The web page displays a colored box for each division, where the color indicates the condition (green = OK, orange = WARNING, red = CRITICAL, gray = UNKNOWN). Pending errors/warnings are shown inside each box, with a brief description and the time at which they occurred. If the

¹ The Dashboard is built in html, css, javascript (which must be enabled in your browser) and PHP and is based on a MySQL database and an Apache server. It has been tested with Internet Explore 7 and 8, Safari and Firefox 2 and 3.

² Python, Perl and PHP scripts are used.

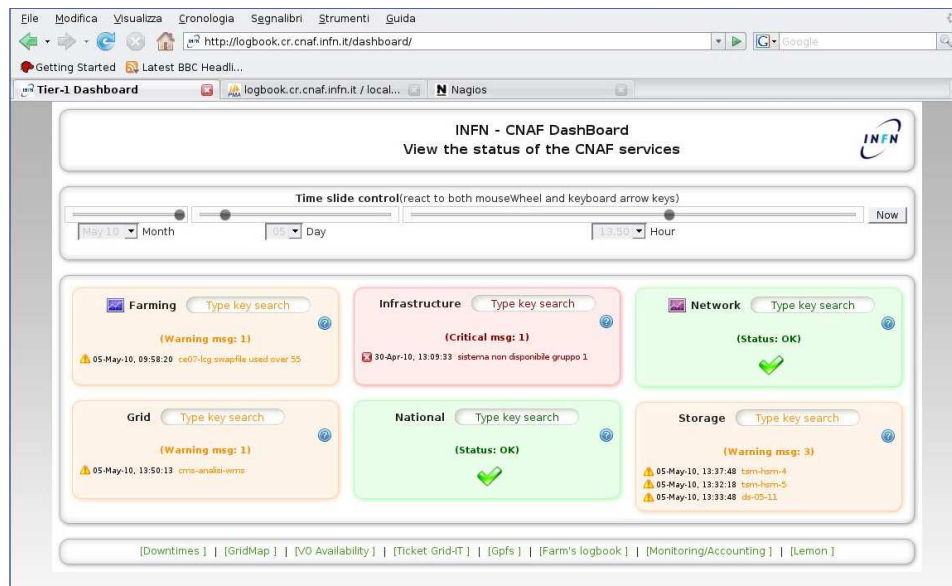


Figure 2. The Dashboard is a web page with a box for each INFN-CNAF internal division. The status of the basic services of the division is shown with a color in the box: green=OK, orange=WARNING, red=CRITICAL and gray=UNKNOWN. Pending Critical and Warning messages are shown inside each box with a brief description and the time at which they occurred.

mouse pointer is moved closer to the line which indicates the problem, a tool tip with more information and a link to a specific web page appears and it's possible to find details about the error and the relative solution. Each box has also a search field which offers the possibility to access to information in the MySQL database if a search is required over host names, services, or a partial string of an error (Fig. 2).

2.1. Dashboard Gateway (Dashgw)

The Infrastructural alarms are also collected into the Dashboard. The infrastructure has a proprietary monitoring system that could send an email or sms message to a specified recipe. The mail message receipt triggers a Python script called Dashgw. This script checks whether the sender, the recipient and the relay host are allowed to update the database. If they are, the mail subject and the body content are checked and if they follow the proper syntax, they are converted into three different possible status tags: OK, CRITICAL and WARNING. The python script then connects to the dashboard database and updates it according to the new status. In this way the database is always up to date with the infrastructure logbook and the status can be seen using the dashboard.

3. Storage

The storage division uses Lemon (LHC Era Monitoring) a CERN developed, server/client based, monitoring system available for Linux. On every monitored machine a lemon agent gathers information using the suitable sensors and forwards information via TCP/UDP to one or more servers. Information from these sensors can be obtained using a command line interface. In addition to a set of Lemon sensors, used for checking the agent or for standard system monitoring, new custom sensors can be written³, for particular purposes. Sensors information are stored in

³ Lemon custom sensors are written in Perl or CPP.

