

The IMS Open Corpus Workbench (CWB) CQPweb System Administrator's Manual

— CQPweb Version 3.2.32 and above —

Andrew Hardie
<http://cwb.sourceforge.net/>

June 2019

Contents

1	Installing CQPweb	7
1.1	What you will need	7
1.2	Hardware requirements	7
1.3	Installing the webscripts	8
1.4	Setting up Corpus Workbench	9
1.5	Setting up R	9
1.6	Setting up PHP	9
1.7	Setting up disk locations	11
1.8	Extra security on disk locations	12
1.9	Setting up your webserver	13
1.9.1	Overview	13
1.9.2	Using HTTPS	13
1.9.3	Specific webserver: Apache	14
1.10	Setting up MySQL	15
1.10.1	Creating the database	15
1.10.2	Known “gotchas” in the MySQL setup	16
1.10.3	Using a separate computer for MySQL	18
1.10.4	MySQL file access	19
1.11	Creating a configuration file	20
1.12	Completing setup	21

2	The CQPweb Configuration File	22
2.1	About the configuration file	22
2.2	Compulsory configuration variables	23
2.3	Optional configuration variables:	24
2.3.1	Locations of programs on the system	24
2.3.2	MySQL features	24
2.3.3	Memory, disk cache, and other hardware resource limits	25
2.3.4	Configuring the user interface	26
2.3.5	Tweaking the look-and-feel	27
2.3.6	User account creation	29
2.3.7	User corpus system	30
2.3.8	RSS feed control	31
2.3.9	Error reporting	31
2.3.10	Miscellaneous configuration options	32
2.4	Using the auto-configuration script	35
2.5	Using the configuration file template	35
2.6	Changes from earlier versions of CQPweb	35
2.7	Obsolete feature: the corpus settings file	36
3	The System Administrator's Interface	37
3.1	Introduction	37
3.2	The Admin Control Panel: Feature list	37
3.3	Corpus Admin Tools: Feature list	39
3.3.1	Corpus settings	39
3.3.2	Manage access	40
3.3.3	Manage text metadata	40
3.3.4	Manage text categories	40
3.3.5	Manage corpus XML	40
3.3.6	Manage annotation	40
3.3.7	Manage frequency lists	40
3.3.8	Manage visualisations	40
3.3.9	Cached queries	40
3.3.10	Cached databases	40
3.3.11	Cached frequency lists	40

4	Managing the CQPweb data cache	41
4.1	Introduction	41
4.2	Some background on the MySQL database system	41
4.3	Explaining the different types of cached data	43
4.4	Disk locations for stored data	43
4.5	Moving the cache location on an existing CQPweb server	44
4.6	Optimising MySQL for cache performance	44
4.7	User-data cache sizes	45
4.8	Finding and fixing cache leaks	45
5	Administering CQPweb from the commandline	46
5.1	Introduction	46
5.2	The main <i>cqpweb</i> script	46
5.3	autoconfig.php	47
5.4	autosetup.php	47
5.5	cli-lib.php	47
5.6	execute-cli.php	48
5.7	force-innodb.php	48
5.8	install-corpus.php	48
5.9	load-pre-3.1-groups.php	49
5.10	load-pre-3.1-privileges.php	49
5.11	load-pre-3.2-corpsettings.php	49
5.12	offline-freqlists.php	49
5.13	upgrade-database.php	50
6	Indexing corpora	51
6.1	Quick checklist	51
6.2	Basic concepts	51
6.3	The notion of a handle	52
6.4	File format	53
6.5	Linking handles and descriptions	53
6.6	Annotation	53
6.7	Annotation templates	53
6.8	XML	54
6.9	XML templates	54
6.10	The indexing process	54
6.11	Using a pre-indexed corpus	54

6.12	The metadata setup process	54
6.13	Building frequency lists	54
6.14	Linking annotation to CEQL syntax notation	55
6.15	Setting up corpus access rights	55
6.16	Further corpus configuration	55
6.17	Putting corpora into categories	55
7	Metadata	56
7.1	Introduction	56
7.2	Corpus metadata	56
7.3	Text metadata	56
7.4	XML metadata	57
7.5	The different possible datatypes	57
7.5.1	Free text	57
7.5.2	Classification	58
7.5.3	Unique ID	58
7.5.4	ID-link	58
7.5.5	Date	62
7.6	Metadata templates	62
7.7	Matadata file format	62
7.8	Installing metadata	63
8	Parallel corpus data	65
8.1	Introduction	65
8.2	Setting up parallel corpora	65
8.3	Naming alignment attributes	65
8.4	Creating alignment attributes	66
8.5	Registering alignment attributes with CQPweb	66
8.6	How alignment attributes can be used	67
8.7	Parallel corpora and user privileges	68
9	Controlling query visualisation	69
9.1	How the primary annotation affects visualisation	69
9.2	Setting up an “alternate” view for context display	69
9.3	Using position labels	69
9.4	XML visualisations	69
9.4.1	Introduction	69
9.4.2	Creating and managing XML visualisations	69

9.4.3	Conditional XML visualisations	70
9.4.4	The embedded variable	71
9.4.5	HTML allowed in XML visualisation code	71
9.4.6	Extra code files	73
9.4.7	Fallback visualisation methods	76
9.5	Field data presentation mode	76
9.6	Field data mode as a workaround for parallel corpora	76
10	User accounts and privileges	77
10.1	Basic concepts	77
10.2	User accounts	77
10.3	Viewing user account details	79
10.4	User groups	79
10.5	Privileges	81
10.5.1	Corpus access privileges	81
10.5.2	Frequency list privileges	82
10.5.3	Database privileges	82
10.5.4	File upload and disk space privileges	82
10.5.5	The CQP binary file privilege	83
10.5.6	Corpus installation privileges	83
10.5.7	Creating and editing privileges	83
10.6	Grants: creating and managing grants of privileges	84
10.7	Running an open server	85
11	Using plugins	86
11.1	What is a plugin?	86
11.2	Types of plugin	86
11.2.1	Annotators	86
11.2.2	Corpus Installers	87
11.3	Installing and registering plugins	87
11.4	Creating plugins	87
11.4.1	Introduction to writing plugins	87
11.4.2	Naming your plugin	88
11.4.3	Methods your plugin must implement	88
11.4.4	Methods you can inherit	92
11.4.5	An API for plugin writers	94
11.5	Builtin plugins	95

11.5.1 BasicTokeniser	95
11.5.2 TreeTagger	95
11.5.3 UcrelTagger	96
11.5.4 BasicVrtInstaller	96
11.5.5 SimplePlaintextInstaller	96
11.5.6 StandardToolInstaller	97
11.6 Permissions for plugins	97
12 Extensible CQPweb	98
12.1 Introduction	98
13 Using the CQPweb API	99
13.1 Introduction	99
14 Updating CQPweb	100
14.1 The update process	100
14.2 Updating the database from very old versions	100
14.3 Updating from version 3.0.16 to version 3.1.0	102
14.4 Updating from version 3.1.7 or earlier to version 3.1.8 or later	102
14.5 Updating from version 3.1.8 or earlier to version 3.1.9 or later	102
14.6 Updating to version 3.2.0	102
14.7 Updating to version 3.2.4	104
14.8 Updating to version 3.2.6	104
14.9 Updating to version 3.2.23	104
14.10 Updating to version 3.2.32	104

1 Installing CQPweb

This chapter contains only the minimum amount of information you need to get CQPweb up and running; it touches on many aspects of the system, but does not go into full detail. Further information can be found in other sections of this manual.

1.1 What you will need

You'll need a machine with a Unix-style operating system. CQPweb has been installed successfully, to our knowledge, on Mac OS X; Sun OS; Debian; Ubuntu; SUSE; and Fedora.

Windows compatibility is planned but not yet achieved, so for the moment if you are on Windows you will need to install CQPweb, and all its dependencies, on top of Cygwin.

Other software you need to have installed:

- Apache or some other webserver
- MySQL (v5.0 at minimum, and preferably v5.7 or higher)
- PHP (v5.6 at minimum, and preferably v7.2+)
- Corpus Workbench (see [1.4](#) for version details)
- R
- Standard Unix-style command-line tools: **awk**, **tar**, and **gzip**; either GNU versions, or versions compatible with them.

A word of warning: Installing CQPweb on a shared server where you do not have full control over the setup can sometimes be problematic. For instance, as will be explained, there are certain options (especially in PHP and MySQL) which need to be set in certain ways for CQPweb to work properly. If you don't have control over these settings, then it will be very difficult to get things working properly. For instance, MySQL servers can be configured to block some of the permissions that you will need to have - so if you can't reconfigure those permissions you will not get very far.

1.2 Hardware requirements

It is difficult to generalise about what hardware you will need. Even a (relatively) low-powered computer should be able to run CQPweb with small-to-medium sized corpora; a modern desktop system should also be fine with pretty large corpora - as should most laptops as long as they have a big enough hard drive (see below). As a general rule, having less-than-ideal hardware will cause things to run slowly, rather than not run at all.

One potential bottleneck is working memory for big database operations. For very intensive operation, MySQL requires lots of memory; if it can't get enough RAM, it will use hard-disk space as temporary storage instead; if it uses up all available hard-disk space, it will fall over mid-operation, possibly with a very uninformative error message (e.g. saying that it can't read from or write to a particular temporary file). These "big database operations" tend to be aspects of the corpus-setup procedure - especially frequency-table creation - rather than anything that the non-administrator user can set in motion.

《*XREF to the CACHE chapter – section on disk space and advising on partitioning etc*》

TODO

How much disk space precisely MySQL might need for some of CQPweb's big setup procedures is difficult to say for certain - it depends, apart from anything else, on how far the MySQL daemon can get just with RAM. However, as a rule of thumb, you should aim to run CQPweb's MySQL database on a disk or partition which has free space equal to a multiple of the raw-text size of the corpus you are working with. (This will also be more than enough for cache space and CWB-indexing of your corpus, if the cache and data directories are on the same disk or partition.)

For example, on one of the CQPweb development machines - a relatively modern server with plenty of RAM - MySQL needed approximately 4 GB of temporary disk space for the process of building frequency tables for a corpus whose raw text took up 1 GB. So if you are regularly dealing with corpora of that size (on the order of 100,000,000 words), it is a good idea to have, say, ten times as much space available as your raw text takes up.

1.3 Installing the webscripts

Download CQPweb by going to the CWB website (<http://cwb.sourceforge.net/download.php#gui>) and doing one of the following:

- Get a release of CQPweb as a compressed download file, then decompress it. This will get you a stable of the program, but one that might be quite old. If your system does not have Subversion installed on it, then this is the only option.
- Use Subversion to *export* a copy of the program from our code repository. You can get either the *trunk* version (cutting-edge version, may contain bugs) or one of the older *branches* (may lack recently-added features but should have fewer bugs), as follows:
 - `svn export http://svn.code.sf.net/p/cwb/code/gui/cqpweb/trunk CQPweb`
 - `svn export http://svn.code.sf.net/p/cwb/code/gui/cqpweb/branches/X.Y.Z CQPweb`
 - (for a branch, replace “X.Y” with the version number of the branch you want; the CWB website will say what branches are available/recommended).
- Use Subversion to *check out* a copy of the program. The difference between checking out and exporting is that a checked-out copy can be automatically updated if the version in the repository changes. This makes updating the system easy. This is, therefore, recommended. Again, you can check out either the trunk or a branch.
 - `svn co http://svn.code.sf.net/p/cwb/code/gui/cqpweb/trunk CQPweb`
 - `svn co http://svn.code.sf.net/p/cwb/code/gui/cqpweb/branches/X.Y.Z CQPweb`

When you first create it, the base CQPweb directory will contain several subdirectories, as follows:

adm Web-directory for the admin interface.

bin This directory contains scripts that can be run offline using command-line access to the machine that CQPweb runs on. See chapter 5.

css Web-directory for stylesheets and other files related to the appearance of CQPweb.

doc Web-directory for manual files.

exe Web-directory for the corpus-query interface.

jsc Web-directory for client-side JavaScript code.

lib This directory contains the actual CQPweb code. CQPweb never runs from this directory, but all the directories where it *does* run operate by calling the code found here.

rss Web-directory for the RSS feed, if enabled (see RSS-related configuration options in section 2.3).

usr Web-directory for the user account interface.

You should *never* rearrange the internal structure of these directories - it will break CQPweb if you do so.

Once you have downloaded/exported/checked out CQPweb, move its base directory into your web server's document tree. Note that the location you choose for the web-script directory will determine the web address of your CQPweb installation relative to your web-server as a whole.

You may then need to adjust the ownership/permissions of the base directory and the files within it to make sure that the user account your webserver runs under has access to it. See section 1.7 for more information on this process.

1.4 Setting up Corpus Workbench

If you do not already have CWB on your system, you will need to install it before going further.

Instructions and links for installing the core CWB system can be found at <http://cwb.sourceforge.net/download.php>.

You should install the most recent available version of CWB from the 3.4.x series. Older versions do not work as well with CQPweb.

When installing CWB, you have assorted options, some of which affect the location of the executables once installed. The default location for the executables under Linux is `/usr/local/bin`, but it is possible to install them elsewhere. Wherever they are, note that it will be important for CQPweb to be able to find them. By default, CQPweb assumes these executables are on the `PATH` variable in the environment the web server runs on. If this is not the case, then you can tell CQPweb explicitly where to find the CWB executables using the `$path_to_cwb` variable in the configuration file (see 2.3).

1.5 Setting up R

There are no special considerations to note in relation to the R statistical software.

1.6 Setting up PHP

You need at least version 5.6 of PHP and preferably version 7.2+.

- The **zlib** extension is required
- The **mysqli** extension is required (see below)
- The **gd** extension is needed if you want the user-account-creation function to be protected by CAPTCHA.

CQPweb crucially requires either the **mysqli** extension to PHP in order connect to the MySQL database. Nearly all versions of PHP are almost certain to include **mysqli**. However, on many Linux distributions (Debian-based or Fedora-based particularly) you may need to select and install a separate

package to get **mysqli**. On other distributions, including possibly on Windows, the **mysqli** extension may be present but not enabled: in which case, it can be enabled by amending the `php.ini` file (discussed in detail below).

In the unlikely event that your version of PHP does not have **mysqli** *at all*, the only fix is to recompile PHP to add it. (It would probably be easier to reinstall PHP from a source that *does* include **mysqli**).

If you are running CQPweb on the internet (i.e. not simply on a standalone computer), then it is also crucial for CQPweb to be able to send email via PHP's `mail()` function. You can check whether this is working by running the following command:

- `php -r "mail('you@somewhere.net', 'it works...', 'Yes it works.');" >/dev/null`

(obviously using your real email address). If you get an email, then PHP can indeed use the `mail()` function. If not, you need to reconfigure your system to allow PHP to send email. It's beyond the scope of this manual to explain how to do this; the PHP manual has quite a lot of information (<http://php.net/mail>) and more can be found via web-search.

Certain aspects of the behaviour of PHP are controlled by its `php.ini` file (see <http://php.net/configuration.file>, <http://php.net/ini>). Your system may have several files with this name, for different environments; you need to find the one used when PHP is run via your webserver; one known “gotcha” on Apple OS X is that the `php.ini` file may not exist, if not, then there should be a `php.ini.default` file which you can copy into the correct directory (under the name `php.ini`) and then amend if need be.

The precise settings in the `php.ini` file may be different depending on how you installed PHP (or your system administrator may have subsequently adjusted them). Most of them will not affect CQPweb.

However, four directives set limits on PHP's use of system resources, and CQPweb has been written on the assumption that these directives are set to at least moderately generous values; if your system's settings are much less than these, you may have problems. The directives in question are as follows:

- `upload_max_filesize` needs to be quite high if you want to upload corpus files for indexing over HTTP; we recommend 20M (for files larger than that relying on HTTP is probably a bad idea anyway)
- `post_max_size` needs to be at least as high as `upload_max_filesize`
- `memory_limit` should be generous as some CQPweb operations are RAM-intensive (e.g. building subcorpus definitions in memory); a reasonable starting point would be 128M (but if the default in your system is higher than that, keep the higher value!), while bearing in mind that certain types of corpora - for instance, those with complex XML structure - may need more RAM for common operations, e.g. 512M.
- `max_execution_time` should be generous as well; we suggest 60

Less “generous” limits for file size, memory use, and execution time may be OK if you are running CQPweb on a standalone computer rather than a shared web server - although on the other hand, on a standalone computer it doesn't matter nearly so much if CQPweb locks up all the system's resources for long periods at a time! Once CQPweb is up and running, the values of all these settings can be checked by looking at “PHP configuration” in the Admin control panel. You can also find out the location of the active `php.ini` file from this screen.

If you are running a version of PHP with the Suhosin patch (which comes by default in some Linux distributions) then there is an additional “gotcha” to look out for. This is that a limit is placed on

the length of individual values in an HTTP request - 512 bytes by default. This can result in very long CQPweb queries failing to work if they overrun this limit. If you have Suhosin and you want to be sure of not running into this problem, you need to add the following line to `php.ini`:

- `suhosin.get.max_value_length = 8000`

However, even with this limit increase, users may have trouble with very long queries, since different browsers may impose alternative, lower limits on HTTP request values. For example, some versions of Internet Explorer are known to restrict HTTP request values to 2000 bytes.

Finally, PHP has a “safe-mode” (see http://php.net/features.safe_mode). This mode restricts what PHP can do in an attempt to provide security. However, it has been recognised as misguided, and is deprecated in PHP v5.3 and removed in PHP v5.4. If your system has safe-mode enabled in spite of this, you may well find that some CQPweb operations do not work:

- PHP safe-mode may not allow CQPweb to run CQP and R as separate child processes if the `cqp` and `R` executables are not in its “safe” directories. If this happens, CQPweb will simply not work at all.
- PHP safe-mode will stop CQPweb from lifting restrictions on execution time for long-running operations; these functions (mostly set-up functions for big corpora, but also processes like collocations on queries with lots of results) will stop working.
- PHP safe-mode will restrict CQPweb to only opening or modifying files that CQPweb has created itself. However, some operations need CQPweb to work with files which it didn't create itself but were, rather, created by the MySQL server or manually by you. All these operations will stop working.

The solution is to turn safe-mode off. Find the line in your `php.ini` file that looks like this:

- `safe_mode on`

and change it to

- `safe_mode off`

1.7 Setting up disk locations

As well as the location of the web scripts themselves, CQPweb needs you to allocate *four other directories* for its use. These are used for the following purposes:

1. CWB corpus index files
2. CWB registry files
3. Temporary files (including the query cache)
4. Files uploaded into the system by users

All four of these directories should be *outside your webserver's document tree* so they are not exposed to the world. Once you have created them, you should leave them exclusively to CQPweb; no one else should add, amend or delete files in any of the four directories.

You will need to enter the paths to these directories in the configuration file: see section 1.11.

If you also use CWB and CQP from the commandline on the same machine that is running CQPweb, then it's worth noting that the locations you choose for the CWB index data and registry *do not* need to be the same as the ones used normally by commandline CQP. CWB has, compiled into it, a default registry path, but CQPweb *does not use that directory*.

The username of the webserver process needs to have full read-write-execute access to all four directories. The username of the *mysqld* process also needs read and write access to the third and fourth (temporary/upload) directories if you want MySQL to use file-access functions, as described in 1.10.4.

The easiest way to accomplish this is to give read-write-execute permissions on these folders to “all”, or - if you are worried about security - to “group” (where the file is assigned to some group that both the MySQL server's account and the web server's account belong to).

How to work out the usernames of the server programs: For *mysqld*, the username is usually *mysql*. For the Apache webserver (process name something like *httpd* or *apache2*) it is usually something like *apache* or *www* or *www-data*. To find out for certain, run the following command (in this example, for *mysqld*):

- `ps -e -o user,comm | grep mysqld`

... and the first word on the output line will be the username you want.

1.8 Extra security on disk locations

A “gotcha” can occur when creating the directories discussed above in systems that extra security software installed. These systems limit the areas of the filesystem that server programs like Apache or *mysqld* can access - blocking access to anything outside the designated areas, *even if* the server has all the necessary filesystem permissions.

The two systems of this type that are often encountered on Linux are **AppArmor** and **SELinux**.

AppArmor disallows file access to programs that it thinks ought not to be manipulating a particular directory, even if the username of that program has all the proper permissions. If it does not allow *mysqld* access to the CQPweb directories, you will not be able to run CQPweb successfully.

You can fix this by adding exceptions to AppArmor's configuration file for *mysqld*. In the systems we have seen this on, the file to edit is:

- `/etc/apparmor.d/usr.sbin.mysqld`

(The filename of the configuration file represents the path to the executable whose filesystem access is being restricted.)

An AppArmor configuration file must, of course, be edited as root. Before the closing brace in this file, add a line like the following for each CQPweb directory:

- `/path/to/the/directory/in/question/** rw,`

Then restart AppArmor (`/etc/init.d/apparmor restart`).

The alternative is to disable AppArmor completely if you do not need the extra security it supplies.

SELinux (<https://selinuxproject.org/>) works differently from AppArmor but can have the same effect of blocking a program's access to files and directories that the program's username has filesystem permissions to read and/or edit. If enabled, SELinux may block Apache (and possibly also MySQL) from accessing/modifying your CQPweb disk locations.

Rather than editing a configuration file, telling SELinux to allow a given program access to a given folder is done via a pair of commandline utilities: `semanage fcontext` and `restorecon`. Please consult the SELinux documentation for full information. However, the following commands should work to prevent Apache from being blocked in most circumstances:

- `semanage fcontext -a -t httpd_sys_rw_content_t "/path/to/directory(/.*)?"`
- `restorecon -R -v path/to/directory`

The first of the commands above adds a “file context” for the path specified which makes the directory, and files/subdirectories, available to Apache for read and write. The second reloads the context for that path so that your added file context will take effect.

1.9 Setting up your webserver

1.9.1 Overview

There are two things that you need to make sure are configured in your webserver:

- It *must* be configured so that files with the `.php` extension are run as PHP (whether via CGI or via a module like Apache's `mod_php`). This is the default on most webservers.
- It *must* be configured so that a file named `index.html` or `index.php` is served as the default when the address of a directory is accessed (so `http://my.server.net/directory` produces the same as `http://my.server.net/directory/index.php`). This is also usually the default situation.
- It *must* allow symbolic links in URLs. CQPweb corpora are addressed via URLs of the general form `http://my.server.net/CQPweb/corpus`, but the actual entry for `corpus` within the CQPweb web-directory is implemented as a symbolic link, not an actual directory.

Other steps that are not necessary, but that might be useful, include the following:

- Block HTTP access to the `bin` and `lib` subdirectories of the CQPweb base directory.
- Turn off the use of `.htaccess` files (Apache webserver only).
- Set up your webserver to use HTTPS instead of HTTP for CQPweb.

These are discussed in further detail below.

1.9.2 Using HTTPS

《write this》

TODO

1.9.3 Specific webserver: Apache

Since Apache is the most commonly used webserver on Unix systems, we have accumulated more experience on installing CQPweb alongside Apache than any other server. Indeed, older versions of CQPweb actually relied on Apache for username/password checks (this was changed in version 3.1). The notes in this section outline some of the most common points of Apache configuration that are important for CQPweb.

Apache configuration is a rather complex topic, and cannot be dealt with in full here (see <https://httpd.apache.org/docs/>). In particular, note that it is impossible to say here specifically what Apache configuration files you need to edit to adjust how Apache treats CQPweb, since this can differ drastically from system to system; especially if you are on Linux, a lot depends on how your distro has decided to package Apache: see <http://wiki.apache.org/httpd/DistroDefaultLayout>.

However, here is some general advice.

Apache's behaviour is controlled by various configuration directives. These directives can be given in the main configuration file, or they can be given in **.htaccess** files that may optionally be added to each directory in Apache's document tree.

If you change an Apache configuration file, you will need to restart the webserver process for the changes to have effect. This is not necessary if you are using **.htaccess** files.

In earlier versions of CQPweb, **.htaccess** files were used to control user access to corpora. *As of version 3.1.0, this is no longer the case.*

The only directories to which it is necessary to control access are the **bin** and **lib** subdirectories of the base CQPweb directory. This *can* be done with **.htaccess** files, but it is better done in one of the main Apache configuration files. That way, the directives are loaded only once (when Apache starts up), and not every time the server is accessed (which is what happens when **.htaccess** files are used, meaning multiple extra disk-reads are required).

The directive which turns on the use of **.htaccess** files is:

- `AllowOverride All`

Conversely, the directive that turns it off is:

- `AllowOverride None`

In Apache's configuration file, this directive may be given globally, or within a **<Directory>** directive. You can create a separate **<Directory>** block for the directory where CQPweb lives, or you can adjust the settings for a higher-level directory, if that will not interfere with other uses of the webserver.

