

# EUROPEAN MIDDLEWARE INITIATIVE

## LOGGING AND BOOKKEEPING – ADMINISTRATOR’S GUIDE

---

Document version:	<b>1.2.9</b>
EMI Component Version:	<b>3.x</b>
Date:	<b>June 21, 2011</b>

---

This work is co-funded by the European Commission as part of the EMI project under Grant Agreement INFSO-RI-261611.

Copyright © Members of the EGEE Collaboration. 2004. See <http://www.eu-egee.org/partners/> for details on the copyright holders.

Licensed under the Apache License, Version 2.0 (the "License"); you may not use this file except in compliance with the License. You may obtain a copy of the License at

<http://www.apache.org/licenses/LICENSE-2.0>

Unless required by applicable law or agreed to in writing, software distributed under the License is distributed on an "AS IS" BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the License for the specific language governing permissions and limitations under the License.

## CONTENTS

<b>L&amp;B DOCUMENTATION AND VERSIONS OVERVIEW</b>	<b>5</b>
<b>1 INTRODUCTION</b>	<b>7</b>
1.1 SERVICE OVERVIEW . . . . .	7
1.1.1 L&B API AND LIBRARY . . . . .	8
1.1.2 LOGGER . . . . .	8
1.1.3 SERVER . . . . .	8
1.1.4 PROXY . . . . .	9
1.2 DEPLOYMENT SCENARIOS . . . . .	9
1.2.1 STANDALONE L&B SERVER . . . . .	9
1.2.2 HYBRID L&B SERVER-PROXY . . . . .	10
1.2.3 L&B SERVER ON WMS NODE . . . . .	10
<b>2 INSTALLATION AND CONFIGURATION</b>	<b>12</b>
2.1 COMPLETE LIST OF PACKAGES . . . . .	12
2.2 COMMON LOGGING FORMAT . . . . .	13
2.3 L&B SERVER . . . . .	14
2.3.1 HARDWARE REQUIREMENTS . . . . .	14
2.3.2 STANDARD INSTALLATION . . . . .	14
2.3.3 MIGRATION TO A DIFFERENT OS VERSION . . . . .	16
2.3.4 MIGRATION OF DATABASE TO SUPPORT TRANSACTIONS . . . . .	17
2.3.5 MIGRATION TO L&B 2.X . . . . .	17
2.3.6 MIGRATION TO L&B 3.X . . . . .	18
2.3.7 CONNECTING TO THE MESSAGING INFRASTRUCTURE . . . . .	18
2.3.8 INDEX CONFIGURATION . . . . .	18
2.3.9 AUTHORIZATION POLICY . . . . .	19
2.3.10 EXPORT TO R-GMA . . . . .	21
2.3.11 DATA BACKUP . . . . .	21
2.3.12 PURGING OLD DATA . . . . .	21
2.3.13 EXPLOITING PARALLELISM . . . . .	22
2.3.14 TUNING DATABASE ENGINE . . . . .	22
2.4 L&B PROXY . . . . .	22
2.5 L&B LOGGER . . . . .	23
2.6 L&B HARVESTER . . . . .	23
2.7 SMOKE TESTS . . . . .	23
2.7.1 JOB REGISTRATION . . . . .	23
2.7.2 LOGGING EVENTS VIA LB-LOGGER . . . . .	24

2.7.3	LOGGING EVENTS VIA LB-PROXY . . . . .	24
2.7.4	NOTIFICATION DELIVERY . . . . .	24
<b>3</b>	<b>MAINTENANCE</b>	<b>25</b>
3.1	CHANGING DEFAULT SETTINGS . . . . .	25
3.2	L&B SERVER AND PROXY . . . . .	25
3.2.1	STANDARD AND DEBUG LOGS . . . . .	25
3.2.2	CHANGING INDEX CONFIGURATION . . . . .	25
3.2.3	MULTIPLE INSTANCES . . . . .	26
3.2.4	BACKUP DUMPS . . . . .	27
3.2.5	PURGING AND PROCESSING OLD DATA . . . . .	27
3.2.6	ON-LINE MONITORING AND STATISTICS . . . . .	29
3.3	L&B LOGGER . . . . .	31
3.3.1	EVENT FILES . . . . .	31
3.3.2	BACKLOG REASONS . . . . .	31
3.3.3	NOTIFICATION DELIVERY . . . . .	32
3.3.4	DEBUG MODE . . . . .	32
3.4	USED RESOURCES . . . . .	32
3.4.1	SERVER AND PROXY . . . . .	32
3.4.2	LOGGER . . . . .	33
<b>4</b>	<b>TROUBLESHOOTING</b>	<b>34</b>
<b>5</b>	<b>FAQ—FREQUENTLY ASKED QUESTIONS</b>	<b>35</b>
5.1	JOB IN STATE 'RUNNING' DESPITE HAVING RECEIVED THE 'DONE' EVENT FROM LRMS . . . . .	35
5.2	WMS CANNOT PURGE JOBS OR PERFORM OTHER PRIVILEGED TASKS . . . . .	35
5.2.1	FOR <b>L&amp;B VERSION 3.0.11 OR HIGHER</b> , USING YAIM . . . . .	35
5.2.2	FOR ALL VERSIONS OF L&B, USING YAIM . . . . .	35
5.2.3	FOR <b>L&amp;B VERSION 2.1 OR HIGHER</b> , WITHOUT YAIM . . . . .	35

## L&B DOCUMENTATION AND VERSIONS OVERVIEW

The Logging and Bookkeeping service (L&B for short) was initially developed in the EU DataGrid project<sup>1</sup> as a part of the Workload Management System (WMS). The development continued in the EGEE, EGEE-II and EGEE-III projects,<sup>2</sup> where L&B became an independent part of the gLite<sup>3</sup> middleware [1], and then in the EMI Project.<sup>4</sup>

The complete L&B Documentation consists of the following parts:

- **L&B User's Guide** [2]. This user's guide explains how to use the Logging and Bookkeeping (L&B) service from the user's point of view. The service architecture is described thoroughly. Examples on using L&B event logging command to log a user tag and change job ACL are given, as well as L&B query and notification use cases.
- **L&B Administrator's Guide** - this document. This administrator's guide explains how to administer the Logging and Bookkeeping (L&B) service. Several deployment scenarios are described together with the installation, configuration, running and troubleshooting steps.
- **L&B Developer's Guide** [3]. This developer's guide explains how to use the Logging and Bookkeeping (L&B) service API. Logging (producer), querying (consumer) and notification API as well as the Web Services Interface is described in details together with programming examples.
- **L&B Test Plan** [4]. This test plan document explains how to test the Logging and Bookkeeping (L&B) service. Two major categories of tests are described: integration tests (include installation, configuration and basic service ping tests) and system tests (basic functionality tests, performance and stress tests, interoperability tests and security tests).

The following versions of L&B service are covered by these documents:

- **L&B version 3.0**: included in the EMI-1 Kebnekaise release
- **L&B version 2.1**: replacement for **L&B version 2.0** in gLite 3.2,
- **L&B version 2.0**: current stable version, in production as part of gLite 3.2,
- **L&B version 1.x**: old stable version, in production as part of gLite 3.1.

The older version of L&B that appeared in gLite 3.0 became obsolete and is not maintained anymore.

L&B packages can be obtained from two distinguished sources:

- **gLite releases**: gLite node-type repositories, offering a specific repository for each node type such as *glite-LB*
- **emi releases**: EMI repository or EGI's UMD repository, offering all EMI middleware packages from a single repository and relying on EPEL for dependencies

<sup>1</sup> <http://eu-datagrid.web.cern.ch/eu-datagrid/>

<sup>2</sup> <http://www.eu-egee.org/>

<sup>3</sup> <http://www.glite.org>

<sup>4</sup> <http://www.eu-emi.eu/>

*Note:* Despite offering the same functionality, binary packages obtained from different repositories differ and switching from one to the other for upgrades may not be altogether straightforward.

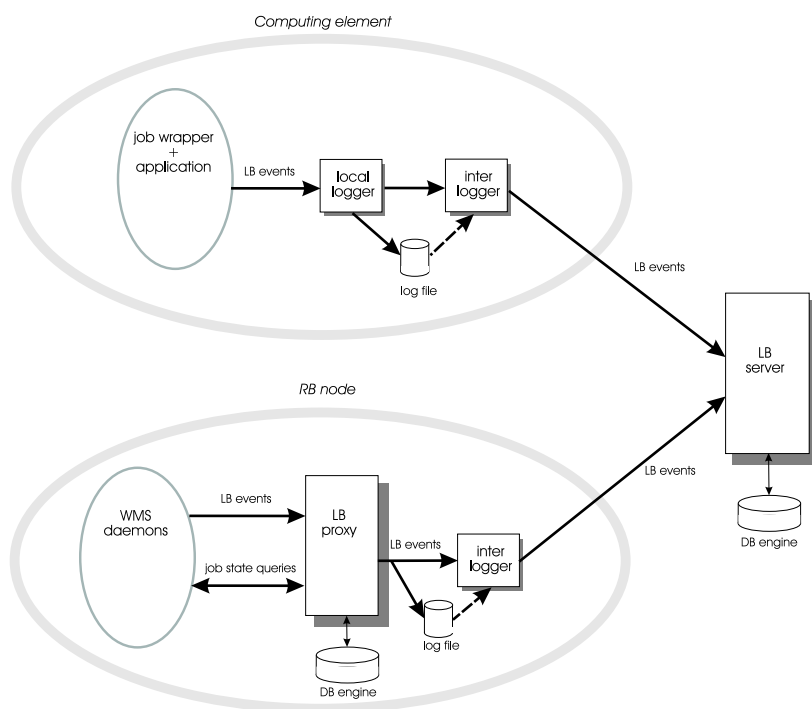
Updated information about L&B service (including the L&B service roadmap) is available at the L&B homepage: <http://egee.cesnet.cz/en/JRA1/LB>

# 1 INTRODUCTION

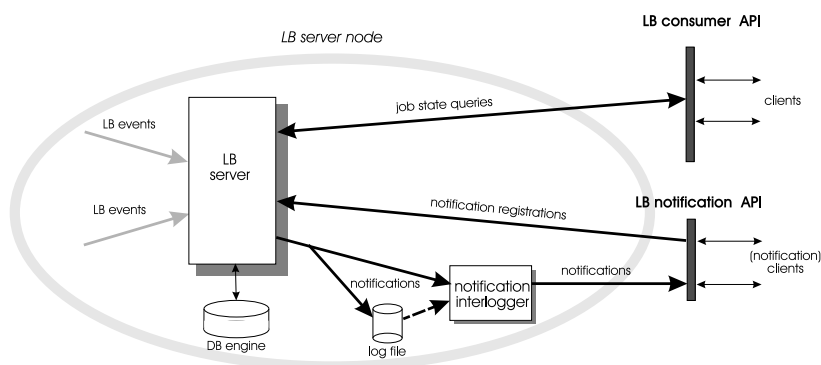
## 1.1 SERVICE OVERVIEW

A fairly complete overview of the L&B service is given in L&B User's Guide [2]. This section is a brief excerpt only, providing minimal information necessary for understanding the rest of this document.

