

Towards a top-level BDII hosted by the GOC

Brain Bockelman,¹
University of Nebraska, Lincoln

Flavia Donno,²
CERN

Dan Fraser,³
Argonne National Laboratory

Soichi Hayashi,⁴ Rob Quick,⁵ Scott Teige,⁶
Indiana University

Burt Holzman,⁷
Fermi National Accelerator Laboratory

and

Xin Zhao⁸
Brookhaven National Laboratory

1 Introduction

We present a plan to implement a top level BDII service at the Grid Operations Center (GOC). The GOC currently operates several services for the OSG including a site level BDII service. It has been previously proposed and presented to OSG management that a North American site host such a service and that proposal is included here as appendix A.

1.1 Overview

It is proposed that the OSG purchase five servers to be housed in the Indiana University (IU) data center in Bloomington. A top level BDII service operated by the GOC will be implemented on these servers. In the following, we describe the proposed system from the hardware level to the expectations of end users.

2 End Point “A”, Hardware Implementation

The servers hosting the service will be located in the IU datacenter in Bloomington. This differs from the typical implementation used by GOC for its other services in that most GOC services are hosted at both Bloomington and Indianapolis. When initially implemented, the Bloomington data center did not exist and this distributed approach addressed reliability concerns. We believe the IU data center is sufficiently reliable that these concerns no longer exist.

¹bbockelm@cse.unl.edu

²Flavia.Donno@cern.ch

³fraser@anl.gov

⁴hayashis@indiana.edu

⁵rquick@iupui.edu

⁶steige@indiana.edu

⁷burt@fnal.gov

⁸xzhao@bnl.gov

2.1 Physical Facilities

The IU data center is designed to withstand category 5 tornadoes. The facility is secured with card-key access and 7 x 24 x 365 video surveillance. Only staff with systems or network administration privileges have access to the machine room. Fire suppression is provided by a double interlock system accompanied by a Very Early Smoke Detection Apparatus (VESDA). Three circuits feed the Data Center, traveling redundant physical paths from two different substations. Any two circuits can fully power the building. A flywheel motor/generator set conditions the power and provides protection against transient events and uninterruptible power supplies protect against failures of moderate (~1 hour) duration. Dual diesel generators can provide power for 24 hours in the event of a longer term power failure. In house chillers provide cooling. Externally supplied chilled water plus city water can be used in the event of a failure of this system.

The IU Bloomington campus connects to the IU/Purdue University at Indianapolis (IUPUI) campus via I-Light⁹. Two redundant fiber bundles following different paths implement this connection. From the Indiana GigaPOP in Indianapolis, IU has the following connectivity to external networks:

- Internet2: 2 x 10 GigE
- NLR: 1 x 10 GigE
- TeraGrid: 2 x 10 GigE
- Commodity Internet: 2 x 10 GigE

IU is responsible for the operations of many national and international networks, including Internet2, National Lambda Rail (NLR), TransPac2, the MAN LAN research exchange point in New York City, the Hybrid Optical Network Initiative (HOPI), the Indiana GigaPOP, the TeraGrid's IP-Grid network, and the CIC OmniPoP in Chicago.

2.2 Server Configuration

Because of the rapidly changing nature of the server market we do not specify the actual server configuration at this time. We present here the specifications and price of a server with capabilities superior to the CERN configuration, see sections 4.1 and 4.2 for descriptions of other implementations.

- Processor: 2x 6 Core, 2.8 GHz, 6MB cache, AMD 2439SE
- RAM: 32 GB
- Hard Drive Configuration: RAID 5 or 10 PERC 6/i Controllers
- Drives: 160GB 7.2K RPM, 3.5in Hotplug Hard Drive (x3)
- Total Storage: 320 GB
- System Cost: \$4370

Five such servers would be required giving a total of 60 cores and 160 GB of RAM. The hard drives are the smallest currently available from DELL and the GOC considers RAID plus hot spare to be the minimum acceptable configuration for a production service.

2.3 Fail-over and Maintenance

Linux Virtual Server (LVS)¹⁰ will be used to load balance and provide fail-over capability.

Experience has shown that DNS round robin (RR) is inadequate for this purpose. Often, external DNSs do not obey the keep-alive time parameter so when a server is removed from the RR users depending on one of these DNSs often get routed to the server removed from service. Additionally, the DNS configuration at IU is not controlled by the GOC and changes to it occur only at the top of the hour. This can lead to an outage of up to one hour.

⁹<http://www.iupui.edu/ilight/>

¹⁰<http://www.linuxvirtualserver.org/>

LVS addresses both of these problems. The service will have a single IP address with the work shared transparently across a collection of servers. The configuration of LVS will be controlled by the GOC, changes can be made at any time.

Control of the LVS configuration will allow the number of servers handling load to be quickly changed during periods of maintenance or server failure. For maintenance a single server can be removed from the service, modified as required and returned without the delay of an externally controlled DNS.

3 End Point “B”, User Expectations

3.1 Availability/Reliability

The current Service Level Agreement¹¹ for the GOC BDII states:

The GOC will strive for 99% service availability. If service availability falls below 99% monthly as monitored by the GOC on two consecutive months a service plan will be submitted to the OSG stakeholders for plans to restore an acceptable level of service availability.