Similarly, you need to make sure that Apache is set to follow symbolic links in URLs (for reasons explained in 1.9). This is enabled by default, but it's advisable to declare it explicitly (because otherwise you are dependent on this setting not being affected by changes elsewhere in the Apache configuration). This *must* be done within a **<Directory>** directive, or in a **.htaccess** file; the declaration is as follows:

- `Options FollowSymlinks`

(with other **Options** values added as required).

Putting all the above together, a typical **<Directory>** block for the CQPweb directory would be:

```
<Directory /path/to/cqpweb>
  AllowOverride None
  Options FollowSymlinks
</Directory>
```

Note that the “path/to/cqpweb” that you need to give is an absolute path on your filesystem (it is *not* relative to the root of the web document tree).

which switches on the `FollowSymlinks` option using an `Options` declaration line.

There are, of course, many more things you can do with Apache to tweak how CQPweb is accessed. For example, you can add Apache-internal password authentication. However, if you do this, the webserver-based authentication will be *separate from and additional to* CQPweb’s own system of usernames and passwords.

Finally: there are some known “gotchas” in Apache’s behaviour under certain configurations. The best approach is not worry about the following until and unless the problems as described happen in your installation!

- Sometimes, Apache will happily serve up the built-in pages (e.g. the admin area) but then give you an “Internal Server Error” when you try to access the pages created during a corpus installation. This appears to be because these files are created with 0664 permissions (group-writeable, world-readable).
 - To fix the problem for an already indexed corpus, open a terminal to the directory containing its script files (`index.php`, `concordance.php`, etc.) and run `sudo chmod 644 *.php`.
 - To prevent the problem from recurring, edit the code file `lib/admin-install.inc.php` and change all instances of `chmod()` where the mode is set to 0664 to set it to 0644 instead.
- The `PATH` environment variable (i.e. the list of locations where Apache will look for executables, if their precise location is not specified) may present a “gotcha”. The `PATH` as seen by scripts running under Apache are *not necessarily the same* as the `PATH` that is available in the login-shell environment of the username Apache runs under. This is because Apache has its own internal system for setting environment variables, using its `SetEnv` directive and related functionality. If CQPweb is having trouble finding the CWB executables, or any other external program such as R, it may be because the `PATH` variable as seen from within scripts running under Apache does not contain their location. This problem is easily overcome by setting the variable `$path_to_cwb` in the CQPweb configuration file (see 2.3).

1.10 Setting up MySQL

1.10.1 Creating the database

CQPweb uses a MySQL database to store most of its ancillary data - that is, everything other than the actual CWB index data.

You will probably need “root” access to the MySQL server in order to set it up. The following instructions are based on the assumption you are accessing MySQL via the command-line client program, but it is also possible to use your preferred graphical interface for this purpose, of course.

First, you must create a new user and a new database for CQPweb to use. The name of the user and the database can be anything you like; for the sake of the example commands in this section we will assume that they are `cqpweb_db` and `cqpweb_user` respectively.

The required MySQL commands are as follows:

- `create database cqpweb_db default charset utf8;`
- `create user cqpweb_user identified by 'cqpweb_password';`

Naturally, instead of “cqpweb_user”, “cqpweb_password”, etc. use the actual username and password that you want the user to have. This username/password combination will be stored in an only-mildly-secure location, so make sure that you do **not** re-use either an account name or a password that is used for any other purpose on the same system.

The reason that we set the default character set to UTF-8 is that this addresses a bug which affects some (but not all!) versions of the MySQL server software – see <http://bugs.mysql.com/bug.php?id=10195>.

Having created the user, we must now give it *all* permissions over the database. If you want MySQL to use file-access functions, the new user also needs to be granted the *file* permission, which is set once-and-for-all on a single MySQL server, rather than at the level of the database. File-access permission is not strictly necessary, but may help speed things up; it can also be useful if the LOAD DATA LOCAL INFILE command is disabled (see 1.10.2 and 1.10.4 below).

- `grant all on cqpweb_db.* to cqpweb_user;`
- `grant file on *.* to cqpweb_user;`

Make a note of the username and password you have used, and of the database name; you will need them for configuration of your CQPweb installation. Also make a note of the server-name needed to access the MySQL server by a webscript. This will probably be **localhost**, assuming that the MySQL server is on the same machine as the webscripts (but see 1.10.3 below).

1.10.2 Known “gotchas” in the MySQL setup

- **Local infile permission:** MySQL can be configured to disable the LOAD DATA LOCAL INFILE command as a security measure – see <http://dev.mysql.com/doc/refman/5.5/en/load-data-local.html> . This command will stop CQPweb working (you will be able to tell this is happening because you will get the error message “ERROR 1148: The used command is not allowed with this MySQL version” when you attempt to set up the metadata for a corpus).

The preferred way to fix this is as follows:

- Edit the MySQL configuration file (usually something like `/etc/my.cnf` depending on your operating system).
- Find the line which deactivates the local infile feature.
- It will be something like `local_infile=OFF, local-infile=0, or set-variable=local-infile=0`
- Change the 0 to 1 (or OFF to ON)
- Restart the MySQL daemon.

Alternatively, the problem should be fixed if you set the CQPweb configuration variable `$mysql_has_file_access` to **true**, because when you do this, the LOAD DATA LOCAL INFILE command is never used. But there are requirements that must be satisfied to use this option – see 1.10.4 below.

If you can neither change the MySQL configuration, nor meet the requirements for `$mysql_has_file_access`, there is yet a third way to solve the problem: set the CQPweb

configuration variable `$mysql_local_infile_disabled` to **true**. This makes CQPweb avoid the `LOAD DATA LOCAL INFILE` command unconditionally. Be warned - doing this should fix things, but it has the potential to be a major performance hit.

- **File privilege with retriected scope:** It is possible to limit the MySQL daemon so that `INFILE/OUTFILE` commands only affect a specific area of the machine's filesystem. This is done by setting a MySQL system variable called `secure_file_priv`.

If `secure_file_priv` is set, then the daemon can only read/write files in the directory that it specifies. This can stop CQPweb from working. There are three ways to fix this problem.

- Use the folder specified by `secure_file_priv` as your *temporary files* directory.
 - * (See 1.7 on the temporary files location).
 - * This solution will not be possible if your MySQL daemon is not solely devoted to CQPweb.
- Change the `secure_file_priv` option to specify the location of your *temporary files* directory.
 - * Edit the MySQL configuration file (usually called `my.cnf`, see above).
 - * Find the line which sets this variable (it will look like `secure_file_priv=/some/path` or `secure-file-priv=/some/path`) and change the path.
 - * Restart the MySQL daemon.
 - * Again, this not practical unless your MySQL daemon is solely devoted to CQPweb.
- Switch off the `secure_file_priv` functionality entirely.
 - * Find the right line of the MySQL configuration file, as per above.
 - * Change the path to an empty string - so the line looks like this: `secure_file_priv=""`
 - * Restart the MySQL daemon.
 - * This solution *does* work if your MySQL daemon is not solely devoted to CQPweb. It is the **preferred and recommended** solution. However, it does have the disadvantage of removing a security measure.

If you are using a shared server and you are not able to change the MySQL configuration, the only remaining option is to stop CQPweb from using the `INFILE/OUTFILE` commands by setting the CQPweb configuration variable `$mysql_has_file_access` to **false**.

A further “gotcha”: prior to about 2015/2016, the default value for `secure_file_priv` was the empty string. But MySQL was changed so that this value now defaults to an operating-system-specific secure path. If you *have no* `secure_file_priv` in your `my.cnf` file, then things will work fine in older versions of MySQL (because the security feature defaults to being switched off) but in newer versions, CQPweb won't work. If you find that features of CQPweb such as collocations, distribution, and so on stop working after a MySQL upgrade, the problem is probably that the default value of `secure_file_priv` has changed. In this case, instead of *changing* the value of `secure_file_priv`, to fix the problem you need to *add* a line `secure_file_priv=""` to your `my.cnf` file to restore what was previously the default state.

- **PHP's connection to MySQL:** PHP needs to know how to connect to the MySQL server via a *socket* in the filesystem. This information is often contained in the `php.ini` file, which contains settings that PHP will load when it starts up. On many systems, connecting to MySQL simply “works” by default, but on some systems you may need to edit your `php.ini` file to tell PHP where to find the socket, by changing the `mysql.default_socket` setting. For instance, if your socket is at `/tmp/mysql.sock` but PHP is looking at `/var/mysql/mysql.sock`, you need to adjust `mysql.default_socket` to `/tmp/mysql.sock`. If you edit a `php.ini` file, make sure it is the one used when PHP is run by the server (whether as CGI or as a module of the webserver itself).

- **MySQL binary logging:** Binary logging is a MySQL feature that can be enabled or disabled, as explained here: <http://dev.mysql.com/doc/refman/5.5/en/binary-log.html>. If enabled, even relatively light use of a CQPweb installation will make MySQL create very large binary log files, ultimately using up all your disk space over time. For this reason, it's recommended that you **disable** binary logging on the MySQL server that will be driving CQPweb. You can disable binary logging as explained here: <http://dev.mysql.com/doc/refman/5.5/en/replication-options-binary-log.html>. In brief, this is done by commenting out or deleting the line containing the `log_bin` command in the MySQL configuration file, or removing the equivalent `--log-bin[=base_name]` directive from the command line that starts up *mysqld*. (In either case you'll need to restart *mysqld*.)

1.10.3 Using a separate computer for MySQL

We normally assume that the MySQL server runs on the same machine as the CQPweb system itself. But it does not have to.

You might want to use two separate machines for the CWB-based and MySQL-based parts of CQPweb, for reasons of performance (for a big corpus and for queries with lots of results, both MySQL and CWB require lots of disk space, disk read/write bandwidth, and processing power).

In this case, CQPweb itself (the web scripts) should be on the same system as CWB, and the MySQL server on a separate system. This affects how you configure CQPweb as follows:

- You will need to insert the correct hostname (or IP address) for your MySQL server machine into the configuration file for CQPweb, instead of the (normal) `localhost` value. See section 2.2.
- Any configuration variable that involves a path to a temporary-storage directory or to the location of a program (see 2.2) needs to refer to *the system with CQPweb on it*, not the system with MySQL on it.
- The optional variable (see 2.3) `$mysql_has_file_access` can only be set to **true** if the paths to the temporary-storage directories are the **same** on both systems (e.g. if they are mounted to the same location). The system cannot check this for you! This is explained in more detail in 1.10.4 below.

You will need to make sure that your MySQL server machine is configured to allow network traffic through the port that your MySQL server is using. How this is done depends on the operating system and firewall software on that machine. Under Linux, you would use a utility such as `iptables` or the more modern `nftables` to modify the operation of the Linux firewall.

The MySQL server *mysqld* usually listens on port 3306, but this can be changed in the MySQL configuration file: if in doubt, check (the MySQL command `SHOW VARIABLES WHERE Variable_name = 'port'` will tell you). Some notes:

- If the CQPweb system is the only user of the MySQL server machine, then for security it would make sense only to open up the port for traffic coming from the IP address of the main CQPweb machine.
- On Linux, it is possible for the firewall to redirect data arriving on other ports to port 3306. So in that case it would not be port 3306 that you would open.

An additional security measure you can implement when creating the CQPweb user in MySQL is to link that account to the particular IP address of the CQPweb machine. The format for this is as follows:

- `create user 'cqpweb_user'@'111.222.333.444' identified by 'cqpweb_password';`
- `grant all on cqpweb.* TO 'cqpweb_user'@'111.222.333.444';`

... instead of the form given above - with the correct IP instead of “111.222.333.444”, of course. When you create the account in this way, only login attempts from the specified IP will be accepted.

1.10.4 MySQL file access

General note: The issues relating to MySQL file access have been discussed in multiple sections of this chapter, where relevant; this section contains an overview.

Many CQPweb operations involve transferring data into MySQL from files in CQPweb's directories (or, conversely, from the database to such files). Specifically, of the directories discussed in section 1.7, the temporary-files location and the uploaded-files location, including any subdirectories, are used in this way.

There are three ways that file data can be transferred between CQPweb and the MySQL server. In rough descending order of speed, these are:

- The MySQL server reads the file directly, using the SQL command `LOAD DATA INFILE`.
- PHP's MySQL client system transmits the file to the server, using the SQL command `LOAD DATA LOCAL INFILE`.
- CQPweb reads the file incrementally and passes its content to the server in small chunks. No `LOAD DATA` command is used.

By default, CQPweb is configured to use the *second* of these methods. This requires the creation of a temporary copy of the file for use by the server, which is by definition slower than the server directly accessing the existing file. However, the difference in performance is probably not huge.

In order for the first method to work, it is necessary for the MySQL server to be able to access the two directory locations mentioned above, which requires certain preconditions to be fulfilled.

If the MySQL server is running on the same computer as CQPweb, the preconditions for the MySQL-server-file-access method are:

- the MySQL server (or, to be precise, the OS-username it runs under) must actually have access to the two directories (see 1.7 for how to make sure of this).
- the MySQL-username used by CQPweb must have `GRANT FILE ON *.*` permissions (see 1.10.1 for more detail on this point).

If the MySQL server is running on a different machine (see 1.10.3), the preconditions for the MySQL-server-file-access method are:

- the two directories in question must be accessible on the machine the MySQL server runs on (e.g. by mounting them as remote drives).

- the two directories must be reachable using the same paths on the MySQL machine as the CQPweb machine (i.e. use the same path for the mount point, or use a symbolic link)
- the MySQL server (or, to be precise, the OS-username it runs under) must actually have access to the two directories - as per above.
- the MySQL-username used by CQPweb must have **GRANT FILE ON *.*** permissions - as per above.

To explain the first two points in a little more detail, let's say that, for example, the temporary-files directory is located at `/var/cqpweb/temp` on the main CQPweb machine. In order for the MySQL server to be able to access files in this folder, first, it is necessary for the same underlying disk location to be accessible from the MySQL machine: e.g. by it being a network drive mounted remotely by both machines, or by the CQPweb machine making that local directory available remotely, by SFTP for instance. Second, the directory needs to be accessible at the same path, either by mounting it at `/var/cqpweb/temp`, or, if it must be mounted elsewhere (e.g. `/mnt/cqpweb/temp`), then by creating `/var/cqpweb/temp` as a symbolic link to its mount-point.

Once you are sure that the relevant preconditions are fulfilled, you can switch to the MySQL-server-file-access method by setting the configuration variable `$mysql_has_file_access` to **true** (see 1.11 and 2.3.2).

When this variable is **true**, the first, fastest method of file transfer is used.

Conversely, as noted above in 1.10.2, if the second, default method does not work due to MySQL's `LOAD DATA LOCAL INFILE` command being disabled, and you cannot use the first method instead, you must fall back on the third, slowest method, by setting `$mysql_local_infile_disabled` to **true**.

N.B.: see note in section 1.10.2 about a possible "gotcha", a MySQL feature that limits **GRANT FILE ON *.*** permissions to a single directory using the `secure_file_priv` option.

1.11 Creating a configuration file

Before going any further with installation, you must create a configuration file. This can be done manually or automatically.

The CQPweb Configuration File is described in its own chapter of this manual (2). As that chapter explains, there are a small number of *compulsory* configuration variables, and a much larger number of *optional* settings. To get CQPweb up and running, you need only create a configuration file with the nine compulsory settings - optional settings can be added later at your leisure.

There are two ways to do this:

- Manually working from a template - see section 2.5
- Using the automatic configuration script - see section 2.4

Either way, you should note that you will need to enter several of the settings you created in the installation steps above:

- The paths to the four directories you created for CQPweb's data (see 1.7);
- The username, password, database name and server name for the MySQL server;
- You will also need to enter at least one system-admin username.

1.12 Completing setup

With your configuration file created, you are ready to run the final steps of the setup process. These include:

- Creating the structure of the MySQL database, and setting up default data;
- Creating accounts for the admin usernames specified in the configuration file;

Although these can in theory be done manually it is much more effective to let the system do it for you. That is the purpose of the auto-setup script, which is called from the commandline as follows:

- `php autoseup.php`

(note that you must be inside the `bin` subdirectory of the base CQPweb directory for this to work). For a general discussion of running commandline scripts see section [5.1](#).

This script makes your CQPweb installation ready to use. It also, as noted above, actually creates accounts for the admin users you specified when creating your configuration file. It will ask you to specify passwords for these users *only at the point when it creates the accounts*. Passwords for CQPweb are not stored in the database - only an encrypted form is stored - and are never saved anywhere on disk.

Once this script has run, you are ready to go. Open a web browser and navigate to:

- `http://your-server.net/path/to/cqpweb/web/directory/`

2 The CQPweb Configuration File

2.1 About the configuration file

This chapter describes the CQPweb *configuration file*.

The configuration file is a text file. It is always called `config.inc.php` and is always placed in the `lib` subdirectory alongside the code files. It contains PHP code creating variables that are used by the rest of the system to control different aspects of how things work. It's important to understand that *none of the things that are set in the configuration file can be changed through the web interface*, even if you are logged on with an admin account. This is for security reasons.

Since the configuration file is a PHP file, you can in theory put any arbitrary code that you like in it - but it is very strongly recommended that you do not do anything other than assign a series of variables as specified here.

The variables can be in any order; though it is convenient to organise the file so that groups of settings that relate to the same part of the system are close to one another, it is not necessary to do so, and in fact the ordering of the variable assignments in the configuration file makes no difference at all to CQPweb.

PHP variable assignment is very simple and will be familiar to anyone who has done a bit of programming. All assignments are of the following form:

- `$variable_name = VALUE;`

In PHP, a file must have `<?php` on its first line so that its contents will be recognised as PHP code. If you use the template file provided (see section 2.5) then this is already present.

There are two kinds of CQPweb configuration variable. First, there is a small group of variables that you *must* set - if you don't, CQPweb just won't work. These can be set up either manually or automatically (see section 2.4).

Then, there is a much longer list of variables that you *can* set, but if you don't, default values will be used. These should be added to the configuration file manually. These two types of variables, compulsory and optional, are listed in sections 2.2 and 2.3.

Every variable has a particular *type*, one of the normal types in the PHP language. Values of different types are specified in different ways. Again, these will be unsurprising to anyone who is familiar with any programming language:

Boolean A value which is either **true** or **false** (which may mean on/off, yes/no, and so on). Use the PHP keywords **true** and **false** to represent these values.

Integer number A whole number, entered as normal decimal digits (no spaces or thousands-separators).

Floating-point number A number with a fractional part, entered in decimal with a `.` for the decimal point.

String A short bit of text. Strings can be surrounded with either single quotation marks or double quotation marks. The difference is that if double quotation marks are used, a wider range of *escape sequences* are available to represent special characters. For the CQPweb configuration file, you are unlikely to need any escape sequences except those for the delimiting quotation marks themselves, which are `\'` and `\"` in a single-quoted string and in a double-quoted string respectively.

2.2 Compulsory configuration variables

Additional information on the variables containing location paths can be found in section 1.7 on “Setting up directories”.

Additional information on the variables relating to MySQL can be found in section 1.10 on “Setting up MySQL”.

Variable name	Description
<code>\$superuser_username</code>	<p>Type: String</p> <p>This should contain the usernames of users who are to be system administrators, separated by the pipe if there is more than one. For instance: <code>"anna bert craig"</code>. You can have as many admin accounts as you like, but you must have at least one.</p> <p>To add a new administrator, first create a normal account for the person (if they don't already have one), and then edit the configuration file to add their username to this variable. To remove an administrator, delete their username from this variable (but note you must never remove the <i>last</i> username!) No other actions are necessary.</p>
<code>\$mysql_webuser</code>	<p>Type: String</p> <p>The username of the MySQL account that CQPweb should use to connect to the MySQL server. It is usual to have a single dedicated MySQL account for use by CQPweb.</p>
<code>\$mysql_webpass</code>	<p>Type: String</p> <p>The password of the <code>\$mysql_webuser</code> account on the MySQL server.</p>
<code>\$mysql_schema</code>	<p>Type: String</p> <p>The name of the database on the MySQL server that you have created for use by CQPweb.</p>
<code>\$mysql_server</code>	<p>Type: String</p> <p>The address of the MySQL server (either a hostname/IP address, optionally followed by a port number, or else a path to a local socket; see the MySQL documentation for more info on this).</p> <p>If your MySQL server is on the same computer as the rest of CQPweb, this variable will usually be the string <code>"localhost"</code>.</p>
<code>\$cqpweb_tempdir</code>	<p>Type: String</p> <p>Location of a directory which CQPweb can use to store its query cache and temporary data files.</p>
<code>\$cqpweb_uploaddir</code>	<p>Type: String</p> <p>Location of a directory to use for CQPweb's “upload area” (storage for files uploaded by you or by users via the web interface).</p>
<code>\$cwb_datadir</code>	<p>Type: String</p> <p>Location of a directory for CWB index data to be stored in.</p>
<code>\$cwb_registry</code>	<p>Type: String</p> <p>Location of a directory that CQPweb can use as its CWB corpus registry. This does not have to be the same as your system's default CWB registry, but it can be if you want.</p>

2.3 Optional configuration variables:

This is a reference guide to the optional configuration variables; many of them are also mentioned elsewhere in this manual.

Some optional configuration variables are not documented yet. This chapter will be expanded over time until it is as close as possible to 100% complete.

《List of undocumented ones can be found in comments in the latex code at this point.》

TODO

2.3.1 Locations of programs on the system

Variable name	Description
<code>\$path_to_cwb</code>	Type: String Default: "" Path of the directory containing the CWB executables (<code>cqp</code> , <code>cwb-encode</code> , and so on). The path can be absolute or relative; if relative, bear in mind that CQPweb always runs from one of the immediate daughter directories of its main directory. If the path contains any space characters, they must be escaped (with an escaped back-slash, e.g. <code>.../Program\\ Files/...</code>). If no path is given, it will be assumed the executables are in the system's usual path.
<code>\$path_to_gnu</code>	Type: String Default: "" Path of the directory containing the GNU (or equivalent Unix-y) utility programs, namely <code>tar</code> , <code>gzip</code> , and <code>awk</code> . These will almost always be on the normal path on a Unix system but may not be on Windows. The same general comments apply as to <code>\$path_to_cwb</code> .
<code>\$path_to_r</code>	Type: String Default: "" Path of the directory containing the R executable; same general comments apply as to <code>\$path_to_cwb</code> .

2.3.2 MySQL features

Variable name	Description
<code>\$mysql_big_process_limit</code>	Type: Integer Default: 5 This variable places a limit on how many <i>big</i> MySQL processes of a single type will be allowed to run at once. There are several types of “big” process (building collocation database, building frequency tables, building sort databases, and building categorised query tables) so more than the limit could run.

Variable name	Description
<code>\$mysql_utf8_set_required</code>	Type: Boolean Default: <code>true</code> Controls how characters are transmitted between the MySQL server and CQPweb. The default is usually OK. If, however, some characters do not display properly in frequency list, keyword or collocation view, then setting this to false may fix things.
<code>\$mysql_has_file_access</code>	Type: Boolean Default: <code>false</code> This variable declares to CQPweb whether or not the MySQL daemon has access to the filesystem on which CQPweb is running, and in particular two key working directories. It is explained in detail in section 1.10.4 .
<code>\$mysql_local_infile_disabled</code>	Type: String Default: <code>false</code> You should set this to true if your MySQL has been set up to disallow the <code>LOAD DATA LOCAL</code> command, and you can't change this set-up. If possible, changing the MySQL configuration to allow <code>LOAD DATA LOCAL</code> is far preferable, it should be noted! See also section 1.10.2 .

2.3.3 Memory, disk cache, and other hardware resource limits

Variable name	Description
<code>\$cwb_max_ram_usage</code>	Type: Integer Default: 50 Some CWB programs allow a RAM usage limit to be set on their activities. When CQPweb calls these programs, it sets the RAM limit to the number of megabytes specified in this variable. The default is 50 megabytes. This variable applies only when CQPweb is run over the web; for the RAM limit that applies when CQPweb is run from the commandline, see the next variable.
<code>\$cwb_max_ram_usage_cli</code>	Type: String Default: 1000 Same as <code>\$cwb_max_ram_usage</code> , but applies when CQPweb is run from the commandline (see 5).