The task of L&B is gathering *L&B events* from various grid middleware components (see Figure 1) and delivering the events to L&B servers where users can query for them (Figure 2).



**Figure 1:** Components involved in gathering and transferring L&B events



**Figure 2:** L&B queries and notifications

### 1.1.1 L&B API AND LIBRARY

Both logging events and querying the service are implemented via calls to a public L&B API. The complete API (both logging and queries) is available in ANSI C binding, most of the querying capabilities also in C++. These APIs are provided as sets of C/C++ header files and shared libraries. The library implements communication protocol with other L&B components (logger and server), including encryption, authentication etc. Since **L&B version 2.0** an experimental Java binding of the logging API is available.

We do not describe the API here in detail; it is documented in L&B Developer's Guide[3], including complete reference and both simple and complex usage examples.

Events can be also logged with a standalone program (using the C API in turn), intended for usage in scripts.

The query interface is also available as a web-service provided by the L&B server (Sect. 1.1.3).

Finally, certain frequent queries (all user's jobs, single job status, ...) are available as HTML pages (by pointing ordinary web browser to the L&B server endpoint), or as simple text queries (since **L&B version 2.0**) intended for scripts. See [2] for details.

### 1.1.2 LOGGER

The task of the *logger* component is taking over the events from the logging library, storing them reliably, and forwarding to the destination server. The component should be deployed very close to each source of events—on the same machine ideally, or, in the case of computing elements with many worker nodes, on the head node of the cluster<sup>5</sup>.

Technically the functionality is realized with two daemons:

- *Local-logger* accepts incoming events, appends them in a plain disk file (one file per Grid job), and forwards to inter-logger. It is kept as simple as possible in order to achieve maximal reliability.
- *Inter-logger* accepts the events from the local-logger, implements the event routing (currently trivial as the destination address is a part of the jobid), and manages delivery queues (one per destination server). It is also responsible for crash recovery—on startup, the queues are populated with undelivered events read from the local-logger files. Finally, the inter-logger purges the files when the events are delivered to their final destination.

### 1.1.3 SERVER

*L&B server* is the destination component where the events are delivered, stored and processed to be made available for user queries. The server storage backend is implemented using MySQL database.

Incoming events are parsed, checked for correctness, authorized (only the job owner can store events belonging to a particular job), and stored into the database. In addition, the current state of the job is retrieved from the database, the event is fed into the state machine and the job state updated accordingly.

On the other hand, the server exposes querying interface (Fig. 2). The incoming user queries are transformed into SQL queries on the underlying database engine. The query result is filtered, authorization rules applied, and the result sent back to the user.

---

<sup>5</sup>In this setup logger also serves as an application proxy, overcoming networking issues like private address space of the worker nodes, blocked outbound connectivity etc.



While using the SQL database, its full query power is not made available to end users. In order to avoid either intentional or unintentional denial-of-service attacks, the queries are restricted in such a way that the transformed SQL query must hit a highly selective index on the database. Otherwise the query is refused, as full database scan would yield unacceptable load. The set of indices is configurable, and it may involve both L&B system attributes (for example job owner, computing element, timestamps of entering particular state, ...) and user defined ones.

The server also maintains the active notification handles, providing the subscription interface to the user. Whenever an event arrives and the updated job state is computed, it is matched against the active handles<sup>6</sup>. Each match generates a notification message, an extended L&B event containing the job state data, notification handle, and the current user's listener location. The event is passed to the *notification inter-logger* via persistent disk file and directly (see Fig. 2). The daemon delivers events either in a standard way, using the specified listener as destination, or forwards them to a messaging broker for delivery through the messaging infrastructure. When using the standard delivery mechanism, the server generates control messages when the user re-subscribes, changing the listener location. Inter-logger recognizes these messages, and changes the routing of all pending events belonging to this handle accordingly.

#### 1.1.4 PROXY

*L&B proxy* is the implementation of the concept of local view on job state (see [2]). Since **L&B version 2.0**, L&B proxy is integrated into L&B server executable. When deployed (on the WMS node in the current gLite middleware) it takes over the role of the local-logger daemon—it accepts the incoming events, stores them in files, and forwards them to the inter-logger.

In addition, the proxy provides the basic principal functionality of L&B server, that is processing events into job state and providing a query interface, with the following differences:

- only events coming from sources on this node are considered; hence the job state may be incomplete,
- proxy is accessed through local UNIX-domain socket instead of network interface,
- no authorization checks are performed—proxy is intended for privileged access only (enforced by the file permissions on the socket),
- aggressive purge strategy is applied—whenever a job reaches a known terminal state (which means that no further events are expected), it is purged from the local database immediately,
- no index checks are applied—we both trust the privileged parties and do not expect the database to grow due to the purge strategy.

## 1.2 DEPLOYMENT SCENARIOS

### 1.2.1 STANDALONE L&B SERVER

This is a recommended standard production deployment.

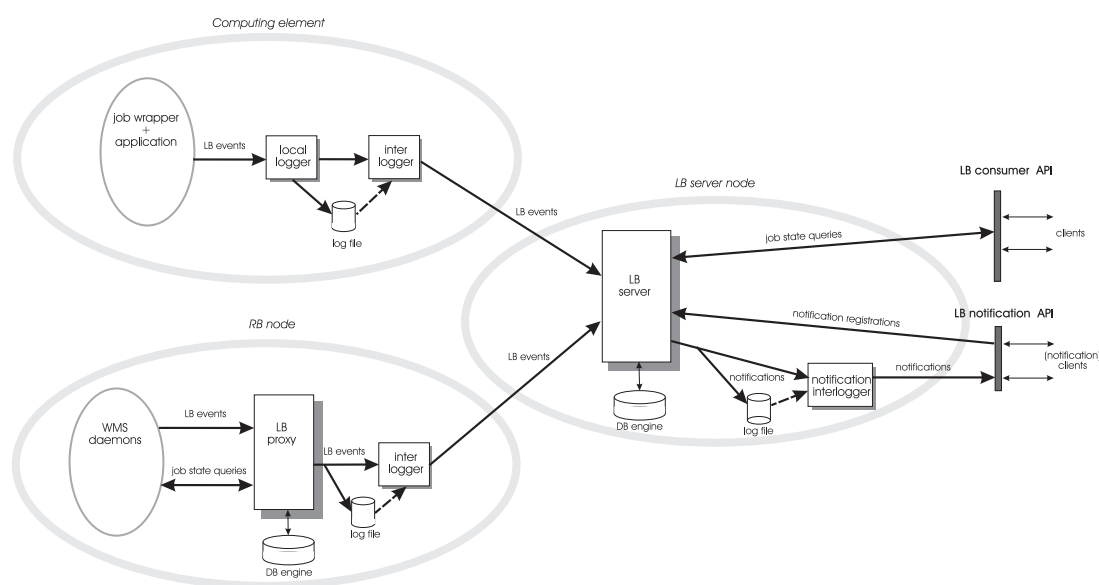
---

<sup>6</sup>The current implementation enforces specifying an actual jobid in the subscription hence the matching has minimal performance impact.

L&B server is installed on a dedicated machine, where no other grid services (gLite WMS in particular) run. Hence user queries and notification processing are offloaded from the WMS, not affecting its performance directly.

In this setup the full reported performance is achieved, currently up to several hundreds thousands jobs per day, with the goal of one million, see [4].

Further performance can be gained with clustering  $M$  WMSs and  $N$  L&Bs while configuring all WMSs to distribute jobs to the L&Bs uniformly. In this setup bottlenecks emerging from L&B proxy to L&B server serialized communication are likely to be eliminated. The optimal  $M : N$  ratio strongly depends on the rate of user queries and number of evaluated notifications, and it must be determined empirically for each specific installation.



**Figure 3:** L&B deployment – overall picture

### 1.2.2 HYBRID L&B SERVER-PROXY

L&B server runs on the WMS node, in combined server-proxy mode, serving both user queries and supporting WMS. Total processing requirements for a single jobs are lower (unlike with separated proxy and server, job state is computed and stored only once).

