

GemsTracker Hosting

Introduction

This document explains the requirements to install GemsTracker on a standard Linux or Windows Webserver.

About Gemstracker

GemsTracker, the Generic Medical Survey Tracker, is an open source web server application available from www.gemstracker.org and hosted by SourceForge. On the GT site you will find more information and examples about how GT can be used in research and clinical care. Moreover, you find more user and developer documentation on our wiki (<http://gemstracker.org/wiki/doku.php?id=start>), which also contains a quick start for setting up GT sites.

GT consists of a core library with a default installation, but is designed from the ground up to enable developers to extend the core functionality for specific projects. There are currently more than 30 different projects in production, each comprising a different website. The projects range from having only the most minimal project code of maybe 20 lines of code; up to the Pulse and ZSD projects that almost rival the GT library in their project specific code-count.

This document describes the requirements for webserver hosting one or more GT installations. It addresses security and explains the GT architecture in order to explain the different deployment scenarios. The document then gradually moves on to more detailed configuration information.

GemsTracker and Security

GT is an application for collecting survey answers containing personal patient data using the web. This is privacy sensitive information and therefore security is an important consideration. In most countries one is required by law to make a serious effort using the best industry practices to keep the information inaccessible to people without a contractual obligation to respect the privacy of the patients.

In the Netherlands the best practices are documented in the NEN 7510 standard. NEN 7510 certification requires that the hosting company should at the very least be ISO 127.0001 certified, therefore a hosting account on a shared webserver is usually insufficient for certification purposes.

The decision to NEN 7510 certify a particular project is up to the owners of that individual project, but in our opinion the default hosting and application environment should pose no obstacles to certification.

GemsTracker Architecture

GT is a library that enables health care staff to make sure the right questions about the right patients are asked at the right time. What GT does not do is asking those questions. For survey creation and entry GT uses third party products through a pluggable architecture. This gives hosting companies the freedom to also use their own (proprietary) questionnaire systems and benefit from the GT software functionality.

GemsTracker Hosting

Currently GT has plugins for the LimeSurvey and OpenRosa and at least one project has a project specific survey plugin. LimeSurvey is the most used survey plugin, but where LimeSurvey is mentioned in the rest of this document the intention is that LimeSurvey can be replaced by another survey system.

Patient identifying information is stored by GT, together with an anonymous random patient identifier that is hidden from the users in the GT interface. This identifier is also stored in the questionnaire software (e.g. LimeSurvey and OpenRosa) with the answers belonging to the patient. Preferably all treatment data is stored in these survey sources in databases separate from the database used by GT. This has the advantage that during the treatment and/or the research GT can show the treatment data for each patient, but during group analysis and evaluation the data cannot be traced to individual patients.

GT installations can also create a bridge table for data warehousing containing all the answers from all the sources in a format easy to use by data mining packages. Again this table can and should be stored in a separate database in order to keep treatment data and identifying information separate.

GemsTracker application environments

Testing or developing a GT application usually requires different project settings than during production. Together these different settings form the application environment of an installation. In the Zend Framework¹ this environment is primarily determined using an environment variable in the underlying operating system.

The application environment determines usage restrictions, caching, folder locations and utility application. E.g. all GT environments except the production environment bounce the e-mails generated by the GT to the sender instead of the test recipient. The reason for this is that all these environments are used to try out the software. When sending an e-mail the user does not want to actually send the e-mail to that address, but instead wants to test how the e-mail looks for the recipient.

These differences in behavior are why GT always tells a user in which environment her or she is working.

¹ The Zend Framework is a PHP framework for web development maintained by the Zend corporation that also maintains PHP and is the framework that was used for the development of GT.

GemsTracker Hosting

To set these modes the Zend Framework uses the operating system `APPLICATION_ENV` environment variable. GT projects can use these application environments:

<code>APPLICATION_ENV</code>	Description
<code>production</code>	<i>Default.</i> For production development, highly cached and high security.
<code>Demo</code>	A demonstration/learning installation, lower security and e-mail bouncing.
<code>acceptance</code>	Identical to <code>production</code> except for e-mail bouncing, but shows that it is a different environment.
<code>Testing</code>	General testing, debug output, lower user password security for easy testing and e-mail bouncing.
<code>development</code>	Lowest security and e-mail bouncing, by default without any caching.

Multiple versions of the same project with different application environments can be installed on a single server using server specific settings e.g. in the `.htaccess` file or by using url parsing code in the `index.php` file of the application.

NEN 7510 certification requires the installations containing the testing and development environments to be on different hardware than the production environment. On the other hand: there are no objections against a demo environment on the same server as the production environment. The acceptance environment *can* be on the same server as long as a new version is first tested on separate hardware. Of course in all cases the content of the data in the production environment must be separate from the data in any other environments, both when the data is on the same hardware and when it is on a different system.