Variable name	Description
<code>\$query_cache_size_limit</code>	<p>Type: Integer Default: 6442450944</p> <p>Controls the size of the query cache (the maximum size, in bytes, to which the temporary directory will be allowed to grow before old cached queries get deleted). Note that this only affects the size of the query cache; anything stored as a MySQL table (such as temporary frequency tables or collocation databases) does not count towards this limit, and are controlled by separate variables also listed in this section. Until the cache limit is reached, the cache will just keep growing! Cached files are never deleted merely due to age, only when the disk space needs to be reused. The default value is 6 gigabytes.</p>
<code>\$db_cache_size_limit</code>	<p>Type: Integer Default: 6442450944</p> <p>Controls the size of the user-database cache (the maximum size, in bytes, to which the MySQL table containing the user-database cache will be allowed to grow before old databases get deleted). This includes data for collocations, sorting, and distribution; categorised query data also counts towards the total but as it cannot be reconstructed, it will never age out of the cache. The default value is 6 gigabytes.</p>
<code>\$restriction_cache_size_limit</code>	<p>Type: Integer Default: 6442450944</p> <p>Controls the size of the restriction cache (the maximum size, in bytes, to which the MySQL table containing the restriction cache will be allowed to grow before old cached restrictions get deleted). This counts only the part of the MySQL database devoted to temporarily-stored restriction data, not any other stored data. Until the cache limit is reached, the cache will just keep growing! The default value is 6 gigabytes.</p>
<code>\$freqtable_cache_size_limit</code>	<p>Type: Integer Default: 6442450944</p> <p>This variable places a limit on how much disk space <i>ad hoc</i> frequency tables in MySQL are allowed to take up. It is expressed as a number of bytes; the default is 6 gigabytes.</p>

2.3.4 Configuring the user interface

Variable name	Description
<code>\$default_per_page</code>	<p>Type: Integer Default: 50</p> <p>The number of results to show per page by default (in concordances, frequency lists, keyword lists, and so on). All tools also allow the results-per-page rate to be altered on a per-query basis.</p>

Variable name	Description
<code>\$default_history_per_page</code>	<p>Type: Integer Default: 100</p> <p>The number of items to show per-page in history-type displays (such as query history, saved queries, and so on). Note that as a general rule you would normally want more of these to display per-page than you would for <code>\$default_per_page</code> - thus the difference in default values.</p>
<code>\$show_match_strategy_switcher</code>	<p>Type: Boolean Default: false</p> <p>The match-strategy switcher will only appear in the Standard Query and Restricted Query forms if this is set to true. This function requires a recent version of the CWB core to work (more recent than 2017-07-01), which is why it is disabled by default. In the future, the default will be changed to true, and eventually the option will be removed entirely.</p>
<code>\$dist_graph_img_path</code>	<p>Type: String Default: <code>"../css/img/blue.bmp"</code></p> <p>This string is the (relative) address of an image file that will be used for the bars in the bar chart mode of the Distribution display. The default is a file internal to CQPweb that creates plain blue bars.</p>
<code>\$dist_num_files_to_list</code>	<p>Type: Integer Default: 100</p> <p>The number of files to display in File Frequency Extremes mode of the Distribution display. The number of files is limited because corpora can easily have thousands or tens of texts, which would make the webpage hard to read and slow to load.</p>
<code>\$uploaded_file_bytes_to_show</code>	<p>Type: Integer Default: 102400</p> <p>When a file uploaded by the admin user (or, in future, by regular users) is displayed, only a certain amount of data is shown: the first section of the file, up to the number of bytes specified in this setting (to the nearest whole line). The default is to show 100 KB. A too-large setting here may cause browser overload, since so many of the files that CQPweb deals with are so very large.</p>
<code>\$hide_experimental_features</code>	<p>Type: Boolean Default: false</p> <p>If set to true, certain new features deemed “experimental” will be hidden in the interface; users will neither see nor be able to use them. What features count as “experimental” will change from version to version.</p>

2.3.5 Tweaking the look-and-feel

Variable name	Description
---------------	-------------

Variable name	Description
<code>\$css_path_for_homepage</code>	<p>Type: String Default: (see below)</p> <p>The relative or absolute URL of a CSS stylesheet to use for the main menu page. Note that, if relative, this must be relative <i>to the homepage</i>, which is one level higher (in the directory tree) than all the other URLs that CQPweb runs from. So if you want to address something in the CSS folder, for instance, you would need to start this variable with <code>css/</code>, rather than <code>../css</code>. By default, a lovely blue-and-grey table effect is used.</p>
<code>\$css_path_for_adminpage</code>	<p>Type: String Default: (see below)</p> <p>The relative or absolute URL of a CSS stylesheet to use for the admin control panel page. By default, a red-and-pink colour scheme is used.</p>
<code>\$css_path_for_userpage</code>	<p>Type: String Default: (see below)</p> <p>The relative or absolute URL of a CSS stylesheet to use for the user-login homepage. By default, a green-and-grey colour scheme is used.</p>
<code>\$homepage_use_corpus_categories</code>	<p>Type: Boolean Default: <code>false</code></p> <p>If this is true, then the list of corpora on the main menu page will be given as a set of lists according to the corpus category, rather than as a single long list of corpora.</p>
<code>\$homepage_welcome_message</code>	<p>Type: String Default: "Welcome to CQPweb!"</p> <p>A little bit of text (which can include HTML formatting) that will appear in the header box of the main menu page.</p>
<code>\$homepage_logo_left</code>	<p>Type: String Default: ""</p> <p>Settings for a logo to display on the left of the header box of the main page menu. This string can contain either (a) a single URL for an image to use as the logo; or (b) two URLs with a tab between them, in which case the logo image becomes a clickable link: the first URL will be used as the address of the image, and the second as the address for the link.</p> <p>URLs can be absolute, or relative to the location of the main menu page (i.e. the URL of your CQPweb base directory). A Corpus Workbench logo that you can use if you wish is located at the relative URL of <code>css/img/ocwb-logo.transparent.gif</code>. If you are running CQPweb on an HTTPS server, note that the logo URL you use <i>must</i> be on the same server (if not, most browsers will warn that the webpage is only partially secure). It is a good idea to add your image files to the <code>css/img/</code> subdirectory.</p>

Variable name	Description
<code>\$homepage_logo_right</code>	Type: String Default: "" Same as <code>\$homepage_logo_left</code> , but whatever you specify appears on the right side of the header box.
<code>\$searchpage_corpus_name_suffix</code>	Type: String Default: "powered by CQPweb" A little bit of text that is suffixed to the name of the corpus in the main search page header. If you don't want anything, set it to an empty string.

2.3.6 User account creation

CQPweb has the facility for users to create their own accounts through a standard “validate-by-email” mechanism. The configuration options in this section control whether this facility is available, as well as various aspects of how it works.

Variable name	Description
<code>\$allow_account_self_registration</code>	Type: Boolean Default: true If this is true, a form will be exposed for anyone to sign up for an account on your server. If false, this form will not be available, and only admin users will be able to create new user accounts. (You can alternatively configure your web server to block or limit access to the signup form if you wish.)
<code>\$account_create_contact</code>	Type: String Default: "" This variable only has any effect if <code>\$allow_account_self_registration</code> is false . In this case, you can supply a snippet of text here that will be added to the web interface to tell prospective users who to contact to request an account. This string can contain HTML markup, for example to add a “mailto” link, e.g. " <code>A. N. Other</code> ".
<code>\$account_create_captcha</code>	Type: Boolean Default: true Determines whether or not the account-creation form is protected by a CAPTCHA challenge. This can help protect your server against web-robots, but you might want to disable CAPTCHA for convenience if your account-creation page is inaccessible to users on the open Internet anyway. If your PHP installation lacks the GD extension to PHP (see: http://php.net/gd), then this setting is <i>always</i> false , regardless of what you specify.

Variable name	Description
\$account_create_one_per_email	<p>Type: Boolean</p> <p>Default: <code>false</code></p> <p>Determines whether or not multiple accounts can be created that are linked to the same email address. By default, multiple accounts with the same email address are allowed. To disallow this, set this option to true. In that case, any attempt to create a second account with an email address already on the system will fail.</p> <p>If you change this setting, it does not apply retroactively: any accounts already in the system that share an email address will not be affected.</p>
\$blowfish_cost	<p>Type: Integer</p> <p>Default: 11</p> <p>CQPweb uses Blowfish to encrypt passwords. The “cost” controls how long it will take for this encryption to run. The higher it is, the harder it is for an attacker to crack a given encrypted password. As of 2019, the default value of 11 is reasonable, but if you are worried about security you might want to try 12 or 13.</p> <p>In the <i>System diagnostics</i> section of the admin interface, there is a tool to stress-test password encryption, to see if your current Blowfish cost is high enough.</p>
\$create_password_function	<p>Type: String</p> <p>Default: <code>"password_insert_internal"</code></p> <p>CQPweb uses “suggest password” functions to ease the creation of user accounts in the admin interface (see 10.2). The default function (<code>password_insert_internal</code>) produces randomised passwords of the form <code>aaaaddaaaa</code>, where a is a lowercase letter and d is a digit.</p> <p>You can optionally create a password function of your own, e.g. to get nicer, more wordlike passwords. This is something you should only attempt if you know how to write functions in PHP! The function you create should take a single argument (an integer), and should return an array of strings, where each string is a suggested password, and the number of strings in the array is equal to the integer argument.</p> <p>Then, set this variable to the name of the function you have created (anything beginning “<code>password_insert_...</code>” is guaranteed not to conflict with any existing function) and add the function somewhere in the file <code>admin-lib.inc.php</code>. The best place to add a function is right at the end of the file (that way, if you are using a checked-out copy of the code from the CWB Subversion repository, your modification to the code is likely to be preserved when you update CQPweb).</p>

2.3.7 User corpus system

Variable name	Description
---------------	-------------

Variable name	Description
<code>\$user_corpora_enabled</code>	<p>Type: Boolean</p> <p>Default: <code>false</code></p> <p>If this is set to true, the system allowing users to upload and index their own corpora will be enabled. This means it will be visible to users. However, users won't be able to actually use it until they have been granted the necessary privileges (to upload files, to use CorpusInstaller plugins, and to use disk space in the upload area and in CQPweb's index data area).</p>

2.3.8 RSS feed control

Variable name	Description
<code>\$rss_feed_available</code>	<p>Type: Boolean</p> <p>Default: <code>false</code></p> <p>If this is set to true, then an RSS 2.0 feed of all the system-administration messages currently on the system will be available; its location will be the rss subdirectory of your CQPweb base directory.</p> <p>(If the RSS feed is activated, an icon linking to the feed will appear in the header bar of the the system-messages block.)</p> <p>The next three configuration variables allow you to define how you want the RSS feed to behave; they will be ignored if <code>\$rss_feed_available</code> is not true.</p>
<code>\$rss_link</code>	<p>Type: String</p> <p>Default: (A URL corresponding to the root directory of the CQPweb installation.)</p> <p>All RSS feeds must have a URL associated with them. By default the CQPweb RSS feed's link is simply the root directory of the installation, but you can configure it to be something else by setting this variable. Note that if you set it to anything other than a valid URL the resulting RSS feed probably won't parse.</p>
<code>\$rss_feed_title</code>	<p>Type: String</p> <p>Default: "CQPweb System Messages"</p> <p>The title of the feed. You can set it to whatever you like, e.g. to mention your institution or organisation; it is useful to do so to disambiguate the RSS feed of one CQPweb server from another.</p>
<code>\$rss_description</code>	<p>Type: String</p> <p>Default: "Messages from the CQPweb server's administrator"</p> <p>The basic description that will pop up in subscribers' feed readers. You can set this to anything you like, but it should not be longer than a short paragraph in order to be effective.</p>

2.3.9 Error reporting

Variable name	Description
<code>\$print_debug_messages</code>	Type: Boolean Default: <code>false</code> If true , messages about what is going on behind the scenes will appear all over various CQPweb pages (totally messing up the layout, but hopefully helping you to diagnose problems).
<code>\$debug_messages_textonly</code>	Type: Boolean Default: <code>false</code> If you set this to true , all debug and error messages will be printed as text with no HTML formatting. Plain text messages are always produced (a) when outputting a textual download; (b) in scripts that run from the command-line.
<code>\$all_users_see_backtrace</code>	Type: Boolean Default: <code>false</code> If you set this to true , error messages will contain a PHP backtrace regardless of which user is logged in. By default, only admin accounts see the PHP backtrace.

2.3.10 Miscellaneous configuration options

The configuration options grouped in this section have not yet been sorted into larger groups; in future they probably will be.

Variable name	Description
<code>\$cqpweb_switched_off</code>	Type: Boolean Default: <code>false</code> If this is set to true , CQPweb appears switched off to the internet: it will never run, and any user who tries to see it will simply be shown the page <code>switched-off.inc.php</code> , located in the CQPweb <code>lib</code> folder. This page contains a default message, which you can amend by editing that file, if you wish; or you can add a short blob of extra HTML in the companion variable <code>\$cqpweb_switched_off_extra_message</code> . When CQPweb is switched off in this way, you can still run operations on the command-line. This setting is mainly of use to ensure nothing that users do can affect the database while upgrades, repairs, and so on are being run: the test to detect switched-off state runs <i>before</i> CQPweb makes a MySQL connection. It would be possible to end up with an inconsistent database state if users are able to access the interface and “do things” while an upgrade is in process.
<code>\$cqpweb_switched_off_extra_message</code>	Type: String Default: <code>""</code> You can set this to contain a short HTML sequence which will be included in the page presented to users when they attempt to access CQPweb at any time when it is switched off.

Variable name	Description
<code>\$cqpweb_root_url</code>	<p>Type: String Default: ""</p> <p>You can set this to an absolute URL (in the style "http://your-server.net/path/to/cqpweb/web/directory") and this URL will be used internally when CQPweb redirects the user's browser from one point in the system to another. If you don't set it, then CQPweb will try to work out its own URL based on information in PHP's global <code>\$_SERVER</code> array. That may produce incorrect results if you are running CQPweb in an environment where there is a lot of server proxying; if it does, things can be fixed by setting this variable appropriately.</p>
<code>\$cqpweb_no_internet</code>	<p>Type: Boolean Default: false</p> <p>If this is set to true, CQPweb assumes it is not available to the open internet and that it is only being accessed locally (i.e. by a web browser sending HTTP requests to localhost). This has two effects. First, CQPweb will disable all email-sending functions. Second, normal user account signup (which relies on email) will be turned off; that is, true here overrides <code>\$allow_account_self_registration</code>, setting it to false.</p> <p>It is convenient to set this as true if you have installed CQPweb for personal use on a desktop/laptop computer.</p>
<code>\$cqpweb_email_from_address</code>	<p>Type: String Default: ""</p> <p>An email address, which if provided will be used as the "From:" and "Reply-To:" address for all emails sent by the system. If this is not provided, defaults will be provided by your system's email-sending program - CQPweb will not attempt to supply a default. Be aware, however, that some mail systems will block emails that lack a clear "From" address, as they are assumed to be spam - so it is wise to set this.</p> <p>The email can be in the form of either a bare address (<i>someone@somewhere.net</i>) or a named address (<i>A. N. Other <someone@somewhere.net></i>).</p> <p>There are two basic possibilities here: you could set this to a real email address, so that if users reply the replies will go to some monitored mailbox; or, you could set it to an obviously fake address that (a) will signal to users they should not reply and (b) will ensure that, if they <i>do</i>, the replies go into a black hole: an example value for the latter would be "CQPweb Server <do-not-reply@cqpweb.anytown.edu>" or something along those lines.</p>

Variable name	Description
\$server_admin_email_address	<p>Type: String Default: ""</p> <p>An email address, which if provided will be exposed in the user interface to logged-in users as a designated means of contacting the server system administrator. If this is left as an empty string, all bits of text in the interface that would have specified a contact email address will simply be omitted. Email addresses on the web are often disguised (e.g. by spelling out “at” and “dot” instead of using the equivalent punctuation marks). It is not normally necessary for the email address in this variable to be disguised in this way, since only logged-on users can access the pages where it is displayed. However, you may wish to put an obscured email address here if you are running an open server.</p> <p>The email address is not converted into a link, so there are no restrictions on the format of the text. HTML is accepted and will be inserted into the user interface as you specify it.</p>
\$cqpweb_cookie_name	<p>Type: String Default: "CQPwebLogonToken"</p> <p>Label by which login cookies will identify themselves to users' browsers. The default value is normally fine, but if you have more than one installation of CQPweb running from the same web domain, you may find that their login cookies get confused with one another unless you set this variable to something different in each installation. The main logon cookie will use this name directly; other cookies will append a suffix. It is best to use just word characters (letters, numbers, underscore) for this setting.</p>
\$cqpweb_cookie_max_persist	<p>Type: Integer Default: 5184000</p> <p>This is the maximum length of time, in seconds, that a user's login will persist if they do not visit the site (the clock is “reset” every time they do visit the site). The default value is 60 days.</p> <p>This does not affect the option given to the user to choose whether or not they stay logged in; if they choose <i>not</i> to stay logged in, their browser will delete the cookie anyway and the persistent login will never be re-used.</p>
\$cqpweb_running_on_windows	<p>Type: Boolean Default: (see below)</p> <p>You can set this variable to true to declare that the operating system is Windows, or to false to declare that it is (some flavour of) Unix. If this variable is not set, CQPweb will use PHP's internal settings to guess the OS - so the only reason to use this variable is if CQPweb guesses wrongly on your system. (Moreover, since Windows compatibility is not yet implemented as of v3.2, this variable does not actually do anything yet.)</p>

2.4 Using the auto-configuration script

One of the administrative tools supplied as part of CQPweb is an auto-configuration script.

The script, also discussed in section 5.3, is an interactive tool for creating a basic configuration file. When you run the script, it will ask you a series of questions. Each question sets one of the compulsory settings (see 2.2).

To use the auto-configuration script, open a command-line terminal and go to the base directory of your CQPweb installation. **Then go into the `bin` directory, and enter the following:**

- `php autoconfig.php`

... and follow the instructions to enter the paths and other information that you made a note of when setting up the various directories, databases, etc. to be used by CQPweb (as discussed in chapter 1). When you are asked to specify admin usernames, enter at least one username (the one you personally will use).

Once you've answered all the questions, the script writes your answers to a configuration file in the correct location. If a configuration file already exists, the script will not overwrite it.

Once you've run this script, you can edit the resulting `config.inc.php` file to add any optional configuration variables that you might want.

2.5 Using the configuration file template

The file `template.config.inc.php` in the `lib` subdirectory of CQPweb is a blank template file which contains all the compulsory configuration variables. The easiest way to create a configuration file manually is using this template file, as follows:

- Make a copy of the `template.config.inc.php` file in the `lib` directory
- Rename the copy to `config.inc.php`
- Use a text editor to edit its contents; for each compulsory variable, insert the correct value for your system
- Add any optional variables you wish to use.

2.6 Changes from earlier versions of CQPweb

The configuration file format described in this chapter is that of version 3.2 of CQPweb. Version 3.1 made changes from the configuration file format used in 3.0 and earlier. If you have a configuration file from an earlier version, here are the things you need to know about the changes that have taken place.

(Not mentioned here: the new configuration variables, of which there are many, which have been added over time; for those, see the table of optional configuration variables in section 2.3.)

- The four compulsory variables that represent CQPweb's storage locations changed their format.
 - In 3.0, they were absolute paths, but with the initial `/` left off. This was a very non-standard way of specifying a filesystem path, and has been abandoned. In 3.1, these variables are treated as relative paths if they do not begin with `/`, and as absolute paths if they do.

- This means that if, for instance, you had the following in 3.0:
`$cqpweb_tempdir = 'var/cqpweb/temp';`
then in order for things to continue to work, you must change it to the following in 3.1:
`$cqpweb_tempdir = '/var/cqpweb/temp';`
- This affects `$cqpweb_tempdir`, `$cqpweb_upload_dir`, `$cqpweb_data_dir`, and `$cqpweb_registry`.
- `$path_to_cwb` and `$path_to_perl` changed as follows:
 - They became optional; if they are not set, then the CWB and Perl executables will be sought in the system path as per usual for programs whose precise location is not specified.
 - Their form was changed in the same way as the storage location variables: previously, they were absolute paths with the initial slash missing, now they are absolute *or* relative paths.
 - This means that if, for instance, you had the following in 3.0:
`$path_to_cwb = 'usr/local/bin';`
then in order for things to continue to work, you must change it to the following in 3.1:
`$path_to_cwb = '/usr/local/bin';`
(alternatively, if `/usr/local/bin` is on the path for the web-server user, as it usually would be, there is no need to specify this variable at all).
- `$path_to_apache_utils` was removed.
- `$password_more_security` was removed.
- `$cqpweb_uses_apache` was removed.
- `$utf8_set_required` was renamed `$mysql_utf8_set_required`.
- `$cwb_extra_perl_directories` was renamed `$perl_extra_directories`.
- `$default_mysql_process_limit` was renamed `$mysql_big_process_limit`.
- `$cache_size_limit` was given a new default value: 6GB rather than 3GB.
- `$mysql_freqtables_size_limit` was given a new default value: 6GB rather than 3GB.
- `$use_corpus_categories_on_homepage` was renamed `$homepage_use_corpus_categories`.

In version 3.2.11 and 3.2.12, the names of certain cache-limit variables were changed to make them more consistent.

The old names will still work until (at least) version 3.3.0.

- `$mysql_freqtables_size_limit` was renamed `$freqtable_cache_size_limit`.
- `$cache_size_limit` was renamed `$query_cache_size_limit`.

In version 3.2.32, all configuration variables relating to Perl were disabled and removed from this chapter (as the dependency on Perl has been ended to make for easier installation).

2.7 Obsolete feature: the corpus settings file

Prior to version 3.2, there existed an undocumented feature whereby any of the variables in the configuration file could be overridden on a per-corpus basis by setting a different value into `settings.inc.php` file for that particular corpus.

As of version 3.2, this is no longer possible.

3 The System Administrator's Interface

3.1 Introduction

The main user interface for system administrators is the *Admin Control Panel*. This contains controls for monitoring and managing many different aspects of CQPweb's behaviour, including indexing corpora, managing user access privileges, and tweaking the look-and-feel of the system.

The Control Panel can be accessed in three ways:

- From the **Admin control panel** link that appears in the side-menu for any corpus's main query screen *if* the logged-in user is an administrator.
- From the **Admin control panel** link on an administrator's user homepage.
- Via direct URL entry: add `/adm` to the base URL of your CQPweb system.

The Control Panel is laid out just like the main query screen and the user homepage: with a side-menu on the left and a main area displaying the selected function.

《add image here》

TODO

Different chapters of this manual explain different parts of what can be done with the control panel. The notation used throughout to refer to features of the control panel is as follows:

- **CP >XXX >YYY**

where “CP” means “go to the control panel”, from where the link for option YYY should be clicked, and this link is found under menu heading XXX. If any further links must be clicked to access some feature, more > are given added.

A general overview of the Control Panel is given in section 3.2.

A secondary interface exists for functions that affect just a single corpus. These are not accessed via the Control Panel, but rather through the corpus query menu.

When a system administrator is logged in, an extra menu section appears in the side-menu of the main query screen, with the heading *Admin tools*. This appears immediately before the *About CQPweb* menu section. The various menu items under *Admin tools* are discussed in section 3.3.

3.2 The Admin Control Panel: Feature list

There follows a listing of all features available from the menu in the Admin Control Panel, with cross-references to the other parts of the manual where discussion of those features can be found. If there is no discussion elsewhere in the manual, a brief explanation is given here.

- Corpora
 - *Show corpora*: list of all installed corpora, together with the size in types, tokens, and texts, and the disk space each uses; plus a link to the delete function. There are disk use totals at the bottom of the screen.
 - *Install new corpus*: see 6.10
 - *Manage corpus categories*: see 6.17

- *Annotation templates*: see [6.7](#)
- *Metadata templates*: see [7.6](#)
- *XML templates*: see [6.9](#)
- Uploads
 - *Upload a file*: This allows you to add files to the administrator's upload area, using a standard web upload form. For very big files, it is better to use an FTP or SFTP client or the like rather than the web form.
 - *View upload area*: list of uploaded files, with size/date info plus controls to view, compress/decompress, or delete a file.
- Users and privileges
 - *Manage users*: see [10.2](#)
 - *Manage groups*: see [10.4](#)
 - *Manage group membership*: see [10.4](#)
 - *Manage privileges*: see [10.5](#)
 - *Manage user grants*: see [10.6](#)
 - *Manage group grants*: see [10.6](#)
- Frontend interface
 - *System messages*: tool to add/remove messages from the list that appears on the homepage and on the standard query page.
 - *Skins and colours*: controls for working with the CSS files that govern CQPweb's appearance.
 - *Mapping tables*: see [《reference here》](#) TODO
- Cache control
 - *Query cache*: see [4](#)
 - *Database cache*: see [4](#)
 - *Restriction cache*: see [4](#)
 - *Subcorpus file cache*: see [4](#)
 - *Frequency table cache*: see [4](#)
 - *Temporary data*: see [4](#)
- Backend system
 - *Manage MySQL processes*: see [《ref todo》](#) TODO
 - *View a MySQL table*: debugging tool, allows you to print out the complete contents of any of tables in CQPweb's MySQL database.
 - *PHP configuration*: displays the PHP configuration settings that are of most relevance to CQPweb.
 - *PHP opcode cache*: tool to monitor and manipulate the PHP opcode cache (explained in the web interface itself)
 - *Public frequency lists*: see [《xref needed》](#) TODO
 - *System snapshots*: **under development, do not use**

- *System diagnostics*: tools to help diagnose problems with CQPweb
- Usage statistics
 - *Corpus statistics*: ranking of corpora by the number of queries run on them.
 - *User statistics*: ranking of users by the number of queries they have run.
 - *Query statistics*: list of the most frequently-run queries (across all corpora).
 - *Clear history*: a control allowing you to clear the Query History, on which the usage statistics are based (thus resetting all stats).
- Exit
 - *Exit to CQPweb homepage*: a link to the main homepage.