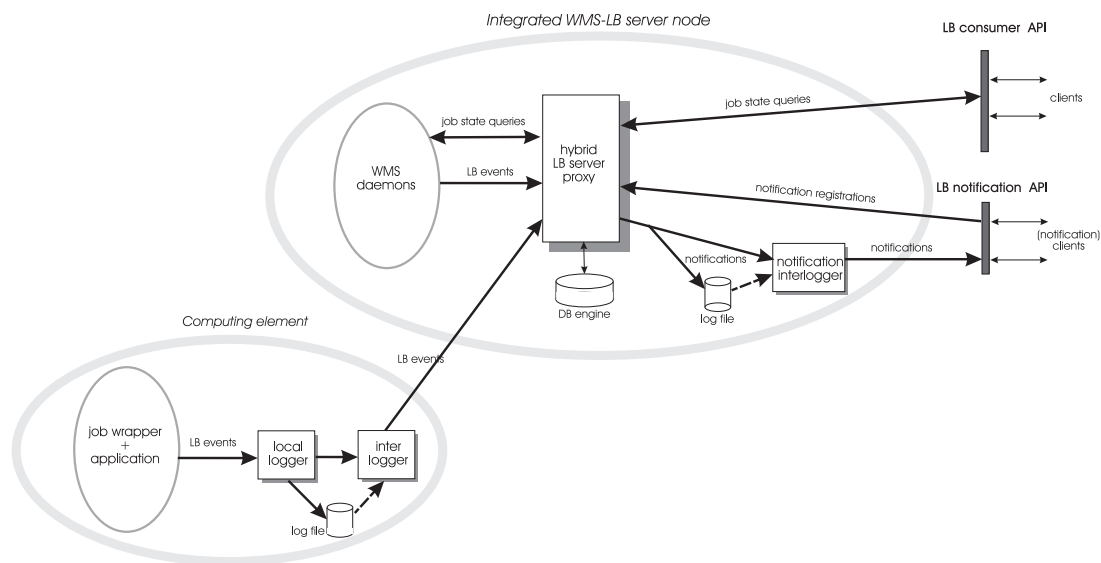
On the other hand, processing user queries is done on the WMS node, limiting its performance then. This setup is suitable for smaller installations with a single (unclustered) WMS, expected load of upto 30–50 kjobs/day, and not very heavy user-generated load.

The functionality is available since **L&B version 2.0**.

### 1.2.3 L&B SERVER ON WMS NODE

**This setup is obsolete and very inefficient, hence discouraged.**

Ancient L&B versions ( $< 1.2$ ), where L&B proxy was not available yet, used to be frequently installed with L&B server on the WMS machine. With the introduction of L&B proxy it makes little sense anymore but,



**Figure 4:** L&B deployment with combined server-proxy

unfortunately, this setup still persists at some sites. Its consequence is doubling both CPU and disk load, yielding observable performance degradation.

Large production sites should consider standalone L&B server (Sect. 1.2.1) instead, while for smaller sites the hybrid setup (Sect. 1.2.2) may be more appropriate.

## 2 INSTALLATION AND CONFIGURATION

### 2.1 COMPLETE LIST OF PACKAGES

L&B is currently distributed mainly in RPMs packages. It is available also in binary form packed as .tar.gz. Recent attempts to multiplatform porting and recent ETICS building system development promise a future possibility to distribute the software in other distribution formats, e.g. DEB packages.

In **L&B version 1.x**, the list of all LB packages was the following:

glite-lb-common	common files
glite-lb-client	client library and CLI tools
glite-lb-client-interface	client library interface (header files)
glite-lb-harvester	enhanced L&B notification client (since <b>L&amp;B version 1.10</b> )
glite-lb-logger	local-logger and inter-logger
glite-lb-proxy	proxy (restricted server used by WMS)
glite-lb-server	server
glite-lb-server-bones	multi-process server building library
glite-lb-utils	auxiliary utilities
glite-lb-ws-interface	web service interface
glite-security-gsoap-plugin	GSS wrapper and GSS plugin for gSoap

Starting with **L&B version 2.0**, the code has been restructured quite a lot, especially the dependencies were lightened, and the new list of packages – applicable also to **L&B version 3.x** is now the following:

glite-lb-doc	documentation
glite-lb-common	common files
glite-lb-client	client library and CLI tools
glite-lb-client-java	Java implementation of the client (since <b>L&amp;B version 2.1</b> )
glite-lb-harvester	enhanced L&B notification client (since <b>L&amp;B version 2.1</b> )
glite-lb-logger	local-logger and inter-logger
glite-lb-logger-msg	plugin for message delivery over ActiveMQ (since <b>L&amp;B version 3.0</b> )
glite-lb-server	server, including merged proxy functionality
glite-lb-state-machine	state machine and LB plugin for Job Provenance
glite-lb-utils	auxiliary utilities
glite-lb-ws-interface	web service interface
glite-lb-yaim	YAIM initialization scripts for L&B (since <b>L&amp;B version 2.1</b> )

More detailed description together with the dependencies can be read directly from each package, for example by issuing the command

```
rpm -qiR <package_name>
```

Some of the LB packages depend also on other gLite packages, different due to the restructuring since **L&B version 2.0**. For **L&B version 1.x** they are:

glite-wms-utils-jobid	gLite jobid management library
glite-jp-common	Job Provenance auxiliary library

And for **L&B version 2.0 and newer**:

<code>glite-jobid-api-c</code>	gLite jobid C API library
<code>glite-lbjp-common-db</code>	database access layer
<code>glite-lbjp-common-jp-interface</code>	interface to the Job Provenance service
<code>glite-lbjp-common-log</code>	glite common logging format implementation
<code>glite-lbjp-common-maildir</code>	persistent request spool management
<code>glite-lbjp-common-server-bones</code>	multi-process server building library
<code>glite-lbjp-common-trio</code>	extended printf implementation
<code>glite-lbjp-common-gss</code>	GSS wrapper (formerly <i>glite-security-gss</i> )
<code>glite-lbjp-common-gsoap-plugin</code>	GSS plugin for gSoap (formerly <i>glite-security-gsoap-plugin</i> )

where all `glite-lbjp-common-*` packages are common both to L&B and Job Provenance (JP).

## 2.2 COMMON LOGGING FORMAT

Since **L&B version 2.1** L&B service follows the **gLite common logging recommendations v1.1**:

<https://twiki.cern.ch/twiki/pub/EGEE/EGEEgLite/logging.html>.

The implementation is done in the `glite-lbjp-common-log` package and it uses `log4c` (<http://log4c.sourceforge.net>) and its configuration file `log4crc`.

There is one configuration file `[/opt/glite]/etc/glite-lb/log4crc` that startup scripts use by setting the `LOG4C_RCPATH` environment variable.

A file `log4crc.example-debugging` may be useful to copy to `$HOME/.log4crc` (or by setting the `LOG4C_RCPATH` environment variable to a directory containing the `log4crc` file) to obtain detailed debugging output. One can debug only specific parts of the LB system, for example by uncommenting `LB.SERVER.DB` category in the `log4crc` file, one gets only the debugging info related to the underlying database subsystem calls.

Currently supported log categories are:

<code>SECURITY</code>	Security matters
<code>ACCESS</code>	Communication issues
<code>CONTROL</code>	L&B server control
<code>LB</code>	<i>unused</i>
<code>LB.LOGD</code>	Messages from the L&B local logger
<code>LB.INTERLOGD</code>	Messages from the L&B interlogger
<code>LB.NOTIFINTERLOGD</code>	Messages from the notification interlogger
<code>LB.SERVER</code>	Processing within the L&B server
<code>LB.SERVER.DB</code>	DB calls made by the server
<code>LB.SERVER.REQUEST</code>	Processing L&B server requests
<code>LB.HARVESTER</code>	Messages from the L&B harvester
<code>LB.HARVESTER.DB</code>	DB calls made by the harvester
<code>LB.AUTHZ</code>	Authentication matters

The following log priorities are recognized, in a descending order of severity: `fatal`, `error`, `warn`, `info`, `debug`. Note that running with any of the logging categories set to `debug` for a prolonged period of time will result in the logfile growing uncontrolledly.

In a typical setup, logging messages generated by the L&B service will be found in `/var/log/messages`.

**Note:** `syslog` is typically configured to drop log messages with priority `debug`. You may want to enable it, for example by using something like this in `/etc/syslog.conf` and then restarting `syslog`:

```
*.debug;mail.none;authpriv.none;cron.none
```

```
-/var/log/messages
```

## 2.3 L&B SERVER

### 2.3.1 HARDWARE REQUIREMENTS

Hardware requirements depend on performance and storage time requirements. Disk space used by LB server consists of database space and working space for backup dumps and temporary files used by exports to Job Provenance and R-GMA tables. Necessary database space can be calculated by multiplying job retention length (job purge timeout), job submission rate, and per-job space requirements (120 KB per job is recommended for current common EGEE usage pattern; jobs can consume more than that with use of very long JDL descriptions, user tags, or very high number of resubmissions). For temporary files, approximately 10 GB is sufficient for LB server setups working normally, more can be needed when backlog forms in data export to any external service. For example, typical setup processing 40 000 jobs per day where all jobs are purged after 10 days needs about 58 gigabytes ( $10 \cdot 40000 \cdot 120\text{KB}$  (per job) + 10) not accounting for operating system and system logs.

For smooth handling of 40 000 jobs/day, this or better machine configuration is necessary:

- 1GB RAM
- CPU equivalent of single Xeon/Opteron 1.5GHz
- single 7200 rpm SATA disk.

In order to achieve higher performance, following changes are recommended:

- Faster disks. Disk access speed is crucial for LB server, couple of 15k rpm SCSI or SAS disks (one for MySQL database data file, the second for DB logs, LB server's working directories, and operating system files) or RAID with battery backed write-back cache is preferable.
- More memory. Large RAM improves performance through memory caching, relative speed gain is likely to be roughly proportional to memory/database size ratio. To use 3 GB or more efficiently, 64bit OS and MySQL server versions are recommended.
- Faster or more processors. CPU requirements scale approximately linearly with offered load.

### 2.3.2 STANDARD INSTALLATION

- Installing from **EMI or UMD repository** (applies to **L&B version 3.0** and higher):  
Install and configure a clean system, the **EPEL** repository, and the **EMI** repository as per instructions given in <https://twiki.cern.ch/twiki/bin/view/EMI/GenericInstallationConfigurationEMI1>. Then install metapackage `emi-lb`.
- Installing from **gLite's node-type repository** (applies to **L&B version up to 2.1**):  
Install and configure OS and basic services (certificates, CAs, time synchronization, software repositories) according to the <https://twiki.cern.ch/twiki/bin/view/LCG/GenericInstallGuide320>. Then install metapackage `glite-LB` from an appropriate gLite software repository.