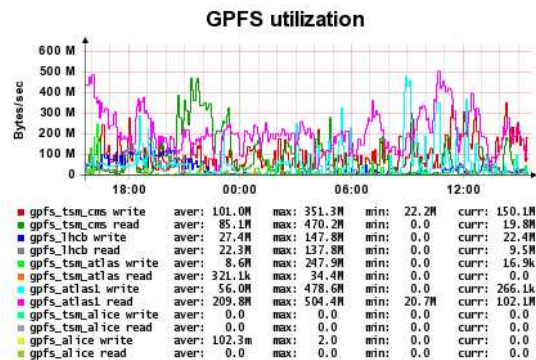


Figure 3. Lemon is one of the monitoring system adopted at INFN-CNAF, used by storage and grid division. The figure represent the GPFS traffic plot.

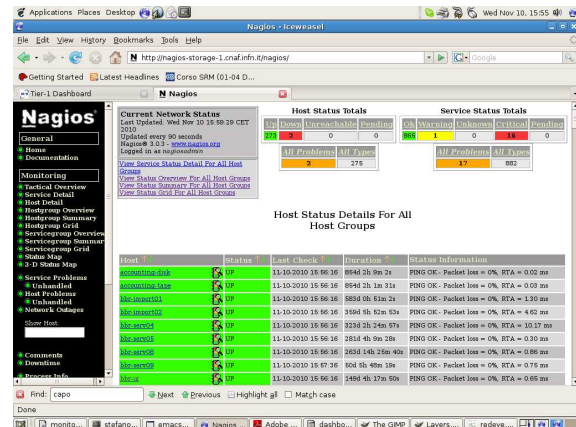


Figure 4. Nagios is an open source computer and network monitoring software application. It supervises hosts and services, alerting users when things go wrong and again when they get better. It is used to monitor Storage, Network, Grid and National Services divisions.

an Oracle DB and can be accessed from a web interface called lemon-web⁴. The web interface displays RRDtool [5] plots from single entities (e.g. hosts or services) or from clusters (e.g. clusters of hosts, clusters of clusters) (Fig. 3).

At INFN-CNAF, specific sensors have been developed for specific needs; in particular General Parallel File System (GPFS, [6]) utilization, Fibre Channel traffic, Network File System (NFS) servers, Tivoli Storage Manager (TSM, [7]) operations and status. The automatic monitoring update of specific grid hosts is possible due to lemon flexibility; each time a new host is added to the system, a lemon configuration file, reading this information from a DB, updates lemon grid cluster structure with the new host.

Nagios [3] (Fig. 4) is a multi devices (hosts, switches, routers, printers, etc.) and services monitoring tool. Nagios has wide set of its own plugins or can be customized with third party plugins. The monitoring daemon runs intermittent checks on hosts and services using plugins which return status information to the server. Current status information, historical logs, and reports can all be visualized in a web browser. Nagios is configured to send alarm messages (warning or critical) via mail, sms and there is a web interface. The Nagios advanced feature has been exploited in order to control the state of hosts and services correlated to the state of different hosts or services (the so called cluster control). In particular, set of plugins have been developed in order to control GPFS clusters/services, TSM server or TSM client clusters and StoRM [8] clusters. In some critical cases Nagios can also execute corrective action like the restart of the service. If a service fails Nagios can try to restart it once and if the check is still wrong it give up the restart and simply notify the failure. This can be useful and fully configurable.

4. Farming

Red-eye [4] is an INFN-CNAF developed software used for monitoring, reporting, alarming and self-acting and it is used by the Farming division. It uses a technology where the clients

⁴ Lemon-web interface exploiting Apache, PHP and Python.

periodically send collected data to one single collector server. In this configuration a single server is enough for monitoring about 1000 client nodes. Every five minutes the Red-eye clients perform the checks and they take the following actions: report the value on an auto refresh web Page and in case of an alarm condition send an e-mail and a report to the Dashboard.