《Make some of the XREFS above more specific when those sections of the manual are finished》

TODO

3.3 Corpus Admin Tools: Feature list

The first item on this menu, **Admin control panel**, is simply a link to the Control Panel. The other items are discussed in turn below.

3.3.1 Corpus settings

This menu item takes you to a screen where you can configure a large set of miscellaneous options. Also to be found here are the controls for corpus-level metadata (《crossref》).

TODO

To modify any option, simply tweak its value in the interface then press the **Update** button.

The options are as follows:

- *Corpus title*: allows you to change the title you supplied when you indexed the corpus.
- *Directionality*: allows you change the value (left-to-right/right-to-left) that you supplied when you indexed the corpus.
- *Corpus requires case-sensitive collation*: see 《XREF to a part of the manual discussing collation & case-sensitivity》
- *Stylesheet address*: allows you to change the CSS file location from what you supplied when you indexed the corpus.
- *Amount of context shown in concordance*: you can specify the width of the concordance in either words, or relative to any XML element.
- *Initial/Maximum words in extended context*: you can specify more/less words here based on, for instance, whether or not your users have the necessary permissions to see extended excerpts of the corpus data.
- *Word annotation for alternative view*: this is explained in section 9.2.
- *Corpus category*: see section 6.17.
- *Visibility of the corpus*: *visible* corpora are listed on the homepage; *invisible* corpora aren't. However, an invisible corpus **is** listed on the “Corpus permissions” page of any user who has been granted the privilege to use that corpus.

TODO

- *External URL*: if you specify a link here, it will be embedded in the left-hand menu as a link under the **Corpus info** heading. The idea is that you would provide here a link to some online documentation about the corpus, if any such exists, to allow users to go straight to the relevant information.
- *Primary text categorisation*: this is explained in the chapter on metadata, specifically section *«XREF»*.

TODO

Below these options are the corpus metadata controls; see *«crossref»*.

TODO

3.3.2 Manage access

This menu item presents some information relating to the access that users have to the corpus.

The corpus-access privilege system is explained in section 10.5.1. Privileges cannot be manipulated from here - this must be done via the Admin Control Panel (links are provided to the appropriate screens in the CP).

However, here you can find two things that are not always easy to spot within the full list of privileges:

- A list of all privileges that affect access to the the corpus, together with a list of the groups and individual users to which those privileges are granted.
- A full combined list of all users with access to the corpus, whether they have it from a membership in one or more groups, or as individuals; duplicates are removed (that is, if Group A and Group B both have access to the corpus, and User X is a member of both A and B, you will see User X listed here only once).

3.3.3 Manage text metadata

«prob just an xref?»

TODO

3.3.4 Manage text categories

3.3.5 Manage corpus XML

3.3.6 Manage annotation

«XREF to section in the indection chapter on linking CEQL to annotation, or XREF from that to here?»

TODO

3.3.7 Manage frequency lists

3.3.8 Manage visualisations

The controls found in this section are explained in chapter 9.

3.3.9 Cached queries

3.3.10 Cached databases

3.3.11 Cached frequency lists

4 Managing the CQPweb data cache

4.1 Introduction

CQPweb is built around a strategy of extremely aggressive caching of dynamically-generated data.

The core CWB system does not employ caching. In general, any part of the corpus data is only stored *once*, in the CWB index; whenever a query is run, or data is requested in any other format, it is retrieved anew from the CWB indexes - even if the same data has been requested recently. (That said, CWB benefits from the fact that modern operating systems, and especially the Unix systems that it is primarily designed for, cache disk content in RAM automatically.) CWB indexes are designed to be maximally compact on disk, even if this complicates or slows down retrieval.

By contrast, CQPweb is designed with the assumptions that (a) the same requests will often be made many times in a row, and (b) speed of response to requests is the most important thing (far more important than minimising the use of disk space). These assumptions are rooted in its origin as a teaching tool: a common use-pattern for CQPweb is a roomful of students, working on the same data, and all doing pretty much the same kinds of queries. In this situation, caching all generated data (queries, ad hoc frequency lists, and more) leads to a substantial performance improvement, as resource-intensive processes only run *once*. For every user other than the first to run a particular process, response time is much faster; the dynamically-generated data does not need to be rebuilt from the underlying corpus indexes, it merely needs to be retrieved from the cache.

This chapter explains the different kinds of data that CQPweb caches, and discusses aspects of cache administration that apply especially to large, multi-user installations.

4.2 Some background on the MySQL database system

CQPweb uses MySQL as its secondary datastore: that is, while the CWB indexes contain the main corpus data, all ancillary data (corpus and text metadata, frequency lists, analysis data for collocation/distribution/etc., cached queries and analyses, user data like categorised queries, as well as CQPweb's system management data) is stored in a MySQL database.

To understand the ways this database - especially its caches - can be configured, some background on MySQL is needed. Beware! This is a very incomplete account of what is actually going on in MySQL; for the full details, see the MySQL manual. This section only covers the points relevant to CQPweb.

MySQL is a client-server database management program: a server or “daemon” called `mysqld` does the actual work of creating, modifying and searching the databases; the daemon is always accessed via a client program which interacts with the user, issues requests to the daemon, and handles the responses sent back by the daemon. CQPweb itself acts as a client to the MySQL daemon, but some actions involved in administering CQPweb involve accessing the daemon via the MySQL command-line client (called `mysql`). An example is the process of creating CQPweb's MySQL database in the first place, as described in an earlier chapter.

Most of the time, a database can be treated as an abstract entity which we control without worrying about the details of what the daemon is actually doing. However, when considering the performance issues that arise on large, heavily-used installations of CQPweb (and similar programs) it *is* necessary to dig in to what the MySQL daemon is doing behind the scenes.

The MySQL system stores all the actual data using one of two “engines”.¹ The older of these is **MyISAM**; the newer is **InnoDB**. Each separate table within a database can be set to use a different engine. The databases and tables are ultimately represented by various disk files (though while the

¹There are actually a whole lot of engines, but MyISAM and InnoDB are the only two that matter in practice.

MySQL daemon is running, there is not necessarily a guarantee that all important information has actually been written to disk at any given moment!), as follows:

- The MySQL daemon has a dedicated directory on disk for its data (its location is set in the daemon's configuration options; the default location varies across operating systems, but on many Linux systems it is something like `/var/lib/mysql`).
- For each database that is created, a folder is created with the same name as the database, directly within that dedicated directory. So, if we create a database called **cqpwebdb**, its folder would be `/var/lib/mysql/cqpwebdb` (with the first part, again, depending on the OS).
- For each table within the database, a `.frm` file is created, which contains the table definition, but not the actual data.

The location of the actual table data depends on the engine; furthermore, InnoDB can locate its data in two different places depending on whether a mode called *file-per-table* is switched on or not.

- **MyISAM** stores actual table data in the same place as the `.frm` file, in two extra files with the same name, but different file extensions (`.MYD` and `.MYI`). So, for instance, if the table **corpus_info** uses the MyISAM engine, the database directory will contain `corpus_info.frm`, `corpus_info.MYD` and `corpus_info.MYI`.
- **InnoDB** with file-per-table mode switched **off** stores all the data for all InnoDB tables from all the different databases in a single location called the “global tablespace”. This takes the form of a single file called `ibdata1`, which will be found directly under MySQL's main data directory. The `ibdata1` file will therefore be very large!
- **InnoDB** with file-per-table mode switched **on** stores the data for each table in a separate `.ibd` file, which is placed (by default) alongside the `.frm` file - although it is possible for a separate location for the `.ibd` file to be specified when the table is originally created. With this setup, there will still be an `ibdata1` file but it will be much smaller as it will not contain actual table data.

MyISAM was the default in older versions of MySQL; later, InnoDB became the default. Initially, InnoDB's default setting was to have file-per-table switched off; in the very most recent versions of MySQL, file-per-table mode is on by default. It is generally agreed that InnoDB with file-per-table switched on is the best way to store table data in most cases, if you are using an up-to-date version of MySQL.

CQPweb is, in general, designed to work with the default settings of MySQL. What that means is that originally it was built with the assumption that most tables would be MyISAM; later it was modified to force most tables to be in InnoDB, once InnoDB became preferred. It did not, however, make any difference to CQPweb whether file-per-table was on or off until version 3.2.14, when it became possible (as described below) to store different types of data in different locations. This new functionality depends on file-per-table mode being switched *on*.

To find out what the default engine is on your system, enter the following command in the MySQL client:

- `show engines;`

This will display a list of available engines with additional information. The key column is “Support”: in this column, the word `DEFAULT` will appear next to one of MyISAM or InnoDB.

To find out whether your MySQL daemon has the InnoDB file-per-table mode switched on, enter the following command:

- show variables like "innodb_file_per_table";

However, it's important to note that these settings only represent the present situation. If your CQPweb server has been running for a while, especially if it's been running for over one year, it's quite possible that these settings could have been different in the past. In that case, any data that was setup under the old settings will not have been updated to the new settings.

- If your CQPweb server used to run on a version of MySQL where MyISAM was the default engine, any old tables created as MyISAM will still be MyISAM.
- On the InnoDB side, if file-per-table mode is switched on but used to be switched off, any tables created in the global table space will still be located there.

In either case, your database will contain tables running on a hodgepodge of engines and table spaces. An administration script has been provided that will fix these issues²: see 5.7 on `force-innodb.php`.

4.3 Explaining the different types of cached data

Let's consider the data cached by CQPweb under broad categories

《explanations of the different caches》 TODO

《explanations of what is cached in corpus setup》 TODO

《on the restriction cache: maybe add note - more RAM may be needed if there are restrictions based on XML elements.》 TODO

4.4 Disk locations for stored data

《discussion of how the locations of these caches affects performance - multiple disks etc. Advising on partitioning etc》 TODO

《the config settings to use for each location》 TODO

²There is an important thing to remember if you formerly used InnoDB without file-per-table, and have switched to InnoDB with file-per-table. This is that forcing all the existing tables out of the `ibdata1` file and into their own files *will not result in the `ibdata1` file getting any smaller* - since it *never* gets smaller.

So, for instance, if you have a 100GB `ibdata1` file, then switch on file-per-table and force the tables out into their own file, the result will be a 100GB `ibdata1` file *as well as* 100GB of separate table files.

The *only* way to reclaim disk space from the `ibdata1` file is to export all your MySQL databases using the `mysqldump` program; drop all the databases from MySQL; delete the `ibdata1` file; and then re-insert the databases from the exported file.

This is, it should go without saying, a somewhat dangerous process. It's also outside the scope of this manual, so no more will be said about it here - but see the following links for some discussion:

- Relevant pages in the MySQL manual (v5.7, other versions linked from there)
 - <http://dev.mysql.com/doc/refman/5.7/en/using-innodb-tables.html>
 - <http://dev.mysql.com/doc/refman/5.7/en/innodb-multiple-tablespaces.html>
 - <http://dev.mysql.com/doc/refman/5.7/en/innodb-configuration.html>
 - <http://dev.mysql.com/doc/refman/5.7/en/mysqldump-sql-format.html>
 - <http://dev.mysql.com/doc/refman/5.7/en/reloading-sql-format-dumps.html>
- Discussions on StackOverflow and StackExchange/DBA:
 - <http://stackoverflow.com/a/4056261>
 - <http://dba.stackexchange.com/a/24963> - on the general process
 - <http://dba.stackexchange.com/a/2227> - on exporting the databases

(all links correct as of April 2016).

«how do we move one of the caches once it's already in a place?»

TODO

It's important to note that setting the location of one of the MySQL-based caches will only work if you have InnoDB's file-per-table mode switched **on**; see section 4.2.

<https://dev.mysql.com/doc/refman/5.7/en/tablespace-placing.html>

4.5 Moving the cache location on an existing CQPweb server

The process for moving a cache folder is different for caches involving files and caches involving MySQL tables.

https://blogs.oracle.com/mysqlinnodb/entry/choose_the_location_of_your

4.6 Optimising MySQL for cache performance

«this might get hived off into a seaprate chapter on optimisation since it's a bit of a shift of topic for this one»

TODO

«this section currently just contains rough and badly typed notes»

TODO

Lots of useful stuff here: <http://dev.mysql.com/doc/refman/5.7/en/converting-tables-to-innodb.html>

Cover ehre things like this:

`innodb_flush_method=O_DIRECT innodb_log_file_size=1G innodb_buffer_pool_size=4G`

flush method – setting this to `O_DIRECT` means that read/write access is direct to disk without using the OS's disk cache in RAM. This apparently improves performance if MySQL has lots of RAM to cache disk stuff in, (in which case the OS cache will be small and twatty.)

It is named after the Linux flag that it causes `mysqld` to use with innodb files. Here is a quote from `man open(2)`: `O_DIRECT` (Since Linux 2.4.10) Try to minimize cache effects of the I/O to and from this file. In general this will degrade performance, but it is useful in special situations, such as when applications do their own caching. File I/O is done directly to/from user-space buffers.

`buffer_pool_size` - this is InnoDB's internal cache of stuff from disk. The default is 128MB. The manual says that on a dedicated box, you'd set it to 80% of the available RAM. The persona performance blog says that ideally you'd have as much RAM in the buffer pool as the DB takes on disk. Obvs in the case of CQPweb that's a no go. But it can certainly be more than 128MB. Given the RAM we have, I gave it 2GB for now

`innodb_log_file_size` – the MySQL manual says it can be up to 1/2 of the buffer pool.

THERE SEEM TO BE TWO CHOICES..... 1 Don't use `O_DIRECT`, use a small buffer pool - InnoDB does not do RAM caching of disk stuff, we rely on OS caching. 2 Use `O_DIRECT` and a moderately hefty buffer - InnoDB will cache some of the disk stuff in its buffer pool and we bypass the OS cache.

But note this helps READS more than it helps WRITES I think!

Discuss relationship of this to OS RAM caching of the CWB index files.

SEE: <http://dev.mysql.com/doc/refman/5.7/en/disk-issues.html>

«end of rough and badly typed notes»

TODO

4.7 User-data cache sizes

《*recommendations on disk space to allocate: explain the default size allocation of 6GB per cache, and point out that a single-user installation might well want to drop this.*》 **TODO**

4.8 Finding and fixing cache leaks

A *cache leak* is the term we use for any situation where something goes wrong in the creation or management of cached data or temporary data such that CQPweb can no longer control the data and delete it where appropriate.

This can happen, for instance, if CQPweb is interrupted halfway through running. For instance, if CQPweb aborts halfway through a cache cleanup, it is possible that the record of some block of cached data might be deleted but the block of data itself retained (or vice versa).

In either case, the data becomes “invisible” to CQPweb’s normal cache-management procedures, but will still be present, not being used but taking up disk space. This can have one of two negative consequences:

- Leaked data blocks can fill up your disk space if they have completely escaped monitoring.
- Conversely, if the leaked data blocks are still counted towards the size of your cache, they constitute a fraction of the cache that can never be used by actual productive data.

The very nature of “leaked” cache data or cache records, with imperfect or incomplete tracking, means that it would not be safe for CQPweb to attempt to clean them up manually - there would be a risk of data loss. However, tools are provided to assist you in cleaning up cache leaks.

To monitor and manually delete leaked items, go to one of the five functions at **CP >Cache control**. There is one for each type of cached data. Each is slightly different depending on the nature of the cache it covers. Each of these functions also gives you some information about the size of the cache and how full it is. Once a CQPweb server has been running for a while, most of the caches will become about 100% full and stay at that level. This is the expected behaviour - it means you are getting the maximum benefit from the disk space allocated to these caches.

The cache control functions often take quite a long time to display. This is because generating information about possible leaks involves searching the whole of the cache. So don’t worry about this.

5 Administering CQPweb from the commandline

5.1 Introduction

Several actions that you can take as a system administrator are more conveniently undertaken outside the web interface, or else cannot be exposed in the web interface for reasons of security. These actions are instead available from the commandline using scripts in CQPweb's `bin` directory (one of the subdirectories from the base CQPweb directory).

All the scripts in this folder are straightforward CLI scripts written in PHP. Some are interactive (that is, they will require user interaction as they run); others will simply run on their own and then finish.

Some of these scripts run within a complete CQPweb environment; that is, they read in the configuration file, create a connection to the MySQL server, and so on. For that reason, they all need to be run from within one of the subdirectories of the base CQPweb directory. If the script is intended to operate on some specific corpus, then you should change your working directory to that corpus' directory, and call the script as follows:

- `php ../bin/name-of-script.php`

Alternatively, for scripts that do not refer to a particular corpus, your working directory should be the `bin` directory itself. The script is then called as follows:

- `php name-of-script.php`

Scripts which write files (e.g. `autoconfig.php`, `load-pre-3.1-privileges.php`) need to have write access to the CQPweb directory. You have three choices:

- Run them as the username under which the webserver runs (i.e. the user that owns the CQPweb directory);
- Run them as some other username, making sure that this user has write access to the CQPweb directory;
- Run them as root, changing ownership of the resulting files to the webserver user afterwards.

In the remainder of this chapter, the details of each of the scripts are explained.

5.2 The main *cqpweb* script

This is an executable that allows you to call any function from CQPweb's internal function library. If you wish to call a function that relies on being run in the environment of a specific corpus, run the script from that corpus' folder; otherwise run it from `bin`.

The syntax is as follows:

- `./cqpweb NAME_OF_FUNCTION ARG1 ARG2 ...`

It is difficult to provide any general comments since this script can do almost anything. Here are some examples of useful calls.

- `./cqpweb add_corpus_to_privilege_scope PRIVILEGE-INTEGER-ID CORPUS-HANDLE`
- `./cqpweb remove_corpus_from_privilege_scope PRIVILEGE-INTEGER-ID CORPUS-HANDLE`
- `./cqpweb create_corpus_default_privileges CORPUS-HANDLE`
- `./cqpweb add_new_privilege 1 "" "Permission to use at retriected level (initially has scope over no corpora, they can be added later)"`
- `./cqpweb add_new_privilege 2 "" "Permission to use at normal level "`
- `./cqpweb add_new_privilege 3 "" "Permission to use at full level "`
- `./cqpweb add_new_privilege 4 5000000 "Permission to create freq lists up to 500K tokens"`
- `./cqpweb grant_privilege_to_user USERNAME PRIVILEGE-INTEGER-ID`
- `./cqpweb grant_privilege_to_group GROUP-NAME PRIVILEGE-INTEGER-ID`
- `./cqpweb remove_grant_from_user USERNAME PRIVILEGE-INTEGER-ID`
- `./cqpweb remove_grant_from_group GROUP-NAME PRIVILEGE-INTEGER-ID`
- `./cqpweb update_corpus_visualisation_gloss CORPUS-HANDLE
1-OR-0-FOR-SHOW-IN-CONCORDANCE 1-OR-0-FOR-SHOW-IN-CONTEXT P-ATTRIBUTE-HANDLE`
- `./cqpweb update_corpus_visualisation_translate CORPUS-HANDLE
1-OR-0-FOR-SHOW-IN-CONCORDANCE 1-OR-0-FOR-SHOW-IN-CONTEXT S-ATTRIBUTE-HANDLE`
- `./cqpweb add_variable_corpus_metadata CORPUS-HANDLE ATTRIBUTE-DESCRITPION
VALUE-CONTENT`
- `./cqpweb update_corpus_title CORPUS-HANDLE "new title goes here"`

5.3 autoconfig.php

Creates a configuration file with just the essential configuration variables.

This script is interactive, that is, it asks you questions and requires you to type in the configuration values that you want to appear in the finished file.

This script is also discussed in section [2.4](#).

Note that prior to version 3.1, this script *also* performed a variety of miscellaneous set-up actions. These are now performed by `autosetup.php`.

5.4 autosetup.php

This script is used in the process of setting up a new CQPweb installation. It is described in section [1.12](#).

Like `autoconfig.php`, it is interactive.

5.5 cli-lib.php

This is not a script, it is a library file used by the other commandline scripts. It is mentioned here simply so that you know not to try and run it! (Nothing bad will happen if you do, in fact nothing *at all* will happen if you do.)

5.6 execute-cli.php

This script allows you to call any function from CQPweb's internal function library. If you wish to call a function that relies on being run in the environment of a specific corpus, run the script from that corpus' folder; otherwise run it from `bin`.

This is the backend behind the `cqpweb` executable (see 5.2). It is used in exactly the same way.

The syntax is as follows:

- `php execute-cli.php NAME_OF_FUNCTION ARG1 ARG2 ...`

It is difficult to provide any general comments since this script can do almost anything. Needless to say, you *really* need to know exactly what you are doing before you start messing with this script. Caveat emptor.

5.7 force-innodb.php

This script forces all tables in CQPweb's MySQL database to use the InnoDB engine, instead of the older MyISAM engine. For an explanation of these engines, see 4.2.

The script operates by running the following command on each table:

- `alter table name_of_table ENGINE=InnoDB`

The effect of this command is discussed in the MySQL manual here <https://dev.mysql.com/doc/refman/5.7/en/alter-table.html> and here <http://dev.mysql.com/doc/refman/5.7/en/converting-tables-to-innodb.html>.

Run this script in the following situations:

- Your CQPweb installation is an old one, from before CQPweb began enforcing the use of InnoDB, and contains MyISAM tables.
- You have switched your MySQL server from using the InnoDB global tablespace to using file-per-table mode and want to force your existing tables out into their own files.

Certain tables require support for fulltext indexes, which was added to InnoDB only in 2013. If you are using an old version of MySQL whose InnoDB engine lacks this feature, any tables that need it will not be touched by this script; they will be “held back” as MyISAM.

5.8 install-corpus.php

This script installs a new corpus. It has an internal manual explaining how to call it, which you can see by running:

- `php install-corpus --help | less`

(for “less” substitute *Your Favourite Pager*)

Fundamentally, each of the options to the script replicates one of the form elements on the relevant page in the Admin Control Panel.

5.9 load-pre-3.1-groups.php

Prior to version 3.1 of CQPweb, information about what user groups exist, and what users belong to each, were stored in Apache `.htgroup` files (and if you were not using Apache, you were out of luck). In version 3.1, this was changed so that groups are instead stored in the database, and Apache is no longer necessary: CQPweb can manage its own user accounts while running on any webserver.

This script migrates groups automatically from the old system to the new system. You should only need to run it once, immediately after you upgrade the code and database to version 3.1 (see section 14.3).

This script will ask you to specify interactively the location of the group file. This is the directory which will have been set as `$cqpweb_accessdir` in your pre-3.1 configuration file.

5.10 load-pre-3.1-privileges.php

Prior to version 3.1 of CQPweb, corpus access privileges were stored in `.htaccess` files in the web directory of each corpus. In version 3.1, this was changed so that privileges and grants are stored in the database, a much more flexible system. However, this means that lists of permissions in an existing v3.0 installation would be lost. This script retrieves them.

The script creates the three default privileges for each corpus on the system, if they do not exist already (access levels normal, restricted and full).

The “normal” access level is intended to be equivalent to the one-and-only access level available in older versions of CQPweb. For this reason, every group which *previously* had *any* access rights for a given corpus is assigned “normal” access to that corpus by this script.

It should be run *after* `load-pre-3.1-groups.php`.

5.11 load-pre-3.2-corpsettings.php

The use of this script is explained in 14.6.

It should be run *after* upgrading the database as far as v3.2.0.

If the version you are upgrading to is higher than 3.2.0, you need to run the database upgrade script a second time after running this script.

5.12 offline-freqlists.php

This script should be called as follows:

- `php offline-freqlists.php lowercase_name_of_corpus`

The script performs all frequency-list setup functions, i.e.e those actions usually performed on the “Manage metadata” page of the corpus interface, after indexing a corpus and installing its text metadata table (*⟨XREF⟩*).

For very big corpora (hundreds of millions of words or more) this process can take a long time, hours or even a day or more. In that case it is not convenient to run it from a web browser, and the commandline version may make more sense.

This script prints a long stream of debug messages to indicate its progress. These can be safely ignored until and unless something goes wrong. This is not an interactive script and no user attention is required.

TODO

5.13 upgrade-database.php

From time to time the CQPweb database format is changed. This script should be run after you update your code to a new version, and it will implement any necessary database version changes.

This is an interactive program - it will sometimes demand that you acknowledge some alert about something that has changed by pressing Enter before it will continue.

Note that if you upgrade the code, but do not run this script to bring the database into line, CQPweb will break (possibly very badly).

The biggest set of upgrades are those between version 3.0.16 and 3.1.0. However, any database upgrade can potentially take a long time.