YAIM configuration can be done then:

```
/opt/glite/yaim/bin/yaim -c -s site-info.def -n glite-LB
```

Available parameters specific to LB server are (mandatory parameters are **highlighted**):

- **MYSQL\_PASSWORD** – root password of MySQL server (mandatory)
- GLITE\_LB\_EXPORT\_PURGE\_ARGS – purge timeouts (default: -cleared 2d -aborted 15d -cancelled 15d -other 60d)  
 According to local retention policy you may want to use different purge timeouts (for example WLCG would need -cleared 90d -aborted 90d -cancelled 90d -other 90d).
- GLITE\_LB\_EXPORT\_ENABLED – set to **true** for export to JP, installed glite-lb-client and glite-jp-client are needed (default: false)
- GLITE\_LB\_EXPORT\_JPPS – Job Provenance Primary Storage where to export purged jobs, required if export to JP is enabled
- GLITE\_LB\_RTM\_ENABLED – enable settings for Real Time Monitor - indexes and additional access (default: false)
- GLITE\_LB\_TYPE – type of the L&B service: server, proxy, both (default: autodetect, L&B node only: 'server', WMS node only: proxy, L&B and WMS: 'both')
- GLITE\_LB\_INDEX\_OWNER – when specified, add (**true**) or drop (**false**) 'owner' index (default: 'owner' index not touched)
- GLITE\_LB\_MSG\_BROKER – URL of the MSG broker, 'auto' for looking in BDII, 'false' for disabling MSG notifications (default: auto)
- GLITE\_LB\_MSG\_NETWORK – required network type when searching in BDII (default: PROD)
- LCG\_GFAL\_INFOSYS – BDII servers (default: lcg-bdii.cern.ch:2170)

Authorization:

- GLITE\_LB\_SUPER\_USERS – additional super-users (default: empty)<sup>7</sup>
- GLITE\_LB\_WMS\_DN – DNs of WMS servers (default: empty)<sup>7</sup>
- GLITE\_LB\_RTM\_DN – DNs using to get notifications from L&B server (default: heppc24.hep.ph.ic.ac.uk machine certificate)
- GLITE\_LB\_AUTHZ\_<category> – more detailed tuning of access grants, see Section 2.3.9 (default: empty, '.' '\*' for logging and job registrations)

Additional helper or legacy parameters:

- GLITE\_LB\_LOCATION – L&B prefix (default: /opt/glite or /usr)
- GLITE\_LB\_LOCATION\_ETC – system config directory (default: /opt/glite/etc or /etc)

<sup>7</sup>The use of this parameter is a FAQ. See section 5.2.

- `GLITE_LB_LOCATION_VAR` – gLite local state directory (default: `/opt/glite/var` or `/var/glite`)
- `GLITE_JP_LOCATION` – can be used when JP subsystem location differs from LB (default: empty)
- `GLITE_LB_HARVESTER_ENABLED` – set to `true` for sending notifications, used mainly for legacy export to MSG publish system (default: `false`)
- `GLITE_LB_HARVESTER_MSG_OPTIONS` – additional options for MSG publish (default: `-wlcg`)
- `GLITE_WMS_LCGMON_FILE` – pathname of file where job state export data are written for use by lcgmon/R-GMA (default: `/var/glite/logging/status.log`). *Note: This feature is now obsolete and only available in **L&B version 1.x**.*

In addition to those, YAIM LB module uses following parameters:

- `INSTALL_ROOT` – installation root, shouldn't be changed (default: `'/opt'` or `'/'`)
- `GLITE_LOCATION_VAR` – directory for local state files (default: `'/opt/glite/var'` or `'/var/glite'`)
- `GLITE_USER` – unprivileged user name (default: `'glite'`)
- **`SITE_NAME`** – site name (mandatory)
- **`SITE_EMAIL`** – site email, for cron scripts (mandatory)

Lists are separated by comma.

### 2.3.3 MIGRATION TO A DIFFERENT OS VERSION

Migration of a LB server to different machine is possible using following simple procedure (just file copy of the MySQL database). We tested the migration from SL4 32bit (mysql 4.1.22-2) to SL5 64bit (mysql 5.0.45-7).

Steps:

- *Prepare a new machine.* The new machine must get the same hostname as the old machine had. It is a part of job ids stored in the database.
- *Move data.* Just stop the MySQL server and move `/var/lib/mysql` data directory directly to the target machine.
- *(optional) Restore file contexts.* You may need to restore file contexts in case of enabled SELinux.

For example, commands on the target machine:

```
service mysqld stop
cd /var/lib
tar xf /tmp/lb.tar
restorecond -R mysql
service mysqld start
```



### 2.3.4 MIGRATION OF DATABASE TO SUPPORT TRANSACTIONS

Started from version 1.4.3 of the `glite-lb-server` package, the L&B server introduced optional use of database transactions for L&B database updates in order to improve their performance. This feature is switched on by default when underlying MySQL database uses transactional InnoDB tables. For new installations, YAIM configuration process will create transactional database automatically. For existing LB server database the migration process is not automatically handled.

Note: If you want to add transaction when migrating to L&B 2.x skip this section and use L&B 2.x migration procedure. The migration of database to support transactions is included in L&B 2.x migration procedure.

Steps:

- *Stop the server.* Stop both a L&B server and a MySQL server. Making a fresh backup copy of database is a good idea.
- *Database conversion.* Use provided SQL script:

```
mysql -u lbserver lbserver20 \  
    <[/opt/glite]/etc/glite-lb-dbsetup-migrate2transactions.sql
```

- *Start the servers.* MySQL and L&B. Check logs.

### 2.3.5 MIGRATION TO L&B 2.x

The migration process of existing L&B 1.x database to the L&B 2.x is not handled automatically. The database schema change is required due to support of merged L&B server and proxy services using single database, pointers to purged jobs ("zombies") and other improvements.

Warning: There are two types of L&B database based on the fact that you can have a L&B server or L&B proxy. For more information about L&B proxy please see [2.4](#)

Steps:

- *Stop the server and upgrade to L&B 2.x.* Stop both a L&B server and a MySQL server. Making a fresh backup copy of database is a good idea. Do the upgrade to L&B 2.x, optionally you can move the database to new OS in this step (see [2.3.3](#)).
- *Before migration some database tuning is required.* Especially parameter `innodb_buffer_pool_size` needs to be increased, to support bigger transactions. For details see [Section 2.3.14](#)
- *Database type.* Check if you have a L&B server or a L&B proxy. In the following step you must properly set the switch `-s` (server) or `-p` (proxy).
- *Database conversion.* Use provided shell script (with the proper switch from previous step):

```
[/opt/glite]/etc/glite-lb-migrate_db2version20 {-s|-p}
```

- *(Optional) Drop unnecessary index.* This operation is likely to take a lot of time when applied to large database.

```
mysql -u lbserver lbserver20 -e "alter table events drop index host"
```

- (Optional) Check the L&B indexes. You may need to add LastUpdateTime for monitoring services. See [3.2.2](#)
- Start the servers. MySQL and L&B. Check logs.

This migration procedure is tested in following environment: LB 1.9.x from RPMs SL4 32bit (mysql 4.1.22-2), LB server node, migration to SL5 64bit (mysql 5.0.45-7) LB2.0 RPM.

### 2.3.6 MIGRATION TO L&B 3.x

There are no specific configuration changes required when migrating from **L&B version 2.x** to **L&B version 3.x**. In case you wish to migrate from **L&B version 1.x**, follow instructions given in section [2.3.5](#) but upgrade directly to **L&B version 3.0** in the first step.

*Note:* Upgrading from **L&B version 2.x** provided by gLite L&B node repository to **L&B version 3.0** or higher provided by the EMI repository has never been tested.

### 2.3.7 CONNECTING TO THE MESSAGING INFRASTRUCTURE

As of **L&B version 3.0**, the L&B server node becomes a producer of messages, which it delivers into the messaging infrastructure.

Correct broker settings are inferred from BDII by YAIM on configuration. By default, messaging-related settings are stored in file:

```
[/opt/glite]/etc/glite-lb/msg.conf
```

Overall, `msg.conf` specifies the following information:

- The messaging plugin to be used by the notification interlogger. Plugin settings should be correct *ab initio* and do not require modification by administrators.
- Broker settings (attribute `broker`). They may require an adaptive change in case the currently configured broker disappears and automatic checks fail to switch the settings to another one on time.
- Permissible topic title prefixes (attribute `prefix`). Registrations for delivery to topics not matching the prefix will be rejected. In case no prefix is specified, L&B applies the default EGI prefix: `grid.emi.`

### 2.3.8 INDEX CONFIGURATION

Initial YAIM configuration creates L&B indexes typically, the actual configuration depends on required features (for example RTM job monitoring). See Section [3.2.2](#) for instructions on changing L&B server index configuration afterwards in order to meet specific needs.

### 2.3.9 AUTHORIZATION POLICY

The L&B server applies a quite strict access control policy on the operations provided to the clients to ensure a sufficient level of data protection. By default, the information about a job is only available to the owner of the job. The job owner can specify an ACL assigned to their jobs specifying permissions granted to other users so that they could access the job records, too. More information about the ACL management can be found in the L&B Users' guide.

Apart from using the ACLs, the L&B server administrator can also set a server-level policy granting rights to perform particular operation on L&B server that are considered privileged. For example, a privileged user can access data about jobs owned by other users, bypassing the default L&B access control mechanism. **L&B version 2.1** specifies several categories of rights that can be granted to the users:

- ADMIN\_ACCESS
- READ\_ALL
- PURGE
- STATUS\_FOR\_MONITORING
- GET\_STATISTICS
- REGISTER\_JOBS
- LOG\_WMS\_EVENTS
- LOG\_CE\_EVENTS
- LOG\_GENERAL\_EVENTS
- GRANT\_OWNERSHIP (since **L&B version 3.0**)

The first action disables all authorization checks. The next four categories concern with acquiring data from the L&B server, while the other ones make it possible to define a web of trusted sources passing events to the L&B server.