5. Infrastructure

Schneider TAC Vista [10] is a commercial software used at INFN-CNAF for monitoring the infrastructure. A set of dedicated software tools (Business Development Manager) tackles the cooling and electrical power issues of a high density data center as our Tier1. Such tools were implemented and customized for the Tier1 data center requirements imposed to provide redundant cooling and power to the hosted IT equipment. They allow to monitor all components involved and alert faults to the dashboard. The power consumption data are used to optimize the efficiency and reduce costs.

6. Grid

INFN-CNAF Grid Operations service is mainly dedicated in providing, managing and supporting the so-called Grid Core services, provided by Glite [11] middleware, and different testbeds assigned to groups of developers for testing and releasing new services. Requirements in terms of monitoring, high availability, national and international core services and testbeds, are different but all of them use the same instruments, from server installation to reports. Nagios service for Grid Operations service is one of the used components for managing about 300 servers.

The information about all the servers are recorded in a database [12] which holds all the data starting from the buying process to employment, maintenance and phase out. This DB has a customized web management interface which is available using a X.509 digital certificate and is the main source of information for the monitoring tools Lemon and Nagios. When a server is added or modified in the DB, the changes are automatically spread in Lemon and Nagios. By default the Nagios checks that are included on all the Grid Operation servers are: ip-hostname resolution check; check of ssh. Checks for the most common services like MySQL, Apache, Tomcat, etc. are automatically configured for every server if required in the DB. Check for other services must be manually configured. For different critical services, Nagios automatically takes actions which have been configured in case of failure, by means of the event handler mechanism. In case these automatic actions don't re-establish the service, an sms advise is scheduled besides e-mail.

7. Network and National Services

INFN-CNAF network monitoring system is made of three main elements: statistical of bandwidth utilization of the single port, flows analysis of aggregate devices and state of health of network device in general. The monitoring of the three different aspects described above is committed to different independent systems which characterize a single control environment. Each of the three used tools relies on an open source software or on a customized developed solution. The open source tool Multi Router Traffic Grapher (MRTG, [13]) takes care of collecting statistic band utilization of all network device. This system, based on Simple Network Management Protocol (SNMP, [14]), is able both to read and to store various network interfaces counters and to create graphs showing network trend. Besides, thanks to saved counters, it has been possible, from MRTG, to develop a kind of Loadmap (Fig. 5), namely a graphics description based on progress bars of loading charge of the up-links of all the devices. In this way, it's easy to find out, for example, which switch is connected to the more band consuming servers and to highlight bottlenecks. On the other hand, the flows analysis give full visibility to the network utilization because it allows to have real time information about network utilization (host/network conversation) and which protocol (TCP/UDP) is used. These information are

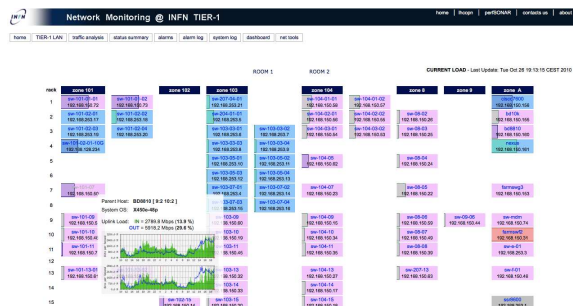


Figure 5. Loadmap is a graphics interface of loading charge of the uplinks of all devices.