A maximum of two non-scheduled outages will be accepted by OSG during each six month period of service. If the GOC experiences more than the allotted outage, a service plan will be submitted to the OSG stakeholders with plans to restore the service to an acceptable level of operations.

It is anticipated the proposed service will be expected to achieve similar performance. Table 1 gives the achieved reliability and availability of the GOC BDII for the past six months.

Service	May	Jun	Jul	Aug	Sep	Oct
OSG BDII-1	100/100	100/100	100/100	99.66/99.66	99.73/99.89	100/100
OSG BDII-2	90.97/100	100/100	99.91/99.91	99.73/99.73	99.91/99.96	99.95/99.95

Table 1: Availability/Reliability (%) of GOC BDII services for the past six months

3.2 Use cases

Flavia Donno is currently preparing a document detailing use cases for ATLAS and CMS. DQ2 uses BDII to get details of storage endpoints and channel information. Panda job submission uses it to identify the compute elements (CE) that support ATLAS. FTS uses it once per day to configure file transfers. Glite and Glideins use it ~once per hour for configuration information. CRAB uses it to find storage elements (SE) and match them with the nearest CE.

3.3 Update Frequency

Ten minutes at the GOC, One hour at CERN, please verify. From the use cases at hand it looks to me like an hour is adequate

4 Current Implementations, Experience and Plans

4.1 CERN lcg

The top level BDII service at CERN currently uses 5 machines of the type used as standard batch nodes (1-2 disks). One has 4 cores and 8GB the others have 8 cores and 16 GB of RAM for a total of 36 cores and 72 GB of RAM. A recent 100 iteration test revealed that this service exhibits occasional long response times. A typical response time to a query from Bloomington was 460 msec but times of 563, 408, 72 and 10 seconds were encountered indicating the system is under-configured for peak demand.

¹¹<https://twiki.grid.iu.edu/bin/view/Operations/BDIIServiceLevelAgreement>

4.2 GOC

GOC hosts a site level BDII service for the OSG. Currently two instances are implemented on physical (as opposed to virtual) machines. A DNS round robin (keep-alive time 60 minutes) makes these available to users. This system handles ~2,000 queries per minute averaged over time. A total of 8 2.66 GHz cores and 16 GB of RAM are used. This system can handle current demand with occasional degradation of service quality during periods of unusually high demand. Running with a single server results in degraded performance sufficiently severe there are query timeouts and service alarms.

GOC is currently involved in an upgrade plan that involves adding additional servers on virtual machines. These VMs will be kept current via an rsync mechanism with an existing server. After implementation of these additional servers LVS will be used to replace the DNS round robin.

A Previous proposal

A.1 WLCG Deployment Strategy for Top-Level BDIIs

The Top-Level BDII services enable the discovery of Grid services along with further information about their structure and state. They aggregate information from all Site-Level BDIIs to provide a single point which can be queried to find the overall status of the Grid. Each instance of the service must be carefully managed as a critical service for WLCG. This not only requires timely support to ensure that issues are dealt with in a timely manner but also that managed fail over is in place for cases of service unavailability. In addition the service must be scalable to the number of queries which are made against it and also to the volume of data describing all the services. This proposal identifies the critical instances of the service for WLCG in order to ensure that those instance are well managed and that support effort is focused only the instances of the service that are mission critical for WLCG.

A.2 Selection Criteria

The first important consideration is query load. Each service instance must be scalable to exceed the typical query load to which it is exposed in-order to handle peak load. As each computing job can potentially query the Top-Level BDII service there is a strong correlation between the number of logical CPUs in a cluster that is configured to use a specific Top-Level BDII service and the number of queries that it experiences. The distribution of the BDII instances should therefore reflect the distribution of computing resources . Queries also originate from other sources such as the WMS, FTS and Grid Monitoring utilities. Although they do not place a high query load on the service, the queries they use are typical expensive such as returning all information (currently 100Mb).

The other consideration is network latency. Network latency significantly affects the query response time and the reliability of the queries. Due to this, compute resources and other consumers of information should query a 'close' instance from the perspective of network latency and also fail-over to a close instance.

The third consideration is available support effort. The service must be managed as a critical service for WLCG, which requires active support with releases being carefully followed. The services must be intensively monitored, to quickly identify failures, and also to measure the query loading in-order to provide additional resources when required. Currently, five physical hosts are recommend for a Top- Level service.

A.3 Identification of Instances

The latency and fail over requirements suggests that there should be at least two instance per continent: two in the North America, two in Europe and two in Asia. The requirement that the number of instances of the Top-Level service should reflect the distribution of computing resources suggests that the number of instances for Europe should be higher. Although sites in the US contribute a significant amount of computing resources, local policy results in less queries being made from the compute resources themselves and as such this requirement can be relaxed. The initial suggestion is;

- North America: Triumf and BNL or FNAL
- Europe: CERN, RAL, FZK, PIC, CNAF, NIKHEF
- Asia: Taiwan and KEK or Tokyo

This list can be extended as the need is identified.

A.4 Configuration for Sites

Sites should configure their primary Top-Level BDII to be the closest instance from the perspective of the network. For configuration options where fail-over is available, the fail-over instance should be the second closest from the perspective of the network.