ADMIN\_ACCESS is the most powerful privilege allowing to bypass any authorization checks on the server. It replaces the superuser role, which existed in **L&B version 2.0** and older. Note, that the `--super-users` command-line option still exists and translates internally into granting ADMIN\_ACCESS.

READ\_ALL enables to access all job information stored on the server. PURGE grants the privilege to ask for purging the L&B database. The L&B server's identity is automatically assigned the READ\_ALL and PURGE so that these operations are available for example to a cron script running on L&B node.

When granted to a user, the STATUS\_FOR\_MONITORING right allows the user to query statuses of all jobs maintained by the server, however only a small subset of the status fields is returned to back. For example, the caller does not obtain the identity of the job owner. The purpose of this right is to provide a way to collect information necessary for overall monitoring while preserving a basic level of privacy.

GET\_STATISTICS allows its bearers to query for on-line statistics generated by the L&B server. See [3.2.6](#) for more information about the purpose of the function.

**L&B version 1.x** allowed anyone possessing a trusted digital certificate to send an arbitrary event to the L&B server. While enabling an easy setup, such an arrangement does not comply with some contemporary requirements, for example job traceability, since the data could be distorted by a malicious

user. In order to strengthen the trustworthiness of the data provided, the **L&B version 2.1** server has introduced an authorization mechanism to control the originators of events. The authorization model presumes a set of trusted components that are only allowed to send "important" events, while other types can be logged from any source. The `REGISTER_JOBS` authorization category specifies clients allowed to register jobs with the L&B server. The `LOG_CE_EVENTS` category makes it possible to define a set of trusted CEs that are allowed to log events originating from within sites (in particular `RUNNING`, `REALLYRUNNING`, and `DONE`). Similarly, the `LOG_WMS_EVENTS` category defines a web of trusted WMS nodes. The `LOG_GENERAL_EVENTS` category comprises events that can be sent from any place on the grid, namely `CURDESC`, `USERTAG`, and `CHANGEACL`. It is important to understand that these access control rules provides additional level to the existing authorization routines. In particular, being granted the `LOG_GENERAL_EVENTS` right is not sufficient to e.g. change ACLs on jobs of other people and obtain an access to the job information.

As of **L&B version 3.0** it is possible for the job owner to hand over the ownership of the payload to another user, which eases handling for example of pilot jobs. In order to set the payload owner, the job owner must log a `GrantPayloadOwnership` event, where the identity of the new payload owner is introduced. The event can only originate from loggers described by the `GRANT_OWNERSHIP` category of the policy file. More information and an example of ownership setting is given in the L&B Users' guide.

The L&B policy is specified in a policy configuration file that must be given in the server configuration. Specifying the policy file also triggers the enforcement of access policy rights, especially the ones describing the event sources. If the policy is not enabled, the L&B server accepts events from any logger with a trusted certificate. The format of the policy is a subset of the Simplified policy language introduced by the Argus gLite authorization service<sup>8</sup>. Unlike the Argus language, the L&B policy supports only certificate subject names and VOMS fully qualified attribute names (FQANs) to describe the users and do not support 'deny'ing of rights. Also, general regular expressions cannot be used in the argument, only the `".*"` wild card is supported. An example of the policy file follows:

```
resource "LB" {
  action "ADMIN_ACCESS" {
    rule permit {
      subject = "/DC=cz/DC=cesnet-ca/O=CESNET/CN=Daniel Kouril"
    }
    rule permit {
      fqan = "/vo/Role=Manager"
    }
  }
  action "STATUS_FOR_MONITORING" {
    rule permit {
      fqan = "/vo/monitoring"
    }
    rule permit {
      fqan = "/TEST/rtm"
    }
  }
  action "LOG_WMS_EVENTS" {
    rule permit {
      subject = "/DC=cz/DC=cesnet-ca/O=CESNET/CN=wms01.cesnet.cz"
    }
  }
}
```

<sup>8</sup><https://twiki.cern.ch/twiki/bin/view/EGEE/SimplifiedPolicyLanguage>

```
}  
action "LOG_GENERAL_EVENTS" {  
    rule permit {  
        subject = ".*"  
    }  
}  
action "REGISTER_JOBS" {  
    rule permit {  
        subject = ".*"  
    }  
}  
action "GRANT_OWNERSHIP" {  
    rule permit {  
        fqan = "/VO/ Pilot_job_factory"  
    }  
}  
}
```

After changing the file, the server has to be restarted.

In order to provide yet additional level of authorization, the LCAS schema<sup>[5]</sup> is still used in the server. If enabled in configuration, it can be used to specify more general settings using the standard LCAS modules and e.g. blacklist particular users. The standard LCAS documentation should be followed to set up the LCAS layer properly.

### 2.3.10 EXPORT TO R-GMA

*Note: This feature is now obsolete and only available in **L&B version 1.x**.*

L&B server can export information on job state changes to R-GMA infrastructure through `lcgmon` in real time. This export is enabled by YAIM by default and uses `GLITE_WMS_LCGMON_FILE` environmental variable to retrieve name of log file which is to be consumed by `lcgmon` (usually `/var/glite/logging/status.log`). The log file has to be rotated regularly.

### 2.3.11 DATA BACKUP

Data stored L&B server can be backed up using backups of underlying database or using `glite-lb-dump` utility. The latter has some advantages, see Section 3.2.4 for details.

### 2.3.12 PURGING OLD DATA

Initial YAIM configuration creates a cron job which runs once a day and purges old data (jobs in Cleared state after two days, Aborted and Cancelled after 15 days, and other jobs after 60 days of inactivity). It is recommended to run the cron jobs more often (in order to purge less jobs during single run) if event queue backlogs form in client WMS machines when the purging cron jobs is running. For details on setting job purge timeouts, see Sect. 3.2.5

### 2.3.13 EXPLOITING PARALLELISM

L&B server uses 10 worker processes (threads) to handle active client accesses (inactive connections are killed when necessary). Each worker process uses separate connection to database server. Number of worker processes can be changed by adding `--slaves` parameter with desired number to servers command line using `GLITE_LB_SERVER_OTHER_OPTIONS` variable.

### 2.3.14 TUNING DATABASE ENGINE

In order to achieve high performance with LB server underlaying MySQL database server has to be configured reasonably well too. Default values of some MySQL settings are likely to be suboptimal and need tuning, especially for larger machines. These are MySQL configuration variables (to be configured in `[mysqld]` section of `/etc/my.cnf`) that need tuning most often:

- `innodb_buffer_pool_size` – size of database memory pool/cache. It is generally recommended to set it to around 75% of RAM size (32bit OS/MySQL versions limit this to approx. 2GB due to address space constraints).
- `innodb_flush_log_at_trx_commit` – frequency of flushing to disk. Recommended values include:
  - 1 (default) – flush at each write transaction commit; relatively slow without battery-backed disk cache but offers highest level of data safety
  - 0 – flush once per second; fast, use if loss of latest updates on MySQL or OS crash (e.g. unhandled power outage) is acceptable (database remains consistent)
- `innodb_log_file_size` – size of database log file. Larger values save some I/O activity, but also make database shutdown and crash recovery slower. Recommended value: 50MB. Clean `mysqld` shutdown and deletion of log files (`/var/lib/mysql/ib_logfile*` by default) is necessary before change.
- `innodb_data_file_path` – path to main database file. File on disk separate from OS and MySQL log files (`innodb_log_group_home_dir` variable, `/var/lib/mysql/` by default) is recommended.

## 2.4 L&B PROXY

All necessary configuration of L&B proxy is done by YAIM, and described with gLite WMS installation elsewhere. Previous L&B server section applies to merged server+proxy setups (since **L&B version 2.0**).

A special care must be taken when an existing L&B proxy database is migrated to **L&B version 2.x**. In general, this is not a typical scenario – **L&B version 2.x** server in proxy mode on WMS node is introduced with a major WMS upgrade, and it is expected to be installed from scratch rather than migrated, preserving L&B proxy data.

If the migration is really needed, `glite-lb-migrate_db2version20` script should be run with `-p` (Sect 2.3.5). However, the L&B database name remains `lbproxy` while the **L&B version 2.x** binaries expect unified `lbserver20` by default. Because renaming a MySQL database is a non-trivial, error prone task, the recommended workaround is to add the following variable setting

```
GLITE_LB_SERVER_OTHER_OPTIONS="--mysql lbserver/@localhost:lbproxy"
```

into the gLite startup environment (`/opt/glite/etc/profile.d`) instead. This setting makes L&B server use the `lbproxy` database instead of the default.

## 2.5 L&B LOGGER

All necessary configuration of normal L&B logger is done by YAIM.

## 2.6 L&B HARVESTER

L&B Harvester gathers information about jobs from L&B servers using efficient L&B notification mechanism. It manages notifications and keeps them in a persistent storage (file or database table) to reuse later on next launch. It takes care of refreshing notifications and queries L&B servers back when some notification expires.

It is not intended for normal usage. You will need Harvester for sending notifications to MSG publish infrastructure, but only for older L&B server releases (gLite 3.1.x, support since gLite 3.1.19). Example of YAIM configuration:

```
GLITE_LB_HARVESTER_ENABLED=true
GLITE_LB_HARVESTER_MSG_OPTIONS="--wlcg-topic=org.wlcg.usage.jobStatus
--wlcg-config=/etc/msg-publish/msg-publish.conf
--wlcg-binary=/usr/bin/msg-publish"
```

Since gLite 3.2.1, L&B Harvester is not needed for MSG publish. Notifications are reliably delivered by interlogger instead. Delivery can be switched on in this case by `GLITE_LB_MSG_ENABLED` and `GLITE_LB_MSG_BROKER` YAIM options.