If you are running an online server (as opposed to one on a standalone computer!) then it is recommended to take the server offline while you do the upgrade - so that users cannot access your server while it is in a half-and-half state (which, depending on what they do, might cause data integrity problems). The easiest way to do this is simply to disable your web sever software temporarily.

See also [section 14](#) for more information on upgrading.

6 Indexing corpora

6.1 Quick checklist

This chapter is not complete. Until it is, users who are not familiar with the corpus indexing process may benefit from using this checklist (courtesy of Philipp Heinrich on the CWB Developers Mailing list).

To create a new corpus, use the following tools in turn:

- In the Admin Control Panel:
 - install new corpus
- Inside the new corpus's UI:
 - manage text metadata
 - (build CQPweb frequency lists either using “manage frequency lists” or the offline script)
 - manage corpus XML
 - manage visualizations
 - manage annotation - setup CEQL bindings
 - check corpus settings
- In the Admin Control Panel:
 - manage privileges
 - manage user / group grants

6.2 Basic concepts

CQPweb is based, naturally, on CWB corpus indexes. Setting up a corpus for use in CQPweb essentially means indexing a CWB corpus. However, the design of CQPweb means that it has some specific requirements for the layout of its corpora beyond the basics of CWB. This section is intended to explain CQPweb corpora to help you prepare appropriate data for input.

A corpus in CQPweb can be characterised in terms of the combination of its *texts*, its *annotations*, its *XML*, and its *metadata*.

- *Texts* are the fundamental organisational unit of CQPweb corpora. They are one type of XML element, but one that must always be present. Every corpus consists of a sequence of one or more texts. Texts cannot overlap, or be contained within other texts. All the words in a corpus must be contained in one of that corpus's texts.
- *Annotations* are the different layers of word-level information. This corresponds to the more general CWB notion of a p-attribute (positional attribute).
- A corpus' *XML* corresponds to the more general CWB notion of an s-attribute (structural attribute).
- *Metadata* is a set of structured information about some aspect of a corpus, typically its texts but potentially other kinds of XML as well. Since metadata is a major topic it will be dealt with in a separate chapter (7).

When indexing a corpus, it is necessary to specify the corpus' design in terms of annotation, XML and metadata, and to make sure that the corpus data to be indexed contains appropriate text tags. Annotation, metadata and XML can be specified from scratch, or their design can be taken from a predesigned *template*.

6.3 The notion of a handle

A *handle* is a short text label used to refer to any of the following entities in CQPweb:

- A corpus.
- An ID code for a text or an XML region.
- An annotation layer.
- A type of XML.
- A field in a metadata table.
- User account usernames.
- A user-specified save-name.

This notion of handle may be familiar as it is fundamentally similar to the labels used for corpora and (p- and s-) attributes in CWB/CQP. However, CQPweb enforces certain requirements for what is and is not a valid handle.

- Each type of handle is limited to a certain, short length (see below)
- Handles can only include certain characters: ASCII letters, ASCII digits, and the underscore character (i.e. the same rules as for “words” in regular expressions, or for identifiers in C and related programming languages).
 - Note that this is different to command-line CWB, which allows the hyphen to be part of an attribute or corpus name in addition to the characters mentioned above.
- Some types of handle (for corpora, annotation/p-attributes, and XML/s-attributes) cannot include any uppercase letters.

It is necessary to limit the length of handles because of features of the MySQL database system. When CQPweb installs a corpus, it creates a number of database tables - and handles associated with the corpus are used to generate many of the table/column namers. Other handles are used as database keys (and some database keys in CQPweb include up to three handles).

However, MySQL has two important built in limits: (1) Table and field names cannot be longer than 64 characters. (2) Database key fields cannot be longer than 333 letters. CQPweb has to limit the length of handles to fit with these constraints. For this reason, there are four broad categories of handle within CQPweb:

- Handles which can be up to 20 characters in length (because they are used as part of a longer table name). This limit is applied to:
 - Corpus handles.
 - Annotation handles (p-attributes).

- Handles which can be up to 64 characters in length (because they may form the name of a table column, but will not be embedded in a longer table name). This limit is applied to:
 - Handles for XML elements or attributes (s-attributes).
 - Handles for metadata field names (any type of metadata).
 - Usernames.
- Handles which can be up to 200 characters in length (because they will never be used as part of a table name, but might be used as part of a longer database key). This limit is applied to:
 - Handles for the categories in classification-scheme metadata.
 - The save-name of a saved query or subcorpus (but see note below about categorised queries).
- Handles which can be up to 255 characters in length (because they might be used as a database key on their own). This limit is applied to:
 - ID codes for texts and for XML elements.

Note that categorised queries, although they are a type of saved query, cannot be given names as long as those of a normal saved query (200). This is because, when a categorised query is separated out to create a saved-query from each of the categories of concordance line, the new name contains the name of the categorised query combined with the name of the category - so the total length of both combined can be no more than 200.

In versions of CQPweb prior to v3.2.0, usernames were limited to 30 characters; they can now be up to 64 characters, as noted above (avoiding the need for a fourth kind of handle).

6.4 File format

«Input file: do by XREF to other docs if possible - the CWB encoding tutorial frinstance.»

TODO

6.5 Linking handles and descriptions

«Explain that wherever there is a handle, it can become a description.»

TODO

6.6 Annotation

«primary annotation - do we introduce this here, or do we leave it for the section on CEQL?»

TODO

6.7 Annotation templates

An **annotation template** is a data structure that describes a set of annotations (that is, p-attributes, corresponding to columns in the input data).

If you wish to index multiple corpora that have the same annotations, then rather than specify the details of each annotation every time, you can create and save a template that stores those details and can then be invoked when a corpus is created.

For instance, one very common pattern of annotations is for the second column of the input file to contain a part-of-speech tag, while the third contains a lemma (recall that the *first* column always contains the wordform of the token). This three-column format is produced by, among others, the **TreeTagger** software.

Rather than enter `pos` and `lemma` over and over again when setting up corpora, you can create a template description which describes this three-column format. When you set up a corpus of this kind, specifying its structure is then as simple as picking your prespecified template.

(In fact, certain useful templates, including the three-column *word*, *pos*, *lemma* format, are actually built into CQPweb for you.

《NOTE: explain template creation here by saying the form is the same as the one for corpus setup. **TODO** BUT don't go into the forms for annotation here. that comes under "indexing-proc" below.》

6.8 XML

(Include here the discussion of "text" as the compulsory element and XREF to "text" as discussed in "basic concepts".)

6.9 XML templates

6.10 The indexing process

(all about the form)

(inc use of the non-template forms)

6.11 Using a pre-indexed corpus

《Turn following rough notes into real manual content》

TODO

Step 1 - index the corpus using command-line CWB (wherever you like on the system, as long as the files/directories you create are in a location on the file system where the web server's user account has permission to read them)

Step 2 - go to the "Install new corpus" page in CQPweb, and click on the link at the top that says "Click here to install a corpus you have already indexed in CWB."

Step 3 - specify the location of the registry file. (this will be copied into CQPweb's own registry if not already there; the index files themselves will not be copied or moved.)

Step 4 - once you've installed the corpus thusly, proceed onto the other installation steps (generate your text metadata from the XML attributes on `!text!`, or else install a metadata file; setup frequency lists; etc.)

6.12 The metadata setup process

《short explanantion here only – direct to the next chapter for the full coverage of metadata》

TODO

6.13 Building frequency lists

Frequency list setup is dependent on the existence of the text metadata table (see (《XREF》)).

TODO

Do automatically small corpora

Do bit by bit in web interface

Do offline with script

6.14 Linking annotation to CEQL syntax notation

《*explain this!*》

TODO

6.15 Setting up corpus access rights

When a corpus is initially set up, no users except the system administrator(s) will be able to access it. To enable access, you need to (a) create a *privilege* that covers use of the corpus, and then (b) grant that privilege to one or more users or groups.

(xref to “user accts” chapter...)

6.16 Further corpus configuration

There are many settings that can be configured for each corpus after it has been indexed. They are available via the individual corpus's interface (*not* via the Admin Control Panel). When you are logged on with an administrator account, an additional section will appear in the left-hand-side menu, labelled **Admin tools**.

Some of the options on this menu have been discussed already in this chapter, as they are involved in the setup procedure. Others are discussed elsewhere: for the **Manage visualisations** option, for instance, see chapter 9. For a full overview, see section 3.3.

6.17 Putting corpora into categories

An optional step in setting up a new corpus is putting it into a specified *category*.

Currently, what category a corpus is in affects only one thing: the layout of the list of corpora on the homepage. If category-based ordering of the homepage is switched on, then the list will be divided into sections based on the corpus category feature. What category a corpus is in will determine where it appears in the list.

See section 2.3.5 for how to set the right configuration variable to switch on category-based organisation of the home page.

Categories are created and managed in the Admin control panel: **CP >Corpora >Manage corpus categories**. The upper part of this screen shows the existing categories. The most important feature of a category is its sort order, which is represented as an integer, and which determines where it appears on the homepage. The sort order is a relative number - that is, it doesn't matter exactly what the value is, it matters what it is relative to the other categories. The **[Move up]** and **[Move down]** controls can be used to adjust the sort order. Categories can be added using the separate control on the lower part of the screen.

In the **Corpus settings** screen for any individual corpus you will find an option to change the category of that corpus. When they are first indexed, all corpora are placed in the default “Uncategorised” category, and will stay there until you move them.

Within each category, corpora are listed in alphabetical order of the corpus handle - note this is not necessarily the same as the descriptive name, which is what is actually displayed!

Finally note that if a corpus is set to be *invisible*, it doesn't matter at all what category it is in.

7 Metadata

7.1 Introduction

《(introductory explanation of what it is, how it is stored, where it appears in the interface)》

TODO

In CQPweb, *metadata* is a covering term for data about an indexed corpus, about the texts in a corpus, or (sometimes) about individual instances of XML elements.

The underlying CWB system does not provide an easy way to store and manipulate metadata. For that reason, CQPwebh uses its MySQL database to store and manipulate all metadata.

Metadata, espically text/XML metadata, is crucial to several core CQPweb functions:

- *Distribution.* 《EXPLAIN》
- *Restricted queries.* 《EXPLAIN》
- *Subcorpus creation.* 《EXPLAIN》

TODO

TODO

TODO

This chapter explains the different levels of metadata (corpus/text/XML), the different datatypes that are available, and how they work; how metadata is installed for a corpus; and the use of metadata templates.

7.2 Corpus metadata

.....

(explain corpus metadata, where it is displayed, how to add it)

.....

7.3 Text metadata

.....

Each corpus stores text metadata in a dedicated *metadata table*.

A metadata table is a database table very much like a table you might create in a spreadsheet. It has a series of columns, where each column represents a particular *field* or *attribute* - that is, a particular bit of metadata. So, for instance, yhou might have columns for the *auhtor*, *title* and *genre* of each text.

The rows represent the items described by the metadata table. The first column always contains the unique identifier for the items. So, in a text metadata table, the first column contains the text ID code. In CQPweb, text IDs and other identifier codes are always handles (see 《XREF》).

TODO

Metadata tables can be loaded into CQPweb from uploaded plaintext files. Alternatively they can be generated from data already present within the indexed corpus. See 《XREF》 below.

TODO

7.4 XML metadata

(Explanation of how XML metadata can be applied in different ways)

...

Unlike text metadata, XML metadata does not normally use a metadata table. So to include XML metadata, you would normally incorporate it directly into the XML tags in the underlying corpus as represented in your input file.

The exception to this is metadata for ID-linked XML attributes. Attributes of this kind which *do* use a metadata table, which thenrefore works in a similar way to the text metadata.

This is explained below XREF.

7.5 The different possible datatypes

Currently, CQPweb supports the following datatypes for text/XML metadata fields.

- Free text
- Classification
- Unique ID
- ID-link
- Date (currently under development)

Note that all *corpus* metadata items are necessarily free text.

7.5.1 Free text



TODO

Free text is the most basic kind of metadata. An item of free text metadata can contain any string (up to a length of XXXXX bytes), although it may not contain tab character or linebreaks.

..... for items where the metadata can vary across every single text (e.g. text title, url)

```
* Freetext metadata fields are allowed to contain certain special forms which indicate
* external resources of one kind or another. These are detected by examining the value's
* "prefix", which is the part of the value before the first colon.
```

```
list($prefix, $url) = explode(':', $value, 2);
```

```
switch($prefix)
{
```

```
case 'http':
```

```
case 'https':
```

```
case 'ftp':
```

```
/* pipe is used as a delimiter between URL and linktext to show. */
```

```
case 'youtube':
```

```
/* if it's a YouTube URL of one of two kinds, extract the ID; otherwise, it should be a code a
```

```

if (false !== strpos($url, 'youtube.com'))

else
$ytid = $url;
$show = '<iframe width="640" height="480" src="http://www.youtube.com/embed/' . $ytid . '" fra

case 'video':
/* we do not specify height and width: we let the video itself determine that. */
$show = '<video src="' . $url . '" controls preload="metadata"><a target="_blank" href="' . $u
break;

case 'audio':
$show = '<audio src="' . $url . '" controls><a target="_blank" href="' . $url . '">[Click here
break;

case 'image':
/* Dynamic popup layer: see textmeta.js */
$show = '<a class="menuItem" href="" onClick="textmeta_add_iframe(&quot;' . $url . '&quot;); r
break;

default;
/* unrecognised prefix: treat as just normal value-content */
$show = escape_html($value);
break;
}

```

7.5.2 Classification

《》

TODO

Classifications – for items where each text falls into one of a limited number of categories (e.g. written vs spoken) In this case the field values must be handles (Unix words as described above), NOT multi-word explanations

7.5.3 Unique ID

《》 Unique IDs are handles, so all the rules of handles apply: see *《XREF》*

TODO
TODO

7.5.4 ID-link

《*This is where the concept of idlink is explained. Use spoken corpus as the explanantion*》

TODO

Notion of INDIRECTION

[datatype = idlink] implies lots of other fields in the associated metadata table, which in turn all have datatypes

But multiple indirection is not allowed.

AN EMAIL I WROTE ON THE LIST TO EXPLAIN THIS *《turn this into coherent paragraphs》*

TODO

The idea is that this offers a layer of **indirection** for the XML.

The paradigmatic case of an ID link is speaker metadata.

In a spoken corpus you have lots of utterances (<u>) and you very often want to do operations within certain utterances and not others based on features of the speakers.

EG you might want to search only within speech by males, or by people in a particular age group.

You COULD add XML attributes for each of these things ie

```
<u speaker_age="12" speaker_sex="male">Hello!</u>
```

But this is not a terribly good design, because age/sex are not features of UTTERANCES, they are features of SPEAKERS.

This speaker will always be male and 12 in this corpus, so why is it necessary to repeat this on every utterance?

The answer is, it is not.

The IDLINK datatype allows us to model this kind of indirection.

Instead of marking speaker features on utterances, we can have a separate table for speakers...

```
ID    age    sex
=====
A001  12     m
A002  65     f
```

.... which is then referred to by the IDLINK attribute.

```
<u who="A001">Hello!</u>
```

So, instead of the data chain going Utterance -> sex , there is another layer of indirection: Utterance -> speaker -> sex.

It's called an IDLINK because once we declare the datatype of s-attribute u_who to be IDLINK, we promise CQPweb that an IDLINK metadata table, with all the right IDs listed, will be available. That is, that the content of the IDLINK (u_who) always LINKS to an ID that exists elsewhere (in the Speaker metadata table, which is therefore an IDLINK metadata table).

All this is then opaque to the general user, who can specify " find instance of word X where speaker is male " in a restricted query, for instance - CQPweb will t

- use the IDLINK table to look up the IDs of the speakers where sex = m
- use the CWB index to find the list of regions in the corpus where u_who is equal to one or other of those IDs (IE utterances by one of those speakers)
- search within only those regions of the corpus for word X

Note this is SIMILAR to how text metadata works (if you search within genre "fiction", then CQPweb looks up the texts where the "genre" column contains "fiction", and searches only within those texts) but not the SAME.

The key difference is that text_id is a unique identifier, ie each text ID occurs in the corpus once and exactly once.

However, IDLINKS aren't unique. There can be many, many utterances where who="A001". A0001 is unique *in the Speaker metadata table*.

This is why we talk about "u who" as an IDLINK rather than an ID:
it is not an identifier, but something that links to an identifier.

=====

Currently, it is not possible to have IDLINKS as a datatype for text metadata.

I was uncertain about this decision, as there is a clear use case: where one author writes many texts within the corpus, it would make sense for the "author" column in the text metadata table to contain an IDLINK to a separate Author metadata table which would contain things like author sex, age, domicile etc.

I decided against this for two reasons.

First, this system for doing Restricted Queries based on things like utterances was already *very* difficult to make work. Making it possible for there to be ANOTHER layer of indirection might have driven me mad. Wibble.

Second, if you look at corpora in practice, people tend not to mind making the sex of an author, say, part of the text metadata rather than having the author-people as a separate data entity. This is the case for the written BNC for instance - in CQPweb's predecessor, BNCweb, "sex of author" is a "written restriction" (IE text metadata). So I just went along with this way of doing it.

now, after all that background, re Chao's question:

Since the assignments are texts, features of the students who wrote them could be included as text metadata columns.

And if one text metadata column in a unique id for the speaker, that makes it easier down the line to track the progress of individual students (you can say things like "create a subcorpus of texts where student=A001 in module=101, and compare a subcorpus of texts where student=A001 and module=201.
- and do many other comparisons, if you have the right metadata for the texts.)

This is what Jiayue recommended, but hopefully the reasons why what you want here is text metadata rather than an idlink now make sense.

-----Original Message-----

From: cwb-bounces@sslmit.unibo.it [mailto:cwb-bounces@sslmit.unibo.it] On Behalf Of Jiayue Wan
Subject: Re: [CWB] Example of metadata file?

Hi

My solution would be to use a metadata file (ascii text file, tab separated values) like this:

```
A201 A 201
B201 B 201
C201 C 201
D201 D 201
```

The first column are the text_id's; the other columns are used to make "text categorisation". In this way the four texts are linked clearly to the same student.

Jiayue

On 10/01/17 03:33, Chao Sun wrote:

```
> Hello all,
>
> First time poster and also want to try if my subscription works. I do
> not have any linguistic background, so please be gentle if I am asking
> silly questions.
>
> I am wondering if any one could provide a comprehensive metadata file
> example with some brief explanation on how CQPWeb can utilise the
> information? I am particularly interest in the LinkID part and assuming
> this could be used for threading different articles in a corpus?
>
> In my example, handed in assignments from various semesters are compiled
> as a corpus, each assignment is a text file with a unique text_id. Is it
> possible to give each assignment various linkIDs to show how the student
> progress through all semesters? For instance, student 201 has four
> assignments in semester A, B, C, D. If I associate four columns of
> linkID (A201, B201, C201, D201) on all his four submissions, will I be
> able to analyse the progress/change in words for this individual student
> in CQPWeb?
>
> Not sure if this is how the linkID and other metadata are designed for,
> besides classification and description. Please correct me if this makes
> non-sense.
>
> Regards,
> Chao
>
```

7.5.5 Date

《(currently under development)》

TODO

7.6 Metadata templates

Just as you can create templates for particular structures of annotation and XML, and reuse those templates across corpora, you can do the same for metadata structut (see 《XREF prev chapter on templates》),

TODO

7.7 Matadata file format

A metadata input file should be formatted as follows:

- It should be a plain text file.
- The encoding should be UTF-8 (or ASCII).
- The line breaks should match the format of the underlying operating system: CR+LF (U+000d, U+000a) on Windows, LF (U+000a) on UNIX-like systems.
- The file should *not* begin with a Unicode “byte-order mark”.
- The data should be arranged in a tabular format, where
 - *columns* are delimited by the tab character;
 - *rows* are demilited by line breaks;
 - and therefore, individual values cannot contain either tabs or libne breaks.
- The first column must contain the unique ID codes (text IDs for text metadata; the ID-link IDs for ID-linked metadata).
- No ID code can appear more than once; every text in the corpus must be listed.
- However, it is not necessary for the text IDs to be in any particular order.
- This first column is implicit in the declaration of the metadata; you **do not** include it explicitly in the definiton of the metadata (whether declared ad hoc or with a template).
- Every subsequent column represents a metadata field, as described when the metadata is inserted into CQPweb.
- The file itself must contain no header row.
- If a column contains values for a field of datatype *classification*, *unique ID*, or *ID-link*, then all its values must be valid *handles*. See section 6.3.
- If a column contains values for a field of datatype *date*, its contents must be formatted using CQPweb’s date formalism. 《XREF!》
- Columns of datatype *free text* are not restricted in what they can contain, except for the general rule (noted above) that they may not contain tabs or linebreaks.

TODO

- Empty values should be represented as zero-length strings - that is, if there is nothing in a given field on a given row, then there should simply be nothing between one tab and the next.

This format is often called *TSV* for “tab-separated value” (by contrast to *CSV*, “comma-separated values”, a text-based file format associated with Microsoft Excel).

In MySQL terms, it is equal to the format required by the `LOAD DATA INFILE` command with escape-sequence interpretation turned off within field values (using the command `FIELDS ESCAPED BY ''`). In fact, this is precisely the command that CQPweb uses to load the metadata! See <http://dev.mysql.com/doc/refman/5.7/en/load-data.html> for more detail.

7.8 Installing metadata

Text metadata needs to be installed for every corpus. Without it, queries don't work. So this should be the first thing you do after the CWB indexing process has run to completion (see *《XREF》*). TODO

There are four possibilities for creating text metadata:

- *From file*. The most common approach: install text metadata from a text file on the server. The structure of the file can be described either by specifying it *ad hoc* or by referring to a pre-determined metadata template 7.6
- *From XML*. This approach translates the values of one or more XML attributes, encoded in the CWB index as s-attributes, into fields of text metadata.
- *Minimalist*. This is the approach to use if you do not have or do not want to use any text metadata. It sets up the text metadata structure in the database, allowing queries to work, but leaves that structure empty.

All

Minimalist

From file Using template

From file ad hoc

Fromn XML

After you install text metadata, the *Manage metadata* page for the corpus will change. Instead of showing options for installing the text metadata table, it will show a single control allowing you to *Reset the metadata table* (i.e. delete it, allowing reinstallation). It should be noted that this involves total loss of the existing data!

《Manage text caTEGOIRIES》 TODO

Under Admin tools – manage text categories

Briefly,

- this page will give you a form for each text-metadata field that was installed with the datatype “classification” - each of these forms will list all the category handles that exist within the given classification - by default, category handles are mapped to a “description” that is the same as the category handle itself

XREF to notion of “descroptin” in prev chap/.

- BUT you can use the forms here to change the category descriptions to something more user-friendly
- since category handles are limited to short codes with no spacing or punctuation this is often useful

- if you do this, then the "descriptions" will show up in a whole lot of different places in the user interface instead of the category handles, including: – restricted query form – concordance header for a restricted query – distribution display – text metadata page.

《installing idlink metadata》

TODO

8 Parallel corpus data

8.1 Introduction

CQPweb supports parallel corpora as of version 3.2.22. This chapter explains how to link parallel corpora together, and how the display works once the alignments are set up.

8.2 Setting up parallel corpora

Parallel corpora are linked by the existence of *alignment attributes* (*a-attributes* for short) in CWB. For a full understanding of how a-attributes work, you are referred to the CWB documentation, especially the **Corpus Encoding Tutorial**.

The key feature to note about a-attributes is that they presuppose the existence of a pair of corpora, a *source* (to which the a-attribute belongs) and a *target* (at which the a-attribute points).

This means that, unlike other corpus data attributes, a-attributes cannot be created when a corpus is indexed into CWB; they must be added afterwards.

CQPweb builds on this procedure. Parallel corpus data is managed as follows:

- First, install the two corpora separately in CQPweb.
- Second, use command-line CWB tools to generate the a-attribute(s).
- Finally, return to CQPweb to register the alignment.

A note. If you install a corpus that has already been indexed via command-line CWB, it's possible that the corpus will already have one or more a-attributes. If so, these a-attributes **will be ignored** by CQPweb when it imports that corpus's registry data. This is because there is no guarantee at the point of installation that the *target* corpora for those a-attributes have also been imported into CQPweb - so it would be dangerous to register the a-attributes. In short, it is *still* necessary in such a case to register the alignment within CQPweb as a separate action.

8.3 Naming alignment attributes

A-attributes always have the name of the target corpus. So if you have two corpora `corpus_a` and `corpus_b`, which represent the same texts in Language A and Language B respectively, then an a-attribute from source `corpus_a` to target `corpus_b` will have the attribute-name `corpus_b` but will *belong to* `corpus_a`.