GemsTracker Deployment Scenarios

A GT deployment always consists of these parts:

- The code specific for the project.
- The GT code.
- The Zend Framework code.
- At least one survey source installation.
- A MySQL database server and a (possibly different) database server for the survey source.

For all current source types the GT installation needs to be able to directly access the database of the survey source – this does not have to be a MySQL database, though a database engine supported by the Zend Framework is necessary.

The survey source itself can be installed on a different server, but this is not required. Usually a subdirectory of the GT installation is used or the source is installed using a separate sub domain.

Here are some example url's. When installing GT on a server with an existing site, it may look like this:

- `www.project.url/gems` – GT installation
- `www.project.url/lis` – LimeSurvey installation

GemsTracker Hosting

When the project server only uses GT and LimeSurvey this is the usual installation:

- `www.project.url` – GT installation
- `www.project.url/ls` – LimeSurvey installation

LimeSurvey is installed in a subdirectory as GT is the default url where users should go to.

Also possible is an installation using sub-domains for the different installations:

- `gems.project.url` – GT installation
- `survey.project.url` – LimeSurvey installation

Using sub-domains may be costly in terms of SSL certificates, as all url's should be used over HTTPS/ Working with sub directories using a common root allows the use of a single standard SSL certificate. However, due to the design of GT individual sub-directories should not be real subdirectories of the web-root. Instead symbolic links or virtual directories should be used – as will be explained later.

We hope this demonstrates that GT deployment is flexible and can be adapted to different needs and to the infrastructure available.

consultation, but that doesn't mean everybody can sign on for free.

Database security

While some projects install all the project data in a single database, this is not the preferred storage configuration for GT. A safe GT installation consists of two or three databases. The number of database users will be the same, but their configuration will be different as their will not be one user for each database, but rather the different users have different roles in accessing the database.

There should be a *project_gems* database containing the patient identifying information, at least one *project_ls* database for each LimeSurvey instance used in the project. When the data should be exported to a data warehouse the extra *project_data* database is needed.

The database user for the GT application must have read/write access to all the databases. The user(s) for the LimeSurvey installation(s) may have only read/write access to their LimeSurvey database. Actually multiple LimeSurvey installations can share a single database using different prefixes for their tables, but we advise to keep this data in different databases. The optional data warehouse user should have read-only access to both the *project_gems* and the *project_gems_data* databases.

Database	Accessible by user	Description
<i>project_gems</i>	Gems, Data warehouse (read-only)	Contains identifying data
<i>project_gems_data</i>	Gems, Data warehouse (read-only)	Optional, contains medical data for data mining
<i>project_ls</i>	Gems, LS	Contains medical data in survey answers

GemsTracker Hosting

When the installed CMS uses a database, then that database should be isolated from the GT databases and the CMS database user should have no access to GT data. The reverse is not required but is good practice.

GT, LimeSurvey and Openrosa in conjunction with GT can all build their own database structure using their web-interface, but the databases should exist prior to initialization.

The best practice for the GT database is to define the database server, name, user and password in the file `var/settings/db.inc.php`.

Email transport

GT applications can send out heavy mail loads and servers may not hinder these mails and should do what they can to make sure the mails are not seen as spam.

GT has built in functionality that routes email through different servers depending on who is sending the e-mail – including secure login when required. E.g. the Erasmus MC does not allow any @erasmusmc.nl mail to be distributed from other servers than their own, so external GT installations sending mail from Erasmus MC users have to login to their server to get the mail distributed. GT does this, though unfortunately it requires storing a password in the database.

In case of organizations that do not allow remote login for sending mail and that do limit the locations their mail is sent from, the IP address of the location where the mail is sent from should be added to the list of allowed originating mail IP addresses.

In all other cases the provider should monitor mail blacklists to make certain that the GT server does not appear on them and take action when this does happen.

Hosting GemsTracker

GT 1.6 is written in PHP 5.6, uses a MySQL 5.1 or higher database and works on most webserver platforms including Windows/IIS, Windows/Apache, Unix/Apache and Unix/Nginx. The application is built using the Zend Framework.

The preferred versions is PHP 7.3 (*not yet 7.4 unfortunately*) and MySQL 5.7.

Technical requirements for hosting GemsTracker

There are no hardware requirements; or rather there is a strong preference to using virtual servers so that we can quickly adjust the hardware to the requirements of a specific project. So we specify only the software requirements for hosting GT, both minimal and preferred.