## 2.7 SMOKE TESTS

Thorough tests of L&B, including performance measurement, are covered in the L&B Test Plan document [4]. This section describes only elementary tests that verify basic functionality of the services.

The following test description assumes the L&B services installed and started as described above.

### 2.7.1 JOB REGISTRATION

Register a new job with the L&B server and check that its status is reported correctly.

**Prerequisites:** Installed `glite-lb-client` package, valid user's X509 credentials, known destination (address:port) of running L&B server. Can be invoked from any machine. <sup>9</sup>

**How to run:**

```
glite-lb-job_reg -m server_name:port
```

A new jobid is generated and printed. Run

```
glite-lb-job_status jobid
```

**Expected result:** The command should report "Submitted" job status.

---

<sup>9</sup>Example scripts or binaries used here can be found either in `/opt/glite/examples` (if installed from a LB node repository) or `/usr/lib64/glite-lb/examples` (if installed from the EMI repository).

### 2.7.2 LOGGING EVENTS VIA LB-LOGGER

Send several L&B events, simulating normal job life cycle. Checks delivery of the events via L&B logger.

**Prerequisites:** Installed glite-lb-client package, valid user's X509 credentials, known destination (address:port) of running L&B server. Must be run on a machine where glite-lb-logger package is set up and running.<sup>9</sup>

**How to run:**

```
glite-lb-running.sh -m server_name:port
```

The command prints a new jobid, followed by diagnostic messages as the events are logged. Check the status of the new job with

```
glite-lb-job_status jobid
```

**Expected result:** Due to asynchronous event delivery various job states can be reported for limited time (several seconds). Finally the "Running" status should be reached.

### 2.7.3 LOGGING EVENTS VIA LB-PROXY

Send events via L&B proxy. Checks the proxy functionality.

**Prerequisites:** Running L&B proxy, in standalone package for **L&B version 1.x**, or L&B server running with **-P** (proxy only) or **-B** (both server and proxy) options.<sup>9</sup>

**How to run:** Similar to Sect. 2.7.2:

```
glite-lb-running.sh -x -m server_name:port
```

And check with:

```
glite-lb-job_status -x /tmp/lb_proxy_server.sock jobid
```

### 2.7.4 NOTIFICATION DELIVERY

Register for receiving notifications, and log events which trigger the notification delivery. Checks the whole notification mechanism. The test is quite complex, though, see notification sections in [2, 4].



## 3 MAINTENANCE

### 3.1 CHANGING DEFAULT SETTINGS

All configurable settings of the L&B daemons (network and local sockets, file paths, modified behaviour etc.) can be set with specific command line options. In the following only relevant options are discussed whenever appropriate. See specific manual pages for complete reference.

### 3.2 L&B SERVER AND PROXY

This section deals with several typical but more peculiar tasks that need more verbose description. It is complemented with the full commands reference that is provided as standard manual pages installed with the L&B packages.

#### 3.2.1 STANDARD AND DEBUG LOGS

In normal operation L&B server sends error messages to syslog. Informational messages are generally avoided in order to prevent syslog congestion.

Since **L&B version 2.1**, the service implements a common logging format<sup>10</sup> (see section 2.2, page 13).

For **L&B version 2.0 and lower**, the following instructions apply:

When tracing problems, `GLITE_LB_SERVER_DEBUG` environment variable can be set to non-empty value when starting the service. Then verbose log `$GLITE_LOCATION_VAR/lb.log` (as well as `$GLITE_LOCATION_VAR/notif-il.log` eventually when notifications are enabled). Beware that these can grow huge easily.

**L&B version 1.x only:** not available for L&B proxy, `-d` and output redirection have to be added manually if necessary.

#### 3.2.2 CHANGING INDEX CONFIGURATION

L&B server only (L&B proxy database is neither so huge nor accessed directly by users).

Inefficient queries, yielding full scan of L&B database tables (up to millions of tuples) would degrade server performance. Therefore L&B does not allow arbitrary queries in general (server option `--no-index` can change this behaviour). On the contrary, a query has to hit a *job index*, build on one or more job attributes. It is left up to the specific L&B server administrator to decide which job attributes are selective enough to be indexed and allow queries (for example for many-users communities job owner can be a sufficient criterion; for others, where only a few users submit thousands of jobs, it is not).

Technically, job indices are implemented via dedicated columns in a database table. These columns and their indices are scanned by the L&B server on startup, therefore there is no specific configuration file. Changing the index configuration is rather heavyweight operation (depending on the number of jobs in the database), it performs updates of all tuples in general, and it should be done when the server is not running.

Indices are manipulated with a standalone utility `glite-lb-bkindex` (see its man page for complete usage reference). A general sequence of changing the indices is:

---

<sup>10</sup>[http://en.wikipedia.org/wiki/Common\\_Log\\_Format](http://en.wikipedia.org/wiki/Common_Log_Format)

<i>IndexType</i>	<i>IndexName</i>	description
system	owner	job owner
	destination	where the job is heading to (computing element name)
	location	where is the job being processed
	network_server	endpoint of WMS
	stateEnterTime	time when current status was entered
	lastUpdateTime	last time when the job status was updated
time	<i>state name</i>	when the job entered given state (Waiting, Ready, ...)
user	<i>arbitrary</i>	arbitrary user tag

**Table 1:** Available index column types and names

1. stop the running server
2. retrieve current index configuration

```
glite-lb-bkindex -d >index_file
```

3. edit `index_file` appropriately
4. re-index the database (it may take long time)

```
glite-lb-bkindex -r -v index_file
```

`-r` stands for “really do it”, `-v` is “be verbose”

5. start the server again

The index description file follows the classad format, having the following grammar:

```
IndexFile ::= [ JobIndices = { IndexList } ]
IndexList ::= IndexDef | IndexDef, IndexList
IndexDef ::= IndexColumn | ComplexIndex
IndexColumn ::= [ type = "IndexType"; name = "IndexName" ]
ComplexIndex ::= { ColumnList }
ColumnList ::= IndexColumn | IndexColumn, ColumnList
```

where eligible *IndexType*, *IndexName* combinations are given in Tab. 1. A template index configuration, containing indices on the most frequently used attributes, can be found in `[/opt/glite]/etc/glite-lb-index.conf.template`.<sup>11</sup>

### 3.2.3 MULTIPLE INSTANCES

Specific conditions (for example debugging, different authorization setup, ...) may require running multiple L&B server instances on the same machine. Such setup is available, however, there is no specific support in automated configuration, the additional non-default server instances must be run manually.

The other server instance must use different ports (changed with `-p` and `-w` options), as well as use different pid file (`-i` option).

<sup>11</sup>The `/opt/glite` prefix applies for L&B installed from gLite's L&B node repository.

The servers may or may not share the database (non-default is specified with `-m`)<sup>12</sup>.

Though it may have little sense to run multiple L&B proxy instances, it is possible too. Non-default listening socket have to be specified via `-p` option then.

### 3.2.4 BACKUP DUMPS

L&B server only, not supported by proxy.

(This functionality should not be confused with per-job dumps, Sect. 2.3.12 and 3.2.5)

Besides setting up L&B server database on a reliable storage or backing it up directly (Sect. 2.3.11) L&B server supports backing up only incremental changes in the data. Advantages of this approach are lower volume of data to be backed up, and possibility to load them to another instance (for example for heavyweight queries which should not disturb normal operation), disadvantage is a more complex and more fragile setup.

Using an external utility `glite-lb-dump` (typical invocation is with a single option `-m server_name:port`, see man page for details) the server is triggered to dump events, which arrived in a specified time interval, into a text file. (Default interval is from last dump till the current time.)

`glite-lb-dump` is a standalone client program, however, the events are stored at server side (that is not transferred to the client, due to performance reasons), in a uniquely named text file prefixed with the value of `-D` server option. This kind of dump contains events according to their arrival time, regardless of jobs they belong to.

It is sufficient to run the dump regularly (from a cron job), with a frequency matching an acceptable risk of losing data (several hours typically), and back up the resulting dump files.

In the event of server crash, its database should be recreated empty, and the server started up. Then the dump files can be loaded back with complementary `glite-lb-load` utility.

Server privileges granting `ADMIN_ACCESS` (see section 2.3.9) are required to run `glite-lb-dump` and `glite-lb-load`. Dumping the events does not interfere with normal server operation.

This backup strategy can interfere with too aggressive setting of old data purging (Sect. 3.2.5). If the purging grace period is shorter than the dump interval, events may get purged before they are captured by the backup dump. However, this interference is unlikely (reasonable purge grace period is several times longer than dump period), and it is not fatal in general (data were purged on purpose either).

### 3.2.5 PURGING AND PROCESSING OLD DATA

Primary purpose of the LB purge operation is removal of aged data from LB database. This is necessary in production in order to prevent ever-increasing database and sustain reasonable performance of the server. Therefore the purge should be invoked periodically.

The purge operation has additional important "side effect" – dumping the purged data into a plain text file. These dumps can be archived "as is" or uploaded to Job Provenance.

**Purge setup** The purge operation itself is performed by a running L&B server (there is no need to shut it down, then). However, it is triggered with `glite-lb-purge` client command (complete usage reference

---

<sup>12</sup>Even when sharing the database, the servers are still partially isolated from one another, for example a job <https://my.machine:9000/xyz> cannot be queried as <https://my.machine:8000/xyz>. However, due to implementation internals, the second job cannot be registered.

is given in its man page). A typical invocation specifies L&B server to purge (`-m` option), and purge timeouts (grace periods) for several job states – options `-a` (aborted), `-n` (canceled), `-c` (cleared), and `-o` (other). L&B versions 2.0 supports also `-e` (done) option. A job falling in one of the four categories is purged when it has not been touched (that is an event arrived) for time longer than the specified category timeout. Suggested values are several days for aborted and canceled jobs, and one day for cleared jobs, however, the values may strongly vary with L&B server policy.