Since a-attribute handles are actually corpus handles, they necessarily have names that follow the corpus handle rules. Note in particular that they cannot contain hyphens. The CWB encoding tutorial recommends the use of hyphenated ISO language codes when using parallel corpora, e.g. `somecorpus-en` and `somecorpus-fr` for the English and French parts of a parallel corpus respectively. For use in CQPweb, however, this practice should not be followed. However, an underscore *can* be used instead if you wish (i.e. `somecorpus_en` and `somecorpus_fr`).

It's especially worth noting that the sample *Europarl* parallel corpus made available for download on the CWB website (<http://cwb.sf.net>) follows this convention. Both the Europarl corpora, and the a-attributes they contain, need to be renamed to use underscores instead of hyphens before they can be imported into CQPweb.

8.4 Creating alignment attributes

There exist multiple methods for indexing alignment attributes. They are discussed in the **CWB Corpus Encoding Tutorial**. All basically revolve around identifying (sets of) ranges of corpus positions which are considered to be “aligned”, i.e. translation-equivalent, across a pair of parallel corpora. An a-attribute contains the data for such a mapping.

CQPweb does not currently have a web-based interface to the alignment generation and importation tools. You must use the command-line tools directly.

The steps are as follows:

- Either create or appropriately reformat the alignment data ready to be imported.
- Create an a-attribute using the appropriate encoding utility.
- Modify the registry file for the source corpus to add a declaration of the new a-attribute (some tools for a-attribute creation do this for you).

In CWB, alignments of a pair of parallel corpora can be, but do not have to be, bidirectional. That is, creating an a-attribute for `corpus_b` within `corpus_a` only creates an A-to-B link; no B-to-A link will exist unless you *also* create an a-attribute for `corpus_a` within `corpus_b`. CQPweb replicates this feature of the underlying system: it supports both scenarios, i.e. things will work just fine with either one-way or two-way links.

Also as in the underlying CWB, it's entirely possible for a CQPweb corpus to have more than one a-attribute, where each one has as its target a different parallel dataset.

Finally it should be noted that the character encodings of any pair of corpora need to be identical or at least compatible (e.g. ASCII/Latin-1). CQPweb normally abstracts away differences in the character encoding of the CWB index - regardless of what encoding the index uses, CQPweb works in UTF-8. However, parallel corpora represent the one case where this doesn't work. If `corpus_a` is in Latin-1 and `corpus_b` is in Latin-2, for instance, then the text of the latter will be treated in all corpus display functions *as if it were Latin-1* - which will be wrong. The preferred solution to problems arising here is to use UTF-8 for all sets of corpora which must be linked together as parallel.

8.5 Registering alignment attributes with CQPweb

Once you have created the a-attributes, you must then register them with CQPweb.

This is done using a one-button control that can be found within the Corpus Admin Tools, on the **Manage parallel alignment** page. The button, labelled *Click here to scan the registry for newly-added alignments*, does exactly what it sounds like it does!

To be more specific, when this function runs, the following happens:

- CQPweb loads the corpus's registry file.
- It searches the registry of a-attribute declarations.
- It subjects any that it finds to two checks.
 - First, does a corpus by the name of that attribute exist in CQPweb?
 - Second, is the alignment already registered within CQPweb?
- If the answers to these two checks are *Yes* and *No* respectively, the alignment is registered.

Once an alignment is registered, it will appear in a table headed *Existing alignments* on this page.

Unlike other types of attribute, it is not possible to add a long description to be used on-screen for an a-attribute. Instead, the “description” for an a-attribute handle is always the same as the target corpus’ on-screen title (as specified when the target corpus was indexed, and managed on its **Corpus settings** page; see 3.3.1).

8.6 How alignment attributes can be used

Once one or more a-attributes has been registered in CQPweb, it has the following effects on the behaviour of the system, as seen by all users:

- Various additional controls appear, allowing the user to turn on or off the display of parallel data in query results, or to switch from one aligned corpus to another.
 - An extra dropdown control listing available parallel corpora appears on the query form.
 - A similar control is present in the concordance display.
 - A similar control is present in the extended-context display.
- When the display of parallel corpus data is activated, it is shown in a separate table row below the main result in either concordance or extended-context display.
- Options are also available to download one or more parallel corpus regions in the *Download query* tool. Parallel corpus data is printed as extra columns in the downloaded text file.

XML visualisations defined for the *source* corpus (see section 9.4) will be applied to the parallel text *if* s-attributes of the same name exist in the target corpus. XML visualisations defined for the *target* corpus are always ignored.

Similarly, when the primary annotation for the *source* corpus is visible, an annotation will only appear for the parallel text if the target corpus has a p-attribute with the same handle - in which case, moreover, the annotation shown is always the one whose name matches that of the source corpus’ primary annotation, *not* the primary annotation of the target corpus.

One final note of warning: the chunk of parallel data displayed is *the entirety of the zone of the target corpus* that is aligned to the zone in the source corpus in which the query match is found. This may be a longer or shorter stretch of text than the co-text shown around the match.

Consider the following two - quite likely! - contingencies.

In concordance view: the width of concordance is very often set in words, whereas alignment is typically built on top of s-attributes. The whole of an s-attribute unit may be shown for the parallel text, even if a span of (say) +/- 10 tokens is showing for the main corpus. This is not avoidable, as there is no way to reduce the size of the unit of parallel co-text; any reduction might easily exclude the part of the parallel text that actually relates to the query match. You are recommended to set the concordance display width to be equal to 1 of whatever unit the alignment is built on (typically sentences), rather than the default word-based width, in order to mitigate this issue.

In extended-context view: only one unit of parallel text can be shown, and unlike the source-corpus context, it can’t easily be widened (because, as noted above, only that unit of the target-corpus directly parallel to the source-corpus unit in which the match appears is returned by CQP when a-attribute display is switched on).

8.7 Parallel corpora and user privileges

To access any data in a parallel corpus (in concordance or extended context, or in a download), a user must have *at least* restricted-level access to that corpus.

That is, assuming Corpus A has an a-attribute linking it to Corpus B, that a-attribute will be visible to the user in the interface for Corpus A only if they have some privilege that allows them access to Corpus B.

Otherwise, it will appear to the user as if the a-attribute for Corpus B does not exist.

Typically, one would set up the different component corpora of a parallel corpus so that *all the same* users and user groups are granted the use of those component corpora.

A user *does not* need to have access to a parallel corpus to use its a-attribute within the body of a CQP-syntax query (e.g. as a limiting factor on the query; see the CQP Query Tutorial); this is allowed even without an access privilege, because it does not involve access to the text of the parallel corpus. However, this possibility is not openly flagged in the interface; the list of attributes that appears on the query form only includes a-attributes that can be viewed, omitting any a-attributes that can be used but not viewed.

9 Controlling query visualisation

In CQPweb, *visualisation* is a covering term for anything to do with how corpus data is rendered in the web interface, but with particular reference to the display of text around a query hit in the concordance display and in the extended context display. This chapter describes how the system administrator can control different aspects of the visualisation process.

9.1 How the primary annotation affects visualisation

《How the primary annotation affects visualisation》

TODO

9.2 Setting up an “alternate” view for context display

《Setting up an “alternate” view for context display》

TODO

9.3 Using position labels

《What they are; quick note on the control》

TODO

《don't forget to explain about the *span-class* they are wrapped in and the possibility of doing things with extra code files, *q.v.*》

TODO

9.4 XML visualisations

9.4.1 Introduction

By default, none of the corpus XML (s-attributes) are displayed in either concordance or extended context view. Nor are they included in downloaded queries. This section explains the use of the XML visualisation system, which allow you (a) to display corpus XML in query results, and (b) transform the raw data of the s-attributes into customised HTML for easier interpretation by users.

An “XML visualisation” is a specified mapping from an s-attribute (XML element or combined element-attribute) to some HTML. These are created through the CQPweb web interface; see 9.4.2. The mapping takes the form of an HTML template written with a small “whitelisted” subset of HTML, with all complex features blocked to avoid the risk of Cross Site Scripting (XSS) vulnerability: see 9.4.5.

You can further manipulate the HTML generated by your visualisation with *extra code files* (see 9.4.6), additional JavaScript or CSS files which you install server-side and which can apply almost unlimited further styling to the visualised XML.

It should be noted that all XML visualisation operates at the corpus level. That is, the visualisations you set up for one corpus do not affect any other corpus on the system. Of course, it is easy enough to copy the HTML code of a visualisation from one corpus' interface to another's!

9.4.2 Creating and managing XML visualisations

The interface for creating and managing XML visualisations can be found within the Corpus Admin Tools, on the **Manage visualisations** page. This page also contains the controls for field-data display mode (see section 9.5) and for position labels (see section 9.3).

The various control forms for XML visualisation are towards the end of the page. In order, you will see:

- The interface for using extra code files: see [9.4.6](#)
- The fallback control: see [9.4.7](#)
- A list of existing XML visualisations
- A form to create a new visualisation command.

To create a new visualisation, you must select the XML element/attribute (s-attribute) that you want to render, and specify whether you are creating a visualisation to appear at the start of stretches of text of that type (start tag) or at the end of such stretches of text (end tag).

Next, you need to enter the actual HTML code. See [9.4.5](#), [9.4.4](#) for the content you are allowed to use here. “HTML code” in this context *includes* any plain text content you might want to use (e.g. a single punctuation code to visualise the boundary indicated by the XML) - and, as noted below, for visualisations intended to be used in downloaded queries, it is usually *preferable* just to use plain text.

Next, you must specify whether the visualisation is to be enabled (a) in the concordance display, (b) in the extended context display, (c) in downloaded queries. It's possible to use the same visualisations in two or three of the possible places. It's also possible to use separate visualisations for the same bit of XML in concordance and/or context and/or downloaded queries.

Finally, you can make the visualisation for a start tag conditional on the value of the XML annotation; this is explained in [9.4.3](#).

Once you have created a visualisation command it appears on the list of existing visualisations. The definition of the target of a visualisation (XML element, start vs. end tag, condition for the value) cannot be changed once it has been created: if you need to change these, delete the visualisation and create a new one. However, you *can* change the HTML code of an existing visualisation, plus of course activate/deactivate its use in the two relevant displays and/or query download.

9.4.3 Conditional XML visualisations

XML visualisations for start tags can be made *conditional* on the value of their annotation. This means that the visualisation only applies when the annotation meets some criterion.

Currently, the only kind of condition that can be applied is a *regular expression match*. That is, when you create the condition, you specify a regular expression that the value is compared to. If it matches, the visualisation is used. If not, not (although if there is another visualisation active with a different condition or no condition, that one might well apply).

The regex flavour used in the conditions is *Perl-Compatible Regular Expressions (PCRE)* (the same regex flavour used by CQP's query mechanisms, as of v3.2 of CWB). However, unlike regexes in CQP, the regexes for conditional visualisations are not anchored: it is not necessary for the regex to match the *whole* value, only that some part of the value must match the regex. If your regex pattern includes any forward slashes, you must escape them with a backslash. This is in addition to all the usual escaping rules for PCRE regexes.

If there is more than one conditional visualisation, then when an XML boundary is rendered, their conditions are checked in order (the conditional visualisation with the longest regex first, and then by descending length). The one that takes effect is the first where the condition is fulfilled (i.e. where the value of the XML attribute contains a match to the regex anywhere within it).

Example use cases for conditional visualisations might include:

- You have an XML attribute with a small number of values - for example, a pragmatic categorisation of annotated chunks of discourse (question/request/command/statement etc.) You want to visualise the beginning and end points as square brackets “[...]” in the interface. If you create a series of conditional visualisations, each with one of the possible values as the regex to be matched, then you could assign different appearances to the opening of each different pragmatic function - e.g. colour them differently depending on the function type.
- You have an XML attribute with two possible values, where one is assumed (the default) and the other applies in only a limited number of cases. By creating a conditional visualisation for the latter but not the former, you can have it appear in the interface only when it has its less common, unexpected value and let it remain invisible when it has the more common, default value.

Be careful with conditional visualisations. It is necessary to consider all possible contingencies, because it will always cause problems to leave any bit of XML uncovered.

For instance, if you have a single visualisation that applies to `s_type` when its value matches “question”, that leaves all other values of `s_type` with no visualisation. They will instead be rendered as nothing. You can specify this explicitly by creating an empty default, that is, something like this:

- Visualise start of `s_type` as [some block of HTML], with condition that the value must match “question”;
- Visualise start of `s_type` as [leave empty], with no condition.

The second condition will apply to all instances of `s_type` with other values, causing them simply not to show up in the display.

Without such an explicit visualisation, any XML boundaries that are not matched by any of the relevant conditions default to being rendered as nothing anyway.

9.4.4 The embedded variable

As well as making the use of a visualisation conditional on the value of the instance of the s-attribute, it is also possible to actually include the value in the HTML rendering.

This is done using the *embedded variable*. This is, quite simply, a sequence of four dollar signs (\$\$\$\$) appearing anywhere in the HTML code of an XML visualisation.

Note that the embedded variable (a) only works for start tags, not end tags; (b) cannot be used for s-attributes that don't have values! When the s-attributes have, as is usual, been created from structured XML, the s-attribute “head” of the “family” lacks values, but the “children” representing the XML attribute-value pairs have values.

So, for instance, if you have `<p>` elements in your corpus with attributes `num` and `type`, then the s-attribute `p` will have no values (so the embedded variable won't work), but the s-attributes `p_num` and `p_type` will have values that can be inserted into the HTML using the embedded variable.

9.4.5 HTML allowed in XML visualisation code

Only certain HTML elements are allowed in XML visualisation code. When you add a new visualisation, all the HTML in your code will be compared to the “whitelist”, and any HTML that is not on the whitelist will be escaped.

So, for instance, `<script>` is not allowed in visualisation code. If you try to install a visualisation containing `<script>`, it will be escaped to `<script>` and appear literally in the browser as `<script>`.

The HTML whitelist is as follows. For convenience, it is divided here into (a) simple codes (just tags with no attributes) and (b) complex codes which require an attribute.

You can use the following simple HTML formatting codes (just tags with no attributes):

- `...` - render text as bold
- `<i>...</i>` - render text as italic
- `<u>...</u>` - render text as underline
- `<s>...</s>` - render text as strikethrough
- `_{...}` - render text as subscript
- `^{...}` - render text as superscript
- `<code>...</code>` - render text as code (usually means a monospace font)
- `
` - add a line break (*hint*, for the appearance of a paragraph break, use `

`)

Of course, nothing forces you to close HTML tags that you open in an XML visualisation; however, if you don't, you may find that your formatting interacts peculiarly with CQPweb's own rendering.

The available complex codes are as follows:

- `` - embed an image; the URL can be relative or absolute
- `...` - apply one or more CSS classes to a span of text (to link a class to one of the other tags, wrap it in a span)
- `...` - create a link (http/https only) (all links will open automatically in a new browser window or tab)
- `<bdo dir="ltr/rtl">...</bdo>` - mark text as left-to-right / right to left: CQPweb uses these tags with right-to-left alphabet corpora, and you may need to use them in your visualisation code to stop odd interactions between your HTML and the direction of the text it appears in

For both simple and complex tags you need to use the lowercase form shown above, i.e. `
` not `
`; full HTML allows either, but CQPweb only whitelists the former. Similarly, you must use double quotes for attribute values, even though full HTML allows single quotes.

As well as HTML tags from the whitelists above, you can also use HTML entities. You don't need to use entities for accented characters or non-basic punctuation, however, although these are two typical uses for entities, because the UTF-8 encoding is used throughout CQPweb, so you can just use the characters directly. One useful entity is ` ` (no-break space), which functions as a unit of "empty" text (note that actual whitespace is ignored).

It's worth pointing out that the following frequently-used HTML codes are *not* usable in visualisations, because they would interact badly with the HTML/CSS that CQPweb itself uses.

- `<p>` and `<hN>`

- `<div>`
- `<pre>`
- `<table>` and other table-structure tags

You can, however, still make use of these and other features of HTML by means of *extra code files*, explained in section 9.4.6.

When you create visualisations for use in downloaded queries, remember that any HTML will normally *not* be rendered: users will see the actual code in the plain text file they download. For that reason, it's normally better just to use a plain text visualisation here (possibly with the embedded variable).

9.4.6 Extra code files

The combination of HTML, CSS and JavaScript in a dynamic web page allows almost limitless variation in how information is presented. However, for security reasons, it is not possible to allow all the richness of HTML to be usable in the XML visualisations that are created via the interface: without the restriction to a limited subset of HTML, an attacker who gained access to a superuser's account would be able to insert arbitrary JavaScript code into the browser of any user who subsequently accessed that corpus.

However, the simple visualisation code permitted by the whitelisted-HTML may well be *too* restrictive for users who want complex rendering of data in their corpus XML. CQPweb has an additional mechanism to allow users to perform more complex formatting: the insertion of **extra code files**.

An extra code file is a file containing either CSS or JavaScript code. You can specify one or many such files for each corpus (with separate lists for concordance view and context view to allow different renderings in each). These files will be linked in the headers when a concordance/context page is generated. The overall appearance of the concordance/context will thus be governed by the combination of CQPweb's built-in styling and the extra CSS/JavaScript you have added.

An extra CSS file allows you to write arbitrary stylesheet code to change the appearance of text within a visualisation. For instance, if you have used a `` element to assign a class to some piece of text, you can use an extra CSS file to apply any style(s) you want to spans with that class.

An extra JavaScript file is even more powerful. You can manipulate the HTML document tree in more-or-less unlimited ways. In addition, CQPweb uses the **jQuery** library, and your extra code can use this library too. jQuery uses CSS selectors to pick HTML elements to operate on, so again, the assignment of classes to particular bits of text in the visualisation code allows you to later write a JavaScript/jQuery function to pick out those specific parts of the document and modify them as you wish. An example of this is provided below.

You can, of course, use extra code files to modify aspects of the concordance and context views other than just the XML visualisations. But it is that combination which offers the most powerful options.

There are three steps to add an extra code file.

- Write the code.
- Insert the file with the code into CQPweb (must be done on the server, cannot be done via the web interface).
 - CSS files should be placed in the `css` subdirectory of the main directory.
 - Javascript files should be placed in the `jsc` subdirectory of the main directory.

- Note that extra code files must have the canonical file-extensions: `.css` and `.js` respectively.
- Finally, activate the code file for the corpus you want to apply it to.

Extra code files are activated using the **Manage visualisations** display, under the heading *Extra code files for visualisation*. Select a file from the dropdown and press “Add this file” to activate it; use the [x] buttons to deactivate code files from a corpus. Deactivating does not delete the file, and only affects the particular corpus in question - any other corpora where the same file is activated will be unchanged.

There are two separate forms for concordance view and context view, allowing you to have very different layouts in concordance and context views. Here are two examples.

In our first example, let's imagine that the corpus contains `<u>` elements for utterances, with a **speaker** attribute, giving the s-attribute `u_speaker`. Let's assume we have created the following visualisation for `u_speaker`:

- `$$$$: `

We can exploit this with an extra CSS file to make speaker labels stand out more distinctively in the display. The CSS file could contain the following:

```
/* ----- THIS IS THE CSS FILE */
span.swanky-speaker-label
{
    color: pink;
    font-weight: bold;
    font-family: "Comic Sans", sans-serif;
    font-size: 16pt;
}
/* ----- END OF THE CSS FILE */
```

When this file is specified as an extra CSS it will apply this highly tasteful additional styling to all speaker labels.

If we want to do something more advanced (for example: make all speaker labels clickable, causing an information box to appear when clicked) this can be done in an extra JavaScript file using either the native JavaScript method of interacting with the DOM (document object model) or the facilities of the jQuery library. The following code exemplifies the latter:

```
/* ----- THIS IS THE JAVASCRIPT FILE */
function show_my_info(speaker)
{
    /* The following is a placeholder for the more complex
       behaviour you would want in a real situation. */
    alert ("The speaker is " + speaker + " !");
}

$(document).ready (function() {
    /* this passes the whole content of the span to the show_my_info function
       * (including any extra matter added around the value, like a <:> (See above).
       * In reality, you would usually want to make it easier for the jQuery code
```

```

    * to access just the value, to do interesting things with it.
    */
    $("span.swanky-speaker-label").click(function () {
        show_my_info( $(this).html() ;
        return true;
    } );
} );
/* ----- END OF THE JAVASCRIPT FILE */

```

One thing that can be done using extra JavaScript files is to create renderings that involve the content of more than one s-attribute. For instance, imagine our original input-format utterance tags look like this:

- `<u speaker="ID_Code" type="QUESTION">`

If we want to make use of *both* the `code` and `type` values to build what appears in the concordance at an utterance boundary, the easiest way is to use a basic visualisation, and then join the two together using JavaScript.

First, we would add these visualisations:

- For `u_speaker` : `$$$$`
- For `u_type` : `$$$$`

The key thing to note here is that the visualisations of attributes within the family of a single XML element will *always* appear directly together. So the above will always generate something like this:

- `ID_CodeQUESTION`

... which is easily manipulable via the HTML DOM.

```

/* ----- THIS IS THE JAVASCRIPT FILE */
$(document).ready (function() {
    /* on document ready, run this function on each of the outer spans */
    $("span.all-info").each(function () {
        var outer  = $(this);
        /* extract the string contents of the inner spans */
        var speaker = outer.children().first().html();
        var type    = outer.children().last().html();
        /* having done so, we can now simply set the inner-html of
        * the outer span to the string we want */
        outer.html(
            "[[This is a <b>"
            + type
            + "</b> spoken by <b>"
            + speaker
            + "</b>]]"
            /* in reality you'd want more elaborate HTML! */
        );
    } );
} );
/* ----- END OF THE JAVASCRIPT FILE */

```

Finally, just to note the obvious: extra code files can't be used with downloaded queries.

9.4.7 Fallback visualisation methods

Many corpora are set up with no XML visualisations (indeed, XML visualisations were not implemented until CQPweb had already been around for many years). In this case, the text from the corpus in the extended context display will all appear as a long, unbroken block.

To make extended context look a little nicer without the need to create a visualisation, CQPweb possesses a *fallback* method. When this is switched on, a double-line-break is rendered in context view *after every token made up only of non-medial punctuation*. This roughly simulates having each sentence in a separate paragraph. This will not work with every language or every type of corpus data. But in most cases it results in a reasonably-OK way to break up the wall of text.

If you have added a visualisation to insert breaks in context view, you might well want to turn off this fallback method. Of course, you might also just want to turn it off anyway! A control can be found to do this on the **Manage visualisations** view, under the heading *Visualisation fallback procedures*.

When a corpus is indexed, it has no visualisations initially. So the fallback is always switched on by default for newly-added corpora.

9.5 Field data presentation mode

《Field data presentation mode》

TODO

9.6 Field data mode as a workaround for parallel corpora

《Field data mode as a workaround for parallel corpora》

TODO

10 User accounts and privileges

10.1 Basic concepts

Access to CQPweb is based on *user accounts*. The user account has two functions. First, it determines what resources (corpora and analysis options) a given user has access to. Second, it provides a basis for (limited) personalised functions. For instance, every user account has its own query history, its own set of saved queries, its own set of subcorpora, its own concordance display preferences, and so on.

If you are running an online CQPweb system, then normally you will want to use CQPweb's functions for creating user accounts via a self-registration process. That is, someone who wants to use your CQPweb server should first visit the homepage, where they will find a link to a form for account creation. This is very similar to the process for signing up for an account on pretty much any website, so should not be challenging. There is also a tool in the Admin Control Panel for you to either create accounts for users yourself, or to send an invitation to a specified email address.

User accounts must be linked to email addresses. This is because knowing the user's email address is the only way to make it possible for them to reset a forgotten password. Again, this is pretty much standard procedure on the web.

Each user account belongs to one or more *user groups*. A user group is exactly what it sounds like: an arbitrary set of user accounts. There are two builtin groups: **superusers** and **everybody** - both of which do exactly what it sounds like they do!

User accounts can be added to groups other than **everybody** in two ways: either automatically at the time of account creation (based on pattern matching against their email address), or manually by you using the Admin Control Panel.

The final core concept is the *privilege*. A privilege simply defines some permission that can be granted to a user. This might be a level of access to some corpus or set of corpora, or permission to initiate a database operation of a certain size. Privileges may then be assigned to users individually, or to user groups. A link between a privilege and either a group or a user is called a *grant*. Each user then has all the privileges granted to them or to any of the groups to which they belong. This set of privileges determines what CQPweb will and will not let that user account do.

It should be noted that any user account in the group **superusers** automatically has every possible corpus-access privilege; other privileges can be granted or not granted to that group, as usual.

The system is very flexible, which makes it rather complex. This chapter explains different aspects of the system, and how you can control it.

10.2 User accounts

By default, users can create their own accounts using the self-registration system. However, if you don't want to allow access to anyone you have not vetted individually, it is possible to switch this system off: see 2.3.6. Self registration runs as follows: first, the user specifies the username and password they want, and supplies an email address (and, optionally, other information about themselves, including their real name, location in the world, and organisational affiliation). They may also have to answer a CAPTCHA challenge, if you have switched that on. The system then sends an email to the address specified with a verification code. The user must click on the verification link: their account is then activated, allowing them to sign in with their username and password. This verification system (a) prevents the creation of accounts that the owner cannot retrieve in case of a forgotten username or password, since it guarantees that a correct email address has been supplied; (b) stops anyone creating an account for someone else.