	Minimal	(Strongly) Preferred
OS Web Server	Windows 2008, Unix	Windows 2018, Unix
Web server	One of these: <ul style="list-style-type: none">• IIS with ReWrite module• Apache 2 with mod_php and mod_rewrite• Nginx	One of these: <ul style="list-style-type: none">• IIS with ReWrite module• Nginx
OS Database	Windows 2008, Unix	Windows 2018, Unix
Database server	MySQL 5.1 Community	MySQL 5.7
PHP	5.6.x	7.3.x
PHP Modules	Default PHP Modules plus: <ul style="list-style-type: none">• Curl• Dom• Fileinfo• JSon• GD2• Ldap• Multi-Byte String• MySQLi• MySQLnd• SOAP• XMLReader• XMLWriter	Minimal PHP Modules plus one of these: <ul style="list-style-type: none">• APC• Memcached• XCache
Other software	Git command line tools and the Composer PHP Package manager (getcomposer.org)	
SSL	Always required for each URL	Always required for each URL
Mail server	Required	Required
LDAP server	Optional	Optional

Server communications

Communication with the server

The webserver is usually only accessed from the outside using https of port 443.

Maintenance by the application management team is done through either SSH (Linux/Unix) or RDP (Windows) from fixed IP addresses.

Communications by the server during normal operations

While running the application will only send out mails and optionally perform LDAP login-tests.

Multiple servers can be used as mail server if required and the server can use any ports and encryption schemes and can also use a username / password combination to login. Usually contact with the main mail server of the organization suffices. This is preferred as then the changes of mail being seen as spam drop.

GemsTracker Hosting

For LDAP authentication we just need the fully qualified domain name of at least one, but preferably at least two LDAP servers. GemsTracker performs no LDAP queries, but just tries to authenticate a user with the password just entered. These passwords are NOT stored locally in any form.

Depending on data being supplied by the hospital additional ports may be required for daily operations.

Communications by the server during maintenance

When the application managers are logged in to the server using SSH or RDP they need access to these URL's to update the software:

- getcomposer.org
- www.getcomposer.org

- github.com
- api.github.com
- codeload.github.com
- status.github.com
- www.github.com

- limesurvey.org
- www.limesurvey.org

- packagist.org
- repo.packagist.org
- www.packagist.org

On Windows system we'll need additional access to extra URL's:

- microsoft.com
- www.iis.net
- update.microsoft.com
- www.update.microsoft.com
- www.microsoft.com

- git-scm.com
- www.git-scm.com

- php.net
- pear.php.net
- windows.php.net
- www.php.net

- notepad-plus-plus.org

- www.notepad-plus-plus.org

Directory structure

The GT code is by default installed in multiple sub directories, only one of which is accessible from the web. The other code, library and storage directories should be stored in a location accessible by the web server, but not from the web.

- `application` – The project specific code.
- `htdocs` – The web root directory containing `index.php`, JavaScript and stylesheets, and possibly code for survey sources.
- `library` – The GT code.
- `var` – writeable by the web server; contains cache, uploaded files and server specific settings e.g. for database access.

Other directories you may find on the server:

- `scripts` – Scripts for the application, e.g. for command-line invocation of GT.
- `test` – Unit tests, not needed on deployments, but part of the development environment.

Unless specifically told otherwise in the `index.php` file, GT applications assume all these directories to be in the same parent directory, including the `htdocs` directory – though this is the one directory that can have a different name without breaking the code.

As this structure would cause problems with multiple installations running from subdirectories on a webserver, the usual deployment scenario is to link the `htdocs` directory to the web root of the webserver using symbolic links with Apache (on Unix using `ln` and on Windows using `mklink`) and virtual directories with Windows IIS.

PHP.INI settings

In general the standard php.ini settings for production and development servers suffice for the production / all other environments. However, some non-standard php.ini settings should be set.

Setting	Example	Comment	Per project alternative
date.timezone	Europe/Amsterdam	Should be set to the time zone the project should use.	In the index.php of the project.
error_log	/tmp/php-error.log	Must already be set on php startup.	Override in the project's application.ini. Starts with the php.ini log, switches during startup to the project log.
include_path	/path/to/zend/1.x	Can be placed in the projects sub-directory, but sharing is preferred.	In the index.php of the project.
soap.wsdl_cache_dir	/tmp/	Must be set to writable directory (when used).	Override in the project's application.ini.
upload_max_filesize	60Mb	The maximum size for uploaded survey PDF's	Override in the project's application.ini.

Scheduled jobs

GT projects use cron or Windows Schedule for automated tasks.

Currently these jobs usually run twice a day at 7:00 o'clock and 19:00 o'clock, but we are developing a sub-scheduler within GT to enable different jobs to run at different moments in time/ These jobs will still be started from the central GT cron function that should then run every 15 minutes or whatever time is appropriate for the specific project.

GT can be started directly from the command-line so no in-between application like `wget` is needed. All the scheduler needs to start web root `index.php` with PHP and with the parameter `cron`:

```
php -f index.php -- cron
```

Of course the location of the PHP executable and the index script must be added to this command.