Optionally, `-s` purge command option instructs the server to dump the purged data into a file at the server side. Its location (prefix) is given by `-S` server option, the purge command reports a specific file name on its output.

It is recommended (and the default YAIM setup does so, via the `glite-lb-export.sh` wrapper) to run the purge command periodically from cron.

Server privileges granting `ADMIN_ACCESS` (see section 2.3.9) are required to run `glite-lb-purge`.

If the server database has already grown huge, the purge operation can take rather long and hit the L&B server operation timeout. At client side, that is the `glite-lb-purge` command, it can be increased by setting `GLITE_WMS_QUERY_TIMEOUT` environment variable. Sometimes hardcoded server-side timeout can be still reached. In either case the server fails to return a correct response to the client but the purge is done anyway.

L&B proxy purges jobs automatically when they reach a state ensuring that WMS will neither query nor log events to them anymore. Therefore routine purging is not required theoretically. However, frozen jobs which never reach such a state may occur in an unstable environment, and they may cumulate in L&B proxy database for ever. Therefore occasional purging is recommended too. **L&B version 2.x** supports `-x` option of `glite-lb-purge`, allowing to purge L&B proxy database too. With **L&B version 1.x** the emergency purge procedure described bellow is the only option.

**Emergency purge** When regular purge was not invoked for some time, it may happen that the database grows huge and the regular (on-line) purge fails. In order to work around such situation we provide an off-line emergency purge script `glite-lb-bkpurge-offline.sh`

The script accepts the same `-acno` options, and adds `-d` for “done” jobs. Via `-p` also L&B proxy database can be purged (all L&B versions).

On startup, a warning message is printed and interactive confirmation requested. Re-check that L&B server (proxy) is not running, and carry on only when you know what you are doing.

**Post-mortem statistics** Once a job is purged from the database, all important data about the job can be processed offline from the corresponding dump file. The idea of post-mortem statistics is the following:

- LB server produces dump files (during each purge on regular basis), see LB server startup script; option `-D` / `--dump-prefix` of `glite-lb-bkserverd`,
- these dumps are exported for the purposes of JP also on regular basis, see LB/JP deployment module; option `-s` / `--store` of `glite-lb-lb_dump_exporter`,
- it depends on the LB server policy if dumps in this directory are used for the statistics purposes or all files are hardlinked for example to a different directory
- general idea is such that data are available for statistics server that downloads and removes dumps after download! Dump files are then processed on the statistics server.

What needs to be done on the LB server:

- `glite-lb-bkserverd` and `glite-lb-lb_dump_exporter` running
- `gridftp` running (allowing statistics server to download and remove files from a given directory)

What needs to be done on the statistics server:

- `glite-lb-utils` package installed
- download and remove files from the LB server see `glite-lb-statistics-gsi.sh` (shell script in the examples directory)
- process dump files using the `glite-lb-statistics` tool see `glite-lb-statistics.sh` (shell script in the examples directory)

all scripts are supposed to be run from a crontab.

**Export to Job Provenance** An important, though currently optional, processing of L&B dumps is their upload to the Job Provenance service for permanent preservation of the data.

When enabled (via configuration environment variables, see below), the export is done in two steps:

- `glite-lb-export.sh` wrapper script, after calling `glite-lb-purge`, breaks up the resulting dump file on a per-job basis. The individual job dump files are stored in a dedicated spool directory.
- `glite-jp-importer` daemon (installed optionally in `glite-jp-client.rpm`) checks the spool directory periodically, and tries to upload the files into JP.

Details, including the configuration variables, are covered at the following wiki page:

[http://egee.cesnet.cz/mediawiki/index.php/LB\\_purge\\_and\\_export\\_to\\_JP](http://egee.cesnet.cz/mediawiki/index.php/LB_purge_and_export_to_JP).

**Persistent Information on Purged Jobs (“Zombies”)** Since L&B version 2.0, the JobID of a purged job is not fully discarded but stored in a separate table. A query for the status of such job returns *Identifier removed*.<sup>13</sup> There are two purposes to this arrangement:

1. Confirming that the job existed and that its details have been exported to Job Provenance (if deployed and configured),
2. Preventing reuse of the JobID in case of flag `EDG_WLL_LOGLFLAG_EXCL` set on registration.

### 3.2.6 ON-LINE MONITORING AND STATISTICS

**CE reputability rank** Rather frequent problem in the grid production are “black hole” sites (Computing Elements). Such a site declares itself to have an empty queue, therefore schedulers usually prefer sending jobs there. The site accepts the job but it fails there immediately. In this way large number of jobs can be swallowed, affecting the overall success rate (namely for non-resubmittable jobs).

<sup>13</sup>As of L&B 3.0. In 2.0 and 2.1 releases an empty job status structure was returned with job state set to *purged*.

L&B data as a whole contain enough information to detect such sites. However, due to the primary per-job structure certain reorganization is required.

A job is always assigned to a *group* according to the CE where it is executed (cf. "destination" job state attribute). Similarly to RRDtool<sup>14</sup> for each recently active group (CE), and for each job state (Ready, Scheduled, Running, Done/OK, Done/Failed), a fixed sized series of counters is maintained. At time  $t$ , the counters cover intervals  $[t - T, t]$ ,  $[t - 2T, t - T]$ , ... where  $T$  a fixed interval size. Whenever a job state changes, the series matching the group and new state is shifted eventually (dropping its expired tail), and the current counter is incremented. In addition, multiple series for different  $T$  values (that is covering different total times) are available.

The data are available via statistics calls of the client API, see `statistics.h` for details (coming with glite-lb-client in **L&B version 2.x**, glite-lb-client-interface in **L&B version 1.x**). The call specifies the group and job state of interest, as well as queried time interval. The interval is fitted to the running counter series as accurately as possible, and the average number of jobs per second which entered the specific state for the given group is computed. The resolution ( $T$ ) of the used counters is also returned.

In gLite 3.1 WMS the calls can be accessed from inside of the matchmaking process via `successFraction(CEId)` JDL function. The function computes the ratio of successful vs. all jobs for a given CE, and it can be directly used to penalize detected black hole CEs in the ranking JDL expression.

The functionality is enabled with `--count-statistics` L&B server option (disabled by default).

The gathered information is currently not persistent, it is lost when the server is stopped. Despite the statistics call API is defined in a general way, the implementation is restricted to a hardcoded configuration of a single grouping criterion (the destination), and a fixed set of counter series (60 counters of  $T = 10s$ , 30 of 1 minute, and 12 of 15 minutes). The functionality has not been very thoroughly tested yet.

**glite-lb-mon** is a program for monitoring the number of jobs on the LB server and their several statistics. It is part of the `glite-lb-utils` package, so the monitoring can be done from remote machine where this package is installed and the environment variable `GLITE_WMS_QUERY_SERVER` properly set. Values like minimum, average and maximum time spent in the system are calculated for jobs that entered the final state (Aborted, Cleared, Cancelled) in specific time (default last hour). Also number of jobs that entered the system during this time is calculated.

A special bkindex configuration is needed. The following time indices must be defined:

```
[ type = "time"; name = "submitted" ],  
[ type = "time"; name = "cleared" ],  
[ type = "time"; name = "aborted" ],  
[ type = "time"; name = "cancelled" ],
```

For more details see man page `glite-lb-mon(1)`.

**glite-lb-mon-db** is a low-level program for monitoring the the number of jobs in the LB system. Using the LB internals, it connects directly to the underlying MySQL database and reads the number of jobs in each state. The tool is distributed together with the server in the `glite-lb-server` package. It can be used to read data also from the database of LB Proxy. For more details see man page `glite-lb-mon-db(1)`.

---

<sup>14</sup><http://oss.oetiker.ch/rrdtool/>

**Subjob states in a collection** can be calculated on demand on the server and returned as a histogram using standard job status query. There are two ways how to obtain the histogram:

- fast histograms, the last known states are returned, see e.g.

```
glite-lb-job_status -fasthist <collectionJobId>
```

- full histograms. the states of all collection subjobs are recalculated, see o.g.

```
glite-lb-job_status -fullhist <collectionJobId>
```

The command `glite-lb-job_status` is a low level query program that can be found in the package `glite-lb-client` among examples.

### 3.3 L&B LOGGER

The logger component (implemented by `glite-lb-interlogd` daemon fed by either `glite-lb-logd` or L&B proxy) is responsible for the store-and-forward event delivery in L&B (Sect 1.1.2). Therefore eventual operational problems are related mostly to cumulating undelivered events.

#### 3.3.1 EVENT FILES

L&B logger stores events in one file per job, named `$GLITE_LOCATION_VAR/log/dglogd.log.JOBID` by default (JOBID is only the part after the L&B server address prefix). The format is text (ULM [6]), one event per line. In addition, control information on delivery status is stored in additional file with `.ctl` suffix.

In case of emergency (for example corrupted file) the files can be examined with `glite-lb-parse_eventsfile`.<sup>15</sup> It is possible to hand-edit the event files in emergency (remove corrupted lines). However, `glite-lb-interlogd` must not be running, and the corresponding `.ctl` file must be removed.

#### 3.3.2 BACKLOG REASONS

**Undeliverable jobid.** In normal gLite job processing, jobids are verified on job submission (via synchronous job registration, see [2]), hence occurrence of undeliverable jobid (that is its prefix does not point to a working L&B server) is unlikely. On the other hand, if it happens, and an event with such a jobid is logged, for example due to a third-party job processing software bug, `glite-lb-interlogd` keeps trying to deliver it indefinitely.<sup>16</sup> The unsuccessful attempts are reported via syslog. The only solution is manual removal of the corresponding files and restart of the service.

**Corrupted event file.** For various reasons the files may get corrupted. In general, corrupted file is detected by `glite-lb-interlogd`, and it is moved to *quarantine* (by renaming the file to contain “quarantine” in its name). The action is reported in syslog. The renamed files can be removed or repaired by hand and renamed back for `glite-lb-interlogd` to pick them up again (in this case, the service needn't be stopped).