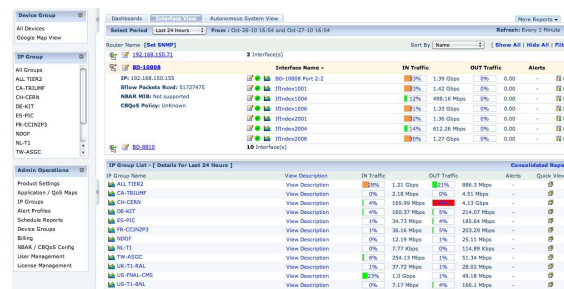


Figure 6. Collector/Analyzer is a network traffic probe that shows the link utilization in terms of network flows. It is useful to analyze network flows spawn by Cisco NetFlow and sFlow capable devices.

available on core switch using sFlow [15] and Cisco NetFlow [16] agents, which save, in a cache, information related to TCP/UDP communications and periodically export them in a collector/analyzer (Fig. 6). Several advantages for network managers are, for example, the possibility to know which applications are more band consuming (P2P, Web, DNS, etc.); to identify unauthorized network activities, to mark out source of possible network attacks and to define profile, priorities and traffic optimization. Another important element is the alarm system. It is based on a Nagios server which monitors the health state of devices and network services and notifies problems using e-mail and SMS. With its modular architecture Nagios allows system managers to monitor any device, metric or service which is “mission critical” for all the infrastructure using integrated applications. For example, integrating “arpwatch” and “splunk” in Nagios, it’s possible to have only one alarm system which handles alerts due to possible duplicated IP addresses in the network. The Nagios installation hosts also national services like: AFS [17] server, video and phone conference server, national mailing lists and LHC Computing Grid (LCG) collaborative video tool.

8. Conclusion

This article describes the control and monitoring system of the Italian Tier1 located at INFN-CNAF, Bologna. The Tier1 hosts services and mass storage systems for LHC experiments and many other experiments. Due to its complexity, a substantial number of monitoring systems is required. The adopted Solutions vary from open source system to commercial software and even to own developed systems. The different monitoring and notification system are joined into an automatic mechanism of problem solving and escalation: in fact, following each alarm, it is possible to take an action able to restore services, to forward the trouble to an expert team or to notify it to an higher level monitoring system. A global dashboard including state of services constituting a modern IT infrastructure (storage system, farm of servers, batch system, etc.) has been developed and represents the highest level of our system. Main advantages arose from using different independent low level monitoring systems, each one of them taking care of a single aspect of the Tier1 activity. First of all flexibility and scalability; as IT infrastructure becomes more complex, monitoring system grows with it, integrating new kind of controls and applications, without modifying the existing structure. On the other hand, malfunctioning of a monitoring element does not compromise other elements availability. However the possibility to consolidate some of the low level systems with similar characteristics should be investigated.

References

- [1] D.Gregori et al., *INFN-CNAF activity in the TIER-1 and GRID for LHC experiments*, Proceedings of IEEE XXIII International Parallel & Distributed Processing Symposium (IPDPS), Rome, 2009.
- [2] LHC Era Monitoring; <http://lemon.web.cern.ch/lemon/index.shtml>
- [3] Nagios: <http://nagios.org>
- [4] Red-eye: <http://tier1.cnaf.infn.it/monitor/>
- [5] RRDtool: <http://www.mrtg.org/rrdtool/>
- [6] GPFS: <http://www-03.ibm.com/systems/software/gpfs/>
- [7] TSM: <http://www-01.ibm.com/software/tivoli/products/storage-mgr/>
- [8] StoRM: <http://storm.forge.cnaf.infn.it/>
- [9] GNUplot: <http://www.gnuplot.info/>
- [10] TAC: http://www.schneider-electric.com/sites/corporate/en/solutions/business_segments/office-buildings/building-equipment/building-management-tac-vista.page
- [11] glite: <http://glite.web.cern.ch/glite/>
- [12] <https://gstore.cnaf.infn.it/NewEntropy/>
- [13] MRTG: <http://oss.oetiker.ch/mrtg/>
- [14] SNMP: http://en.wikipedia.org/wiki/Simple_Network_Management_Protocol
- [15] sFlow: <http://sflow.org>
- [16] Cisco NetFlow: http://www.cisco.com/en/US/products/ps6601/products_ios_protocol_group_home.html
- [17] AFS: <http://www.openafs.org/>