- A known “gotcha” here: CQPweb’s automated emails are known to be blocked by the spam filters used by Yahoo’s free email service. There does not appear to be anyway around this. So if you have users on Yahoo mail, you will probably need either to validate their email addresses manually, or to instruct them to use a different free email service (Google Mail seems generally to work fine).

Whether or not self-registration is enabled, there are also tools available in the Control Panel for you to create accounts for people, as well as view (or delete) existing accounts. Go to **CP >Users and Privileges >Manage users** to access these tools.

At the top of this screen is a summary of the existing accounts on the system. *«Note that the tool for looking at a list of unverified accounts, linked at the bottom of the summary table, has not been written yet.»* **TODO**

Next is the user-account search form. This is how you view the status of an individual account. There are two search tools. The first is the quick username search. Simply start typing the username of the account you wish to view: a list of suggestions will be provided as you type, and you can simply click on one of the links as soon as you have narrowed the list down enough. (Pressing TAB from the quick search box will move you through the suggestions list; press ENTER to view the selected account.)

The second search tool, the *Full search*, checks for the term you enter in three different fields: username, real name, and email address. The results include accounts where your search term appears *in the middle* of one or more of those fields. This is especially useful if you are looking for an account whose username you are not certain of: searching for a fragment of a person’s name here is very likely to find their account. This tool has a more conventional search interface than the quick search - you must press the **Search** button to go to a table of results, each of which gives you a link to view the account details (see section 10.3).

Underneath the search tools are the account creation tools. The main account creation form is an abbreviated version of the user’s self-registration form. The difference is that only a username, password and email address are needed. The optional fields (real name, affiliation, location) are not set here; the assumption is that if the administrator creates an account and tells the user what their credentials are, they will log in later and add these details themselves.

This form generally has much less security than the corresponding self-registration form, based on the assumption that the administrator(s) know what they are doing:

- There is no CAPTCHA
- The password is not masked
- The password does not need to be typed twice

Similarly, when you create an account through this interface, validation by email is optional: you can cause the account to be automatically validated if you are certain that the email you’ve entered is valid; and you also have the option to leave it unvalidated without sending an email.

«It is intended to add an “invitation” tool here, to allow people to be sent a “please sign up” link to the account creation form (prepopulated to the extent possible). However this has not yet been implemented» **TODO**

One final note on account creation relates to security.

If you’re using a normal web server, CQPweb passwords are transmitted in plain text via HTTP. This is very insecure. Theoretically there is little danger in a user’s CQPweb password being stolen. User data stored on CQPweb is seldom highly-sensitive, though of course there may be exceptions.

However, if self-registration is enabled, there is a high level of danger - since users select their own passwords, and it is well known that **most users tend to reuse passwords on multiple sites**. This means that a user's CQPweb credentials being stolen potentially exposes their accounts on other systems - in the worst case, highly consequential systems like social media or online banking.

To protect users from the consequences of password reuse, you have two options:

- Disable self-registration. Generate all passwords for users yourself, to prevent reuse, using highly non-memorable passwords to deter the user from going on to reuse the password you have chosen.
- Use HTTPS instead of HTTP. This is the recommended course of action.

One option that is known *not* to work is instructing users not to re-use passwords. Such instructions are routinely ignored - if the user even notices them in the first place.

CQPweb stores passwords internally in an encrypted form, so that they are protected even if your server is hacked. However, this does nothing to protect users if their credentials are transmitted in plain text via HTTP.

10.3 Viewing user account details

When you view a user's account details, you will see most of the basic information recorded about the user in CQPweb's internal database: their real name, email, and affiliation (as supplied when they signed up), plus also the time of account creation and the time the user last accessed CQPweb.

- If the date of account creation is either "0000-00-00" or a date in 1970, it means that the account was created before CQPweb started tracking the ages of user accounts (in version 3.1.0).
- If the last-visit date is either "0000-00-00" or a date in 1970, it means that the user has never logged in.

manually validating an account

deleting an account

resetting the password

corpus stats

setting the database size

《*account expiry*》

《*password expiry*》

《*bulk log out, and the log in list in the user profile. Here or in the prev section??*》

TODO

TODO

TODO

10.4 User groups

A group is exactly what it sounds like: a group of user accounts. You can organise user accounts into groups manually based on whatever principle you see fit; alternatively, you can organise users into groups automatically using patterns in their email addresses.

The point of user groups is that it is usually much more convenient to grant privileges, such as the privilege to access a certain corpus, to a group of users all at once, than to grant such a privilege to every user, one by one - especially if your server is open to the internet and anyone can create an account on it.

As noted in section [10.1](#), there are two builtin "special" groups, whose membership is not controlled in the usual way. They are:

- **superusers**, which contains all and only those users declared as system administrators in the configuration file (see 2.2).
- **everybody**, to which all users belong automatically, and from which no one can be removed.

Groups are created, managed and deleted via the Admin Control Panel: go to **CP >Users and privileges >Manage groups**.

- To add a new group, use the form at the bottom of the screen; all you need to specify to create a group is its name, which should be letters/digits/underscore only.
- The list of existing groups are tabulated at the top of the screen, in alphabetical order of group handle.
- This table allows you to enter a longer description for the group, which would normally be some *aide-memoire* to the group's nature and purpose.
- It also allows you to enter an *Auto-add regex*, which is explained below.
- After entering or modifying the description or regex, press **Update** to save your changes.
- Also in this table you will find [x] buttons to delete groups.

To add or remove users individually from a group, go to **CP >Users and privileges >Manage group memberships**. Next to each group on this screen is a pair of controls - one for adding users who are *not* already members, and one for removing users who *are* currently members.

A user can be assigned to as many groups as you like.

Users can also be added automatically to groups at the point of account creation. This works as follows. It is possible to associate a regular expression (regex) with each group (this is the *Auto-add regex* discussed above). When a new user account is created, CQPweb checks each group regex against the new user's email address. If there is a match, then the user is added to the group in question.

The typical usage case for this function is the situation where you have created a group that is associated with a particular institution. In that case, you can set up the group regex to detect email addresses "@*anytown.edu*" that institute's domain.

The regex flavour used is PCRE (Perl Compatible Regular Expressions). You can use any PCRE feature in an auto-add regex. The auto-add regex can be very long - up to 64 KB. Needless to say, it is not recommended to use such a long regex! The form limits you to 1024 characters; longer regexes can be inserted by direct manipulation of the MySQL database if necessary.

Changes to the regex are not retroactive. So if a user A signs up on Monday with an email address ending *@anytown.edu*, and on Tuesday you modify the regex for some group G so that it matches addresses ending in *@anytown.edu*, the existing user A will not be added automatically to group G.

(Similarly, users are never removed from groups based on a regex changing. So if user A is a member of group G, and group G's regex changes so that it no longer matches user A's email, user A will *not* be removed from group G in consequence.)

You can, however, force a re-application of the regex to all existing non-members via the **Bulk Add** tool. This can be found at the bottom of the **Manage group memberships** screen.

The first of the two **Bulk Add** tools is *Apply group's stored pattern-match to existing users*. To use this, select a group and press the "...run group regex against existing users" button. All existing user accoutns that are not already in the group will be checked, and added to the group if their email

address matches. But note this is not exclusionary: if there are users *already in* the group whose email addresses do not, or no longer, match the regex, they will *not* be removed from the group by application of this tool.

The second **Bulk Add** tool works in exactly the same way, except that instead of using the stored auto-add regex of the group to which you wish to add members, it uses an *ad hoc* regex that you must supply via the form. This adds a further level of flexibility to the group system.

10.5 Privileges

Privileges represent things a user can do in CQPweb. There are multiple types of privilege, each controlling access to a different type of “thing” that can be done. The most important type is the corpus-access privilege, which determines which of the available corpora a user is allowed to use, and what they are allowed to do with that corpus.

This section first explains the different types of privilege, and then goes on to explain how you can create, monitor, and edit privileges in the CQPweb system.

10.5.1 Corpus access privileges

There are three levels of corpus access privilege. In order of ascending level of access, they are *normal*, *restricted*, and *full*. The latter two levels are defined relative to *normal*.

Normal access is the level of access that you would give to users who have all the necessary licences/IP rights to both use and reproduce the data. For example, if the corpus is one for which a licence must be paid, you would give normal access to any users known to have paid for a licence or to be affiliated to an organisation which has a licence. Alternatively, if the corpus is one that is openly accessible to all, you could grant normal access to *all* users on the system, via the **everybody** group.

When a user has normal-level access, they can do more or less everything. They can perform (and download) concordances without restriction, and they can use the extended-context function. These two functions open the possibility of users being able to access the complete underlying text of (one or more texts from) the corpus - by downloading, or copy-pasting, overlapping extended context chunks. By default, extended context can be expanded out to around 2,000 tokens, so it is in theory possible to extract the whole text of a corpus from the concordance of any evenly-distributed item with a frequency greater than about 0.5 per thousand tokens. This is why normal access should not normally be given to users who do not have the necessary IP rights to the corpus.

Restricted access implements certain limitations on what users can do that are intended to make it much harder, if not impossible, for them to be able to extract the complete underlying text from the CQPweb interface. The idea is that restricted access to a corpus can potentially be granted to users who don't have a licence - even if the corpus contains copyrighted material - because the restrictions stop them getting at any stretch of underlying text longer than concordance-length snippets, reproduction of which is more likely to fall under the “fair use” or “fair dealing” provision of copyright law.

- *Necessary disclaimer:* no one who works on CQPweb is a lawyer and nothing in this manual constitutes legal advice. The restricted access privilege is provided in the hope that you may find it useful as a *partial* tool for making sure you comply with whatever the copyright law is in your jurisdiction; but without any warranty whatsoever that it is either necessary or sufficient for that purpose. You use it at your own risk.

The restrictions implemented for that level of access are as follows:

- Users can't view extended context if they only have restricted access.
- Users are blocked from using the *Download tabulation* function.
- If a query has too many hits, it will automatically be thinned down to a random subset. (Too many = more than half the square root of the size of the corpus or subcorpus in tokens.)

To explain the last one: as noted above, if you do a concordance for something really common, and then download it, you could in theory reconstruct the whole corpus from the concordance lines - even without access to the extended context view. The “half the square root” limit tries to counteract this, while not making concordances in small corpora totally useless.

So, for instance, half the square root gets you 5,000 examples in a 100MW corpus, but 1,000 examples in a 1MW corpus. The limit is always rounded upwards to a whole number of thousands.

《*create table of some common sizes*》

TODO

The idea is that the random subset of examples will be spread out in the corpus so there will be gaps between them no matter how common the thing searched for is – so the underlying text can't be reconstructed.

Full access lets users access extra functions that involve total access to the data in one way or another. Currently, the only such function is the **Export corpus** tool, which allows the raw text of the corpus to be exported in plain-text format for analysis using other software.

A corpus access privilege is defined in terms of its *level* (restricted, normal or full, as per above) and its *scope* - where its scope is the corpus or set of corpora that it grants access to.

So, for instance, you might define a privilege granting *restricted* access to corpora A, B and C, and a second privilege allowing *normal* access just to corpus A. As explained further in section 10.6, CQPweb always uses the highest applicable privilege. Users granted both these privileges would therefore have normal access to A and restricted access to B and C.

Adding a new corpus to an existing scope is an easy way to apply the access pattern you have defined for one corpus to another. Continuing the example, if you create a new corpus D which is governed by the same licensing/IP conditions as corpus A, you could simply add corpus D to the privilege which allows *normal* access to A. All users granted that privilege would then have access to corpus D without you needing to separately configure any grants of new privileges for corpus D.

10.5.2 Frequency list privileges

《*Explain these*》

TODO

《*Explain necessity of giving “everyone” at least one privilege of this sort*》

TODO

《*explain that the “scope” here is a number:*》

TODO

10.5.3 Database privileges

《*Implement these; explain*》

TODO

10.5.4 File upload and disk space privileges

《*Implement these; explain*》

TODO

10.5.5 The CQP binary file privilege

Users who have this privilege are allowed to do two things that users without the privilege are not:

- Download binary-form CQP query files from saved queries
- Create a saved query by uploading a binary-form CQP query file

These two mirror-image functions are designed for advanced users who make extensive use of very large saved queries. Such users may exhaust their allowance of saved-query space. An easy way for them to clear out their disk space allowance is simply to export the binary files of the saved queries for external storage, and then reinsert those binary files if they want to use the queries in CQPweb later. This avoids extensive use of corpus-position dumps.

These functions should normally be restricted to advanced users simply because most users will not understand what CQP binary saved-query files are or how they work!

There is only one privilege of this type, and it is one of the default generated privileges. However, it is not initially granted to anyone, so the functions in question will only be available to superusers.

10.5.6 Corpus installation privileges

《Implement these; explain》

TODO

These privileges are required for users to be able to install their own corpus data on the server.

To allow users to install their own corpora, the first step is to enable the user-corpus system, by setting the `$user_corpora_enabled` configuration option to **true** (see 2.3.7).

This will only make the system *visible* to users. They will not be able to actually install corpora until their account has a privilege allowing them to.

Multiple privileges affect user-corpus installation:

- First, they must have a file-upload privilege (see above).
- Second, they must have a disk-space privilege giving them sufficient space in their upload area to store the input files for corpus installation.
- Third, they must have the privilege to use at least one of the available corpus-installer plugins.
- Fourth, they must have a disk-space privilege for user corpora.

The *third* privilege above governs how much data can be installed into a single corpus with a given installer plugin. The *fourth* governs the total amount of space taken up by all of the user's corpora (this gives the system administrator defence against some user or users taking up all the available disk space).

10.5.7 Creating and editing privileges

《Make sure that the right links in the admin-ui chapter point here》

TODO

《explain the interface》

TODO

《Explain how privileges are presented: unique ID number, then a short prose explanation》

TODO

《Explain the edit interface》

TODO

10.6 Grants: creating and managing grants of privileges

A *grant* is a link between some user, or group, and a privilege. The system works out what a user is and is not allowed to do based on the privileges granted to them.

A user has *both* any privileges that have been granted directly to them; *and* all the privileges that have been granted to *any one or more* of the groups of which they are a member.

Higher privileges override lower privileges. For instance, let us imagine there is a user A who is a member of group B, and we want to know what level of access this user has to a corpus C.

- If group B is granted restricted access to C, but user A is individually granted normal access to C, then the effect is that A has normal access to C: the grant to the user overrules the grant to the group.
- If group B is granted normal access to C, but user A is individually granted restricted access to C, then the effect is the same - A likewise has normal access to C: the grant to the group overrules the grant to the user.

The same is true with other types of privilege: when more than one applies, the one that has effect is the most expansive. This means it is always safe, and usually a good idea, to grant a lowest-common-denominator set of privileges to the “everybody” user group; users and groups with higher privileges will always benefit from those higher privileges, even though the users in question are (by definition) also members of “everybody”.

When a user is removed from a group, they lose the privileges associated with that group, *unless*, of course, they are linked to that privilege in another way - individually, or by virtue of membership in another group that is granted that privilege. To put it another way, the transfer of grants from a group to its members is dynamic and ongoing, not static and persistent.

Finally, if a privilege is edited, then all users and groups who were granted the privilege pre-edit are *still* granted that privilege post-edit. (This makes it easy to give a new corpus the same access pattern as some existing corpus, by putting it within the scope of the privileges that govern the existing corpus. It is then not necessary to add any more grants on the system.)

There are two menu options in the Admin Control Panel which allow you to create, monitor, and delete grants: **CP >Users and privileges >Manage user grants** and **CP >Users and privileges >Manage group grants**.

These menu options lead to very similar screens. At the top is a form for making new grants. To grant a privilege to a user or group, select the user/group from the first dropdown; select the privilege from the second dropdown; and press the *Grant privilege...* button.

«*in future it will be possible to set a grant to expire on a specified date. But this is not yet implemented*» **TODO**
 «*by the present design, expired grants will just be deleted - they won't be preserved inactive. Or should they be?*» **TODO**

Lower down you will find a list of existing grants. Press the [x] button next to any grant to delete it.

The **Manage group grants** screen has one extra function, which is the self-explanatory *Clone grants* tool.

Cloning grants from one group to another is useful if you have just created a new group that is going to have broadly similar privileges to an existing group: cloning the existing group's grants to the new group, and then editing those grants as necessary, can be quicker than granting each of the privileges to the new group one by one through the interface.

10.7 Running an open server

《more》

TODO

Even if all your corpora are completely open-access, it is often still a good idea to get users to sign up for individual accounts, so that they have access to their own private saved/categorised queries, their own query history, and so on.

《more》

TODO

《note the gotcha of blog spam in the macros form》

TODO

11 Using plugins

11.1 What is a plugin?

A plugin is a small program written to operate within the framework of a larger system. Many applications offer a plugin framework, which allows advanced users to add capabilities to the system which it does not possess out-of-the-box.

CQPweb is such an application. A CQPweb plugin is a chunk of code that is added to the system, which is then accessed in a predefined way by CQPweb to provide extra capabilities for users.

Some plugins are supplied with CQPweb. You can also write your own.

The different types of plugin do not have much in common, they just represent a set of things that we thought users of CQPweb might want to have an easy way to customise!

In all cases, to make use of a plugin, you must *register it* (see

In earlier versions of CQPweb, the *Custom Postprocess* type of plugin was implemented. Following a rework of the system for version 3.2.32, *Corpus Installer* and *Annotator* plugins are available. This chapter explains some general things about plugins, and some specific things about the currently-implemented plugins.

11.2 Types of plugin

11.2.1 Annotators

An Annotator plugin is one that tags a file. Normally, this will be a case of interfacing with an external program such as a POS tagger or lemmatiser (or even both!) but the plugin class could also be written to do the job itself.

Normally, Annotators will produce output in CWB vertical format - that is, p-attributes as columns, with XML tags for s-attributes on separate lines. That is because, when the user-corpus system is enabled, users can select Corpus Installers to run over their uploaded files, and the Corpus Installers are able to use Annotators. So the typical sequence would be:

- User uploads a text file or files;
- User selects a Corpus Installer and specifies their uploaded file(s) as input;
- CorpusInstaller calls an Annotator to run over the specified files;
- CorpusInstaller passes the output files (in vertical format), plus information on what attributes are used, back to CQPweb to set up the corpus.

However, ideally Annotator plugins should be written to work independently of this process; the Annotator should only “care” about the tagging of the text files, leaving the other details above to the Corpus Installer plugin.

To write an Annotator plugin, you need to write methods for (a) tagging files, and (b) describing the format of the resulting tagged file.

11.2.2 Corpus Installers

Corpus Installer plugins act as the controllers for installing user corpora. When a *system* corpus is installed, all information about its p- and s-attributes, the input files, and so on are provided by the admin user. For *user* corpora it is different: to make things more usable and friendly, they only need to select the Corpus Installer to use, which will then (a) manage any necessary tagging, by passing off to an Annotator plugin; (b) provide to CQPweb the information necessary to install the corpus.

Users' ability to install their own corpora is determined by which Corpus Installers they have permission to use.

《the other types of plugin, as the code to support them is completed》

TODO

11.3 Installing and registering plugins

《reiterate lib/plugins》

TODO

Once you have added your plugin file to the correct place on the system, you must *register* it so that CQPweb is aware of it, and can activate it where necessary. Note that CQPweb has a number of builtin plugins, but these won't be active in the system unless you register them.

The plugin registry is managed through the Admin Control Panel. Go to **CP >Plugins >Manage Plugins**. At the top of this screen is a form you can use to register a new plugin. You must specify three things:

- What PHP file contains the plugin class.
- A short description of the plugin (mostly for mnemonic purposes; users won't see it)
- Configuration data for the plugin, consisting of a set of zero or more key-value pairs.

The configuration data, if any, is passed to the constructor method of the plugin class. This allows the plugin to be flexible to a degree (you can register a plugin multiple times, with different configuration data, to take advantage of this).

Below the *Register a new plugin* form is a table containing details of all the registered plugins. You can delete the registry entries here (deleting a registry entry **does not** delete the PHP file which contains the plugin code).

11.4 Creating plugins

11.4.1 Introduction to writing plugins

To write a plugin, you will need a reasonable knowledge of programming in general, and at least a basic acquaintance with PHP specifically. It's beyond the scope of this manual to explain programming/PHP from the ground up; the PHP programming language is extensively documented at <http://php.net>.

A plugin takes the form of a PHP class (see <http://php.net/class>) that performs one or more defined tasks. For example, a *Custom Postprocess* plugin takes a defined CQPweb query and changes it ("postprocesses" it) in a certain way.

Each time you create a plugin, you should place it in a single file which has the same name as the plugin itself. So, for instance, to create a plugin called **MyPlugin** you should create the PHP file **MyPlugin.php**. You should then put the plugin file into the **plugins** subdirectory of the **lib** directory. Like all CQPweb code files, this new file must be readable by the username your webserver runs under.

The file should contain nothing except the class that represents your plugin. Each type of plugin has a separate PHP *interface* and your class must implement that interface to be recognised as a plugin.

Implementing an interface in PHP, as in some other object-oriented systems, means that the class must have methods that meet certain defined signatures. These are explained below. If your plugin does not meet the definition of the interface it implements, things will stop working.

The interface that your class needs to implement depends on what type of plugin you are writing, as follows:

Type of plugin	Interface to implement	Symbolic constant
Annotator	Annotator	PLUGIN_TYPE_ANNOTATOR
Corpus Installer	CorpusInstaller	PLUGIN_TYPE_CORPUSINSTALLER

So the overall shape of the plugin file will be like this:

```
<?php
class MyPlugin implements Postprocessor
{
    // you will need to add the methods here...
    // plus any member variables you want to use
}
```

There should be no whitespace before the leading `<?php` delimiter. If there is, CQPweb may stop working properly.

The builtin plugins provided with CQPweb provide examples of how plugins can be written.

11.4.2 Naming your plugin

All plugin names must be legal PHP class names (see <http://php.net/class>). They share a namespace with the internal CQPweb classes, so try to make sure you don't clash (PHP will crash if you do). To be safe, always prefix your plugin classes with a unique element (e.g. your name). Also, class names beginning with `My_` are guaranteed to be safe.

11.4.3 Methods your plugin must implement

This section describes each of the methods that you must implement to fulfil the demands of the interface for each type of plugin. You can, of course, have other methods if you want, but nothing outside the plugin will call them.

It is quite possible, and indeed expected, that in many cases you will want the major work of the plugin to be done outside PHP - for example, by a command-line utility or by a Perl or Python script. In which case, the PHP methods will be a thin wrapper round a call to an external system. CQPweb does not care about this, as long as it can use the methods in the interfaces to get the information it needs!

Note the documentation here is largely derived from the within-code documentation in `lib/plugin-lib.php`. If the comments there are different from what you find here, then the former should be considered definitive as they may be more recently updated or more extensive.

11.4.3.1 Methods required by every type of plugin These methods are part of the CQPwebPlugin interface, which is the parent of the interfaces for specific types of plugin.

- `public function __construct($extra_config = []);`

This method initialises the plugin. The `$extra_config` is an array of key-value pairs, taken from the plugin's entry in the plugin-registry (see [《crossref》](#)). This allows you to define multiple entries in the plugin registry from the same class, but with different configuration (for instance, the built-in Annotator plugin for the TreeTagger needs to be told what language to tag in). What the plugin *does* with these variables is determined by the code of the `__construct` method.

TODO

- `public function description();`

This method should return a string containing the title or short description of this plugin. This should be relatively short, and not contain any HTML.

- `public function long_description($html = true);`

This method should return a string describing the plugin. This may (but need not) include HTML formatting, where requested; and may either be the same as the string returned by `description()` or be longer. Line breaks ("`\n`") in the string may be rendered as HTML line breaks in some contexts, unless HTML is requested.

The method should respect the `$html` parameter by including HTML code in the string that is returned only if if this parameter is true.

- `public function status_ok();`

This method should return boolean **true** if the plugin has not encountered any error conditions, or boolean **false** if one or more error conditions has been encountered. If this method returns **false**, there should be something readable via the `error_desc()` method.

- `public function error_desc($html = true);`

This method should return a string describing the last encountered error.

If there has been no error, then it can return an empty string, or a message saying there has been no error. It doesn't matter which.

The method should respect the `$html` parameter by including HTML code in the string that is returned only if if this parameter is true.

11.4.3.2 Methods required by Annotator plugins

- `public function process_file($path_to_input_file, $path_to_output_file);`

Calling this method should tag the file specified as `$path_to_input_file`, placing the output at `$path_to_output_file`.

Both arguments should be relative or absolute paths. The method **SHOULD NOT** use CQPweb global variables. The input file **MUST NOT** be modified.