<sup>15</sup>Not fully supported tool, installed by `glite-lb-client` package among examples.

<sup>16</sup>Unless event expiration is set, though it is not done for normal events.



**Slow delivery.** Either glite-lb-interlogd or the target L&B server(s) may not keep pace with the incoming stream of events. Amount of such backlog can be quickly assessed by looking at timestamps of the oldest event files (with `ls -lt $GLITE_LOCATION_VAR/log/`, for example). Fully processed event files are deleted in approx. one minute intervals, so files last modified one minute ago or newer do not constitute backlog at all. Unless the backlog situation is permanent, no specific action is required—the event backlog decreases once the source is drained. Otherwise hardware bottlenecks (CPU, disk, network) have to be identified (with standard OS monitoring) and removed, see sections 2.3.1 and 2.3.14 for L&B server tuning tips. In order to maximize performance of `glite-lb-interlogd`, it is recommended to distribute jobs to multiple L&B servers that can be shared by multiple WMS setups in turn.

### 3.3.3 NOTIFICATION DELIVERY

When `glite-lb-logger` package is optionally installed with L&B server, a modified `glite-lb-notif-interlogd` is run by the server startup script. This version of the daemon is specialized for L&B notifications delivery; it uses the same mechanism, however, the events (notifications) are routed by *notification id* rather than *jobid*, and targeted to user's listeners, not L&B servers. See [2, 3] for details.

On the contrary to normal events, it is more likely that the event destination disappears permanently. Therefore the notification events have their expiration time set, and `glite-lb-interlogd` purges expired undelivered notifications by default. Therefore the need for manual purge is even less likely.

The event files have different prefix (`/var/tmp/glite-lb-notif` by default).

### 3.3.4 DEBUG MODE

All the logger daemons, that is `glite-lb-logd`, `glite-lb-interlogd`, and `glite-lb-notif-interlogd`, can be started with `-d` to avoid detaching control terminal, and `-v` to increase debug message verbosity. See manual pages for details.

## 3.4 USED RESOURCES

### 3.4.1 SERVER AND PROXY

**Processes and threads.** By default L&B server runs as one master process, and 10 slave processes.

Threads are not used.

Master pid is stored in the file `$HOME/edg-bkserverd.pid` (**L&B version 1.x**) or in the file `$HOME/glite-lb-bkserverd.pid` (**L&B version 2.x**) respectively. Number of slaves can be set with `-s`, `--slaves` option, pid file location with `-i`, `--pidfile`.

Slave server processes are restarted regularly in order to prevent memory leakage.

**Network and UNIX ports.** L&B server listens on port 9000 for incoming queries, 9001 for incoming events, and 9003 for WS interface queries. The former two can be changed with `-p`, `--port` option (incoming events are always one port higher than queries), the latter with `-w`, `--wsport`.

When L&B notifications are enabled (automatically in server startup script when `glite-lb-logger` is also installed with the server), `/tmp/glite-lb-notif.sock` is used for communication with notification interlogger. It can be changed with `--notif-il-sock`.



L&B proxy communicates on two UNIX sockets: `/tmp/lb_proxy_server.sock` (queries) and `/tmp/lb_proxy_store.sock` (incoming events).

**L&B version 2.x:** as proxy and server are merged, the mode of operation determines the actual server ports (network or UNIX). Notifications are not delivered from proxy-only mode.

**IPC. L&B version 1.x only:** mutual exclusion of server slaves is done via SYSV semaphores. The semset key is obtained by calling `ftok(3)` on the pid file.

**L&B version 2.x** does not use semaphores anymore.

**Database.** Server data are stored in MySQL database. Normal setup assumes “server scenario”, that is the database IS NOT protected with password and it is not accessible over network. Option `-m user/password@machine:database` can be used to change the default connect string `lbserver/@localhost:lbserver20`.

**Dump files.** Backup dump (Sect. 3.2.4) files are stored in `/tmp/dump` by default, however, gLite startup script uses `--dump-prefix` option to relocate them into `$GLITE_LOCATION_VAR/dump`. Similarly dumps resulting from purge operation (Sect. 3.2.5) go to `/tmp/purge` by default, and startup script uses `--purge-prefix` to set `$GLITE_LOCATION_VAR/purge`.

When export to Job Provenance is enabled (Sect. 3.2.5) job registrations are exported to directory `/tmp/jpreg` overridden in startup script with `--jpreg-dir` to `$GLITE_LOCATION_VAR/jpreg`. In addition, dump files are further processed in `GLITE_LOCATION_VAR/jpdump`.

Normal operation of JP export should not leave any files behind, however, abnormal situations (JP unavailable for longer time etc.) may start filling the disk space. Therefore periodic checks are recommended.

Also when purge and dump operations are invoked resulting files are left in their directories and it is the responsibility of the operation caller to clean them up.

**Notifications** (see above with Network ports) are stored in files with `/var/tmp/glite-lb-notif` prefix. It can be changed with `--notif-il-fprefix`.

Normal operation of notification interlogger purges these files when they are either delivered or expire.

### 3.4.2 LOGGER

**Processes and threads.** `glite-lb-logd` uses one permanent process and forks a child to serve each incoming connection.

`glite-lb-interlogd` (and derived `glite-lb-notif-interlogd`) is multithreaded: one thread for handling input, another for recovery, and a dynamic pool of others to deliver events to their destinations.

**Network and UNIX ports.** `glite-lb-logd` listens on 9002 (can be changed with `-p` option). The daemons communicate over UNIX socket `/tmp/interlogger.sock`, can be changed with `-s` (for both daemons simultaneously).

**Event files** are stored with prefix `/var/glite/log/dglogd.log` (changed with `-f`).

## **4 TROUBLESHOOTING**

**TODO:** *gather information and finish the section*

## 5 FAQ—FREQUENTLY ASKED QUESTIONS

### 5.1 JOB IN STATE 'RUNNING' DESPITE HAVING RECEIVED THE 'DONE' EVENT FROM LRMS

Jobs stay in state *Running* until a *Done* event is received from the workload management system. *Done* events from local resource managers are not enough since the job in question may have been resubmitted in the meantime.

### 5.2 WMS CANNOT PURGE JOBS OR PERFORM OTHER PRIVILEGED TASKS

In short, WMS has not been given adequate permissions when configuring the L&B server. You need to modify your configuration and restart the server:

#### 5.2.1 FOR L&B VERSION 3.0.11 OR HIGHER, USING YAIM

Modify your `siteinfo.def`, specifying the DN of your WMS server in YAIM parameter `GLITE_LB_WMS_DN`; for instance:

```
GLITE_LB_WMS_DN=/DC=cz/DC=cesnet-ca/O=CESNET/CN=wms01.cesnet.cz
```

Then rerun YAIM: `/opt/glite/yaim/bin/yaim -c -s site-info.def -n glite-LB`

This will give your WMS exactly the right permissions to carry out all required operations.

#### 5.2.2 FOR ALL VERSIONS OF L&B, USING YAIM

Modify your `siteinfo.def`, specifying the DN of your WMS server in YAIM parameter `GLITE_LB_SUPER_USERS`; for instance:

```
GLITE_LB_SUPER_USERS=/DC=cz/DC=cesnet-ca/O=CESNET/CN=wms01.cesnet.cz
```

Then rerun YAIM: `/opt/glite/yaim/bin/yaim -c -s site-info.def -n glite-LB`

This will give your WMS adequate rights to perform its operations and requests (running purge, querying for statistics, etc.) but it will also grant it additional administrator rights (such as granting job ownership). On newer installations, the method explained in section [5.2.1](#) is preferable.

#### 5.2.3 FOR L&B VERSION 2.1 OR HIGHER, WITHOUT YAIM

L&B's authorization settings can be found in file `[/opt/glite]/etc/glite-lb/glite-lb-authz.conf`. Permit actions `PURGE`, `READ_ALL` and `GET_STATISTICS` for your WMS and restart the L&B server. This will lead to results equivalent to [5.2.1](#). For instance, change the adequate sections in `glite-lb-authz.conf` to:

```
action "READ_ALL" {
    rule permit {
        subject = "/DC=cz/DC=cesnet-ca/O=CESNET/CN=wms01.cesnet.cz"
```

```
        }  
    }  
  
    action "PURGE" {  
        rule permit {  
            subject = "/DC=cz/DC=cesnet-ca/O=CESNET/CN=wms01.cesnet.cz"  
        }  
    }  
  
    action "GET_STATISTICS" {  
        rule permit {  
            subject = "/DC=cz/DC=cesnet-ca/O=CESNET/CN=wms01.cesnet.cz"  
        }  
    }  
}
```

## REFERENCES

- [1] E. Laure, F. Hemmer, F. Prelz, S. Beco, S. Fisher, M. Livny, L. Guy, M. Barroso, P. Buncic, P. Kunszt, A. Di Meglio, A. Aimar, A. Edlund, D. Groep, F. Pacini, M. Sgaravatto, and O. Mulmo. Middleware for the next generation grid infrastructure. In *Computing in High Energy Physics and Nuclear Physics (CHEP 2004)*, 2004.
- [2] A. Křenek et al. L&B User's Guide. <http://egee.cesnet.cz/en/JRA1/LB/>.
- [3] A. Křenek et al. L&B Developer's Guide. <http://egee.cesnet.cz/en/JRA1/LB/>.
- [4] A. Křenek et al. L&B Test Plan. <http://egee.cesnet.cz/en/JRA1/LB/>.
- [5] Site Authorisation and Enforcement Services: LCAS, LCMAPS, and gLExec. <https://www.nikhef.nl/pub/projects/grid/gridwiki/index.php/LCAS>.
- [6] J. Abela and T. Debeaupuis. Universal format for logger messages. IETF Internet Draft, 1997.
- [7] G. Avellino et al. The DataGrid Workload Management System: Challenges and Results. *Journal of Grid Computing*, 2(4):353–367, Dec 2004.