This function should return **false** if the output file was not successfully created. If the output file is partially created or created with errors, it should be deleted before **false** is returned.

If all goes well, the method should return **true**.

- `public function process_file_batch($input_paths, $path_to_output_file);`

This method should process a group of input files, placing the output into a single file. All the comments for `process_file()` apply here too. The return value should be boolean **true** or **false**, just as for that method.

The `$input_paths` paramter will be an array of strings, where each string is the path to an input file.

- `public function output_size();`

This method should return the size of the last output file created as an integer count of bytes, or zero if no file has yet been specified.

11.4.3.3 Methods required by Corpus Installer plugins

- `set_max_input_tokens($max);`

Corpus Installers are controlled by permissions that specify how much text the user is allowed to index at one time. CQPweb will use this method to tell the Corpus Installer what the maximum size is that it should allow. The plugin should implement this method so as to store this value, and utilise it at the appropriate time (usually: after tagging, before indexing).

The parameter is a count of tokens; if zero or a negative limit is set, no restriction at all should be applied.

- `public function set_corpus_name($name);`

CQPweb will call this method to pass in to the Corpus Installer plugin the name (lowercase CQPweb handle) of the corpus it is begin used to create. The method should store this value. The name of the corpus is needed to create some of the corpus-installation setup SQL statements.

- `public function add_input_file($path);`

This method must add the file at the path provided to its list of files to be used as input for installation. (The Corpus Installer plugin may pass these through to an Annotator, or else use them directly.)

The method should also accept an array of strings - all representing paths to input files.

- `public function do_setup();`

```
/**
 * Run setup - that is, anything that needs to be done to get files
 * ready to be encoded. This might include tagging, or even building a corpus.
 *
 * @return bool    True if setup worked OK; false if not.
 */
```

- `public function do_cleanup($delete_input_files = false);`

```
/**
 * Run cleanup, e.g., deleting temporary files, if any.
 * @param bool $delete_input_files If true, the files specified using
 *                               CorpusInstaller::add_input_file() will be deleted.
 */
```

- public function get_charset();

This method should return the CWB string indicator for the charset of the corpus text.
(These days, usually “utf8”.)

- public function get_p_attribute_info();

```
/**
 * Gets information about the p-attributes for cwb encoding
 * (an array of strings for use with -P with cwb-encode).
 * @return array
 */
```

- public function get_s_attribute_info();

```
/**
 * Gets information about the s-attributes for cwb encoding
 * (an array of strings for use with -S with cwb-encode).
 */
```

- public function get_annotation_bindings();

```
/**
 * Gets the annotation bindings that will be used in the created
 * corpus (for Simple Queries).
 * @return array A hash with one or more keys referring to
 *              'special' syntax in CEQL, each mapping to
 *              the p-attribute that will be searched.
 *              For instance, 'primary_annotation'=>'pos'.
 *              The "tertiary_annotation_tablehandle"
 *              is the only non-p-attribute binding; it
 *              must be for a table that actually exists
 *              on the system, if not, declare_maptable()
 *              can be used.
 */
```

- public function declare_maptable();

```
/**
 * Hash of simple-pos to p-attribute regex. CQPweb will add it
 * as a mapping table to the system. If you don't want to bother
 * with this, just have an empty function (which is what the
 * CorpusInstallerBase does).
```

```
*
* @return array
*/
```

《*change code comments to writeup!*》

TODO

11.4.4 Methods you can inherit

The plugin system includes “base classes” which contain implementations of several important tasks for various types of plugin. If you write your plugin to inherit from such a base class, then you do not have to write code for these tasks yourself. **It is highly advisable to make use of these base classes.** (But if you don’t want to, you don’t have to, even if you inherit from a base class: any method in a base class can be overridden in a child class simply by creating a method of the same name.)

The base classes all have the same name as the corresponding interface with the addition of **Base** at the end. The facilities they provide are listed below.

11.4.4.1 CQPwebPluginBase CQPwebPluginBase is the parent of the other base classes. Therefore, its affordances are passed on to those base classes. This base class provides...

- A default `__construct()` method, which simply stores each key-value pair from the `$extra_config` as an object variable; the variable name is the key string, the variable value is the value.
- A default implementation for the `long_description()` method (making it return the same as `description()`, which all plugins must still provide).
- A default error-report system, which provides implementations for `status_ok()` and `error_desc()` as well as an extra method, `raise_error()`, which sets the error state (and logs the message passed as its sole parameter to be returned by subsequent calls to `error_desc()`). `raise_error()` can only be called from *within* a child class.

11.4.4.2 AnnotatorBase This base class provides...

- A text ID generation system. All CQPweb input data must have `<text>` tags with `id="..."` attributes. Annotators generally need to add these, as they are likely not to be present in the input, or to have been removed by the tagger. The AnnotatorBase provides multiple ways to generate an ID for each input text, that can then be incorporated into the output. The function to access this is `AnnotatorBase::get_next_text_id()`.
- A pair of support functions, `validate_read_paths()` and `validate_write_paths()`, each of which checks an array of file paths to make sure they are, respectively, readable/writable.
- A simple implementation of `output_size()` which returns the value of an internal variable, `$bytes_in_output`. Child classes which wish to use this need to actually set that class property!
- Null implementations of `output_annotation_list()` and `output_xml()`, which always return **false**, allowing child classes to not bother implementing these methods.

11.4.4.3 CorpusInstallerBase This base class provides...

- Multiple systems for creation of the necessary SQL statements / CWB declarations for annotations and XML (p-attributes and s-attributes).
 - Via `declare_annotation()` and `declare_xml()`: methods that the child class can use to pass in information about the attributes.
 - Via `declare_annotations_from_template()` and `declare_xml_from_template()`: methods which can be passed template ID numbers, from which the information will be taken.
 - Via `declare_content_from_annotator()`: a method to which a child class can pass its Annotator object, whose interface will then be interrogated to get information about the attributes.
- A default system for setting, and then retrieving, CEQL bindings: the `set_binding()` method allows them to be set, and the implementation of `get_annotation_bindings()` returns them.
- Default implementations of the following methods, which will work alongside any of the different systems mentioned above, and report the declared structures to CQPweb:
 - `set_max_input_tokens()` - stores a token limit in an internal variable, which the child class can use, or which can be applied with `restrict_input_data()` (see below).
 - `set_corpus_name()` (plus an internal utility function, `check_for_corpus_name()`, to make sure it has been called)
 - `get_sql_for_corpus()`
 - `get_sql_for_abort()`
 - `get_p_attribute_info()`
 - `get_s_attribute_info()`
 - `get_infile_info()`
 - `get_charset()` - returns the internal variable `charset`; the child class will need to set this
 - `get_xml_datatype_check_needed()` - returns the internal variable `xml_datatype_check_needed`, which is `false` by default; the child class can set this to `true` if desired.
- Default implementations of the following additional methods:
 - `add_input_file()` - for use alongside `get_infile_info()`
 - `do_cleanup()` - which deletes temporary files (and does nothing else).
- A utility function, `restrict_input_data()`, designed to be called by `do_setup()`, which makes sure the input files do not exceed the limit set by `set_max_input_tokens()`
- An alternative to `set_max_input_tokens()`, the method `set_restriction_from_privilege()`, which can be passed a privilege object, from which the object will extract a token limit.
- A null implementation of `declare_maptable()`, allowing child classes to not bother implementing this method.

11.4.5 An API for plugin writers

When you write a plugin, the whole of CQPweb's internal function library is theoretically available for you to call. However, unless you *really* know what you are doing, it is not recommended to just start calling functions all over the place. Something may go wrong.

That said, there are clearly bits of information that you might need inside a plugin that are not provided via the methods' parameters. For instance, in a Custom Postprocess where you need to keep or reject each concordance hit, it would be nice if you could find out what each concordance hit actually contains! You could do this by accessing CQPweb's interaction layer with the CQP back-end for yourself, but that's prone to all those problems discussed above

So CQPweb provides a set of helper functions that you can call to access other bits of info in such a way that it's "guaranteed" not to mess things up (barring bugs we don't know about yet). There are currently three such functions (not tested as of the current version), all designed to be of use in writing Custom Postprocesses. Their function prototypes and internal documentation are provided below.

```
/*
 * Returns a path that can be used as a temporary filename by a plugin.
 * The plugin is responsible for making sure it gets deleted.
 * (It will be pre-created as a temporary file.)
 */
function pluginhelper_get_temp_file_path();

/**
 * Gets a full concordance from a set of matches.
 *
 * The concordance is returned as an array of arrays. The outer array contains
 * as many members as the $matches argument, in corresponding order. Each inner array
 * represents one hit, and corresponds to a single group of two-to-four integers.
 * Moreover, each inner array contains three members (all strings): the context
 * before, the context after, and the hit itself.
 *
 * The $matches array is an array of arrays of integers or integers as strings,
 * in the same format used to convey a query to a custom postprocess.
 *
 * You can specify what p-attributes and s-attributes you wish to be displayed in the
 * concordance. The default is to show words only, and no XML. Use an array of strings
 * to specify the attributes you want shown in each case.
 *
 * You can also specify how much context is to be shown, and the unit it should be
 * measured in. The default is ten words.
 *
 * Individual tokens in the concordance are rendered using slashes to delimit the
 * different annotations.
 */
function pphelper_get_concordance($matches,
                                $p_atts_to_show = 'word',
                                $s_atts_to_show = '',
                                $context_n = 10,
                                $context_units = 'words'
                                );
```

```

/**
 * Determines whether or not the specified corpus position (integer index) occurs
 * within an instance of the specified structural attribute (XML element).
 *
 * Returns a boolean (true or false, or NULL in case of error).
 */
function pphelper_cpos_within_structure($cpo, $struc_attribute);

/**
 * Gets the value of a given positional-attribute (word annotation)
 * at a given token position in the active corpus.
 *
 * Returns a single string, or false in case of error.
 */
function pphelper_cpos_get_attribute($cpo, $attribute);

```

11.5 Builtin plugins

Some plugins are provided with the CQPweb distribution. However, they will not be available in the web interface until you add them to the system, as explained in section 11.3.

There are two types of builtin plugin: those supplied purely as examples of how to write a plugin, and those supplied because they are expected to be generally useful to lots of different users. All can be found in the `lib/plugins` directory.

But beware! These plugins may be under development (or may relate to types of plugin whose integration into CQPweb is not yet complete).

11.5.1 BasicTokeniser

This is a very simple, but usable, Annotator. It processes input files to produce tokenised output in CWB vertical format.

Extra configuration:

- Does not make use of any extra configuration values.

11.5.2 TreeTagger

This Annotator is designed to interface with the TreeTagger software (see <http://www.cis.uni-muenchen.de/~schmid/tools/TreeTagger/>).

It makes certain assumptions about how the TreeTagger installation files are to be found on the system; if your TreeTagger installation is laid out differently, then you might not be able to use this plugin.

The class contains configuration information for a large number of the parameter-files made available on the TreeTagger website for various languages. Again, if you haven't downloaded the parameters for a given language, you won't be able to use the plugin to tag files in that language.

Extra configuration:

- `tt_no_s_tags` - boolean, if **true**, sentence tags (**s**) will not be added to the TreeTagger output.

- `tt_show_unknown_lemma` - boolean, if **true**, words whose lemma is unknown will be tagged as `<unknown>`, rather than supplying the wordform as “its own lemma”. The default behaviour is for both this and the previous setting to be **false**, which is contrary to the default TreeTagger behaviour, but better fitted for use with CQPweb.
- `tt_bin_path` - string, path to the TreeTagger installation folder (i.e. the directory which contains the `bin`, `cmd`, and `lib` subdirectories).
- `language` - string, sets the language to tag in. The valid language labels are abbreviated from the descriptions of the corresponding parameter sets as distributed via the TreeTagger website. Possible languages are those used as keys in the `LANG_INFO` array (which can be found at the bottom of the TreeTagger class; a utility function (static class method) is provided to check language-identifier strings (`TreeTagger::is_valid_language()`)).

11.5.3 UcrelTagger

This Annotator is a very thin wrapper around the Lancaster University UCREL research centre's dual taggers, namely CLAWS (part-of-speech) and USAS (semantic), via a controller script that generates CQPweb-ready output.

See:

- <http://ucrel.lancs.ac.uk/claws/>
- <http://ucrel.lancs.ac.uk/usas/>

Extra configuration:

- Many possible items, not documented here (yet).

11.5.4 BasicVrtInstaller

This installer assumes input files that are already in vertical format (and therefore do not need to be tagged).

Extra configuration:

- Does not make use of any extra configuration values.

11.5.5 SimplePlaintextInstaller

This is a simple corpus installation manager assuming plain-text input files. It calls a configurable Annotator in a standard way.

Extra configuration:

- `annotator_plugin_id` - integer ID for the plugin registry entry of the Annotator plugin to use to tag the files.
- `vrt_file_path` - (optional) path specifying where to place the output file; if not given, a temporary location will be used.

11.5.6 StandardToolInstaller

This is a corpus installer for use with the UcrelTagger and TreeTagger Annotators.

Extra configuration:

- `tool` - a string (“UCREL” or “TreeTagger”) to tell the plugin which annotator to use.
- `language` - a string that will be passed through to the TreeTagger (see above).
- `semtag-resources` - a string that will be passed through to the UcrelTagger (see above).
- `annotation_template_id` - integer ID of the annotation template describing the tagger output.
- `xml_template_id` - integer ID of the annotation template describing the tagger output.

The last two of these are optional. If they are not supplied, the plugin will find the templates itself (this will work fine as long as you have installed the default annotation and XML templates; see [6.7](#), [6.9](#)).

11.6 Permissions for plugins

《*write this!*》

TODO

12 Extensible CQPweb

12.1 Introduction

This chapter explores the various ways of extending CQPweb by adding your own code to support features that it doesn't possess out-of-the-box.

13 Using the CQPweb API

13.1 Introduction

Use of the CQPweb API is, strictly speaking, not a system administration topic. However, effective use of the API requires an in-depth knowledge of how the system works, and so it is convenient to include the topic in this manual.

14 Updating CQPweb

14.1 The update process

In general, to update CQPweb to a new version, there are three steps to take. The first step is always necessary, but the second two steps will only be needed when there have been major changes. The steps are:

- Check for new dependencies;
- Update the code;
- Update the database;
- Update the configuration file.

New dependencies on other pieces of software are rarely added. They will always be noted here when they are added:

- In version 3.1.7, a requirement for the R statistical software to be available was added (it was previously optional).

To **update the code**, simply get a new copy of the code from the CWB website, and copy the files over your existing installation. Be careful not to alter any of the folders relating to corpora you have installed. If your CQPweb is running from a Subversion checkout, then the following command, when run from the base directory, will typically complete all requisite actions:

- `svn update`

To **update the database**, you need to run the admin script `upgrade-database.php`. There is further information on this script in section 5.13. If you are upgrading a very old version of CQPweb, you may need to perform some updates manually; see section 14.2 below.

Updating the configuration file rarely needs to be done, as we make efforts to keep it consistent. So far the only change in the format of the configuration file has been between 3.0.16 and 3.1.0. The changes are discussed in section 2.6.

Some updates may require additional actions as well as these three steps; if so, these are explained below.

14.2 Updating the database from very old versions

Early versions of CQPweb required MySQL schema changes to be implemented manually. A full list of the changes from version 2.15 to version 3.0.16 is given below. In each case, the MySQL command given needs to run against the CQPweb database when logged in either as root or as the CQPweb database user. You need to run all the commands between the version you are upgrading *from* and the version you are upgrading *to*, inclusive (i.e. the list that follows is cumulative). Sometimes other steps are necessary, and they are listed too.

- Going from 2.15 to 2.16
 - `alter table user_settings drop key `username`;`

- alter table user_settings add primary key (`username`);
- alter table user_settings add column `password` varchar(20) default NULL;
- Going from 2.16 to 3.0.0
 - No database changes.
- Going from 3.0.0 to 3.0.1
 - CREATE TABLE `corpus_categories` (`idno` int NOT NULL AUTO_INCREMENT, `label` varchar(255) DEFAULT '', `sort_n` int NOT NULL DEFAULT 0, PRIMARY KEY (`idno`)) CHARACTER SET utf8 COLLATE utf8_general_ci;
 - ALTER TABLE `corpus_metadata_fixed` MODIFY COLUMN `corpus_cat` int DEFAULT 1;
 - Since this upgrade will wipe out your existing corpus categories, you must re-add them.
- Going from 3.0.1 to 3.0.2
 - DROP TABLE IF EXISTS `xml_visualisations`;
 - CREATE TABLE `xml_visualisations` (`corpus` varchar(20) NOT NULL, `element` varchar(50) NOT NULL, `xml_attributes` varchar(100) NOT NULL default '', `text_metadata` varchar(255) NOT NULL default '', `in_concordance` tinyint(1) NOT NULL default 1, `in_context` tinyint(1) NOT NULL default 1, `bb_code` text, `html_code` text, key(`corpus`, `element`)) CHARACTER SET utf8 COLLATE utf8_bin;
 - Go to System diagnostics and run the check for ``PHP inclusion files''
- Going from 3.0.2 to 3.0.3
 - ALTER TABLE `xml_visualisations` DROP KEY `corpus`;
 - ALTER TABLE `xml_visualisations` ADD COLUMN `cond_attribute` varchar(50) NOT NULL default '';
 - ALTER TABLE `xml_visualisations` ADD COLUMN `cond_regex` varchar(100) NOT NULL default '';
 - ALTER TABLE `xml_visualisations` ADD PRIMARY KEY (`corpus`, `element`, `cond_attribute`, `cond_regex`);
 - ALTER TABLE user_settings add column thin_default_reproducible tinyint(1) default NULL;
 - UPDATE user_settings set thin_default_reproducible = 1;
- Going from 3.0.3 through to 3.0.16
 - Nothing.

From version 3.1.0 onwards, the database template can be updated automatically by running the admin script `upgrade-database.php`. However, before you do this, you must manually apply all the relevant updates from the list above if you are moving from a version earlier than 3.0.16.

14.3 Updating from version 3.0.16 to version 3.1.0

This is a major upgrade, and careful adjustment will be needed. The following steps should be followed *in order*.

Step 1. Update the code (see section [14.1](#)).

Step 2. Update your configuration file, to take account of the changes in the format listed in section [2.6](#).

Step 3. Update the database using the upgrade script described in section [5.13](#).

The above steps are essential. Some further steps are useful if you previously managed users/groups and their access rights using CQPweb's interface to Apache:

- Restore your previous groups using `load-pre-3.1-groups.php` (see section [5.9](#))
- Restore your previous privileges using `load-pre-3.1-privileges.php` (see section [5.10](#))
- Either (a) remove all `.htaccess` files from the web folders for particular corpora and/or (b) turn off the use of `.htaccess` files within the CQPweb web folder completely (see section [1.9](#) for a full account of setting up Apache under CQPweb 3.1 and higher).

Note that the first two of the above steps must be followed *in that order*, and *after* you have upgraded the database.

14.4 Updating from version 3.1.7 or earlier to version 3.1.8 or later

In version 3.1.8, the limit on how big a frequency list a user can create was changed from a single global value to a configurable privilege.

This means that, having upgraded to 3.1.8 or higher, you will need to use the Admin Control Panel (see [3](#)) to add at least one privilege of this sort to your system, and assign it to users/groups.

The four default privileges of this sort enable users to create frequency lists for subcorpora of 1 million, 10 million, 25 million and 100 million tokens.

At least one privilege of this sort should normally be assigned to the “everybody” group, or else these users will not be able to create frequency lists for subcorpora at all.

14.5 Updating from version 3.1.8 or earlier to version 3.1.9 or later

In version 3.1.9, the “Analyse Corpus” function was added to CQPweb. In order to add the webpage supporting this function to existing corpora, you should go to the Admin Control Panel (see [3](#)), and run the “Check corpus PHP inclusion files” function (found under *System diagnostics*).

14.6 Updating to version 3.2.0

There were substantial architectural changes in version 3.2.0, which is why it was a major version change. Although these changes caused little to be different on the surface, they were an essential step towards future developments.

If all goes well in the upgrade you will never need to know what the changes are. However, if something goes wrong, you may need the following information on the architectural changes to be able to effect a manual repair. The differences are as follows:

- In earlier versions, each corpus had a separate web folder inside the main CQPweb directory, with a set of PHP scripts in it. In 3.2.0 this was changed so that there was a single location for the corpus-interface scripts (the built-in sub-directory `exe`) and each corpus's web folder is now a symbolic link to `exe`.
- In earlier versions, a lot of important information about each corpus was stored in its “settings” file, stored in its web folder (filename `settings.inc.php`). In 3.2.0 all this information is transferred to the database, and the settings files are removed.
- In earlier versions, CQPweb relied on the CWB registry to keep track of s-attributes (XML elements/attributes). In 3.2.0, CQPweb has a database structure that keeps track of this information.

To upgrade from version 3.1.16 (or earlier) to version 3.2.0 (or later), you should follow these steps.

Step 1. It's recommended to take a backup copy of the entire web-directory containing the code of CQPweb and the web folders of the individual corpora before starting.

Step 2. Update the code (see section 14.1).

Step 3. Update the database using the upgrade script described in section 5.13.

Normally, the database upgrade script will bring the system right up to the current version of the code. However, it will always stop at version 3.2.0, to allow you to map across the corpus/XML settings from the earlier format. If your code version is above 3.2.0, you will need to run the database upgrade script *again* after finishing the rest of these steps.

Step 4. Run the special script to transfer existing corpus/XML settings to the new format.

The script for this step is listed in 5.11. To run it, go into the `bin` subdirectory, and enter the following command:

- `php load-pre-3.2-corpsettings.php`

The script goes through your list of corpora, and for each corpus it finds, it takes the following actions:

- First, it loads that corpus's settings file and inserts the information it finds into the MySQL database.
- Second, it attempts to replace the corpus's web folder with a symbolic link to the `exe` folder.
- Third, it interrogates the CWB registry to discover the corpus's s-attributes, and creates a record of each attribute in the database.

Step 5. You may see error messages from the special script if any step of the process does not complete correctly, so the next step is to address these messages by making manual adjustments.

- If no settings file is found for a particular corpus, this is probably not a problem: the setup of the corpus in question was already broken.
- If the script reports that, for a particular corpus, it could not replace the web directory with a symlink, then the manual fix for this is to run the following shell commands within the CQPweb main directory:

– `rm -r corpus`

```
– ln -s exe corpus
```

but with the actual corpus handle instead of “corpus”! You may then need to adjust the ownership/permission of the resulting symlink (see notes on web-directory ownership and permissions in 1.3).

Step 6. In earlier versions, as noted above, the “settings” file stored key information. It was possible for system administrators to manually add variables/code to these files, to add extra tweaks to the interface on a per-corpus basis. This has long been *highly inadvisable* due to the increasing complexity of the system, but from version 3.2.0 onwards, it is no longer possible. So your final step, *if and only if you have made any such modifications*, is to review your old `settings.inc.php` files (from your backup created as per above!) and double-check the effects of the loss of your manual tweaks. Some of them may be replicable through the usual administrative tools described in the rest of this manual.

If you never manually edited any of the settings files, then you have nothing to do under this step.

14.7 Updating to version 3.2.4

When you update to version 3.2.4 all users who are currently logged in will be automatically logged out. This is due to a change in the database format regarding the storage of login tokens.

14.8 Updating to version 3.2.6

Running the database upgrade script for the update to 3.2.6 can take a long while, especially if you are updating a server with lots and lots of users, or if CQPweb has been installed for a very long time. DO NOT abort the update script - let it run to completion. If you abort it before it has finished, your database may end up in a half-and-half state, in which case it would become very difficult to repair it without losing some of your users’ data.

14.9 Updating to version 3.2.23

In version 3.2.23, the limit on how large a file a user can upload was changed from a single global value (hardcoded as 2 MB) to a configurable privilege.

This means that, having upgraded to 3.2.23 or higher, you will need to use the Admin Control Panel (see 3) to add at least one privilege of this sort to your system, and assign it to users/groups.

The three default privileges of this sort are for file size limits of 0.5 MB, 1 MB, and 2 MB.

At least one privilege of this sort should normally be assigned to the “everybody” group, or else these users will not be able to upload files at all.

14.10 Updating to version 3.2.32

In version 3.2.32, the STTR statistic was added to the information stored for corpora.

You should run the following command in order to add the STTR to your existing corpora:

```
php execute-cli.php update_all_missing_sttr
```

Be aware this can take a LONG time to run, and must only be done *after* you have upgraded the database!

(A message to this effect is also printed when you run the database upgrade script.)

The STTR for *new* corpora will be calculated at corpus-installation time.

This version also adds the system allowing users to install their own corpora. However, before this is possible, you will need to set up the appropriate plugins, and assign upload and corpus creation privileges; likewise you will need to set the `$user_corpora_enabled` configuration variable (see [2.3.7](#)).