

Extensibility as a focus for corpus analysis software

The CQPweb plugin framework

Andrew Hardie

ESRC Centre for Corpus Approaches to Social Science

Lancaster University UK

@HardieResearch



The background

- Corpus Workbench (CWB)
 - Open project: <http://cwb.sf.net>
- CQPweb : web-based UI to CQP
 - Interface design follows BNCweb
 - 2008: Internal at Lancaster <https://cqpweb.lancs.ac.uk>
 - 2010: open to other users
 - As part of CWB: open source
 - **Many** other servers across the world

Why EXTENSIBILITY?

Features in demand

- Users to be able to upload their own corpora
- Interface to corpus construction systems
- More options to customise download formats
- Wider range of analyses
 - Of concordances, of corpora
 - Traditional, experimental
 - Statistics, visualisations

Plugins

- What is a **plugin**?
- CQPweb modular elements
 - Download
 - Install corpus
 - Annotate text data
 - Process concordance to another form
- ... places where the normal process could be replaced *as a unit* by another unit that connects up the right way

Downloader Plugin

The screenshot shows a dropdown menu on the left with 'Download...' selected. The main interface has a red header 'Download concordance' and a light blue area with two buttons. Below that is a red header 'Detailed output options' and a grey 'Formatting options' section. The settings include:

- Choose operating system line-break style: Windows
- Print short handles or full values for text categories: full values
- Mark query results as <<< result >>>: No
- Size of context: 10 words each way

Now plugins!



Downloader plugins activated in Admin interface

Existing Plugins					
ID	Type	Description	Class	Extra setup	Acti
7	QueryDownloader	Filemaker download button	TypicalForFilemaker		[Activate]
8	QueryDownloader	Word/Excel download button	TypicalWordExcel		[Activate]

Annotator and Corpus Installer

Existing Plugins				
ID	Type	Description	Class	Extra setup
1	CorpusInstaller	Nexis Corpus Importer	NexisCorpusCreate	<pre> semtag-resources => "/opt/ucrel/usas/share/semtag" annotation_template_id => "4" xml_template_id => "11" nexis_script_location => "/home/hardiea/code-checkouts/build-tools/nexis-to-corpus.php" </pre>
2	CorpusInstaller	Install corpus with CLAWS/USAS	StandardToolInstaller	<pre> tool => "UCREL" semtag-resources => "/opt/ucrel/usas/share/semtag" annotation_template_id => "4" xml_template_id => "2" </pre>
3	CorpusInstaller	Install corpus with TreeTagger (English)	StandardToolInstaller	<pre> tool => "TreeTagger" language => "english" annotation_template_id => "2" xml_template_id => "2" </pre>
4	CorpusInstaller	Install corpus with TreeTagger (Portuguese)	StandardToolInstaller	<pre> tool => "TreeTagger" language => "portuguese" annotation_template_id => "2" xml_template_id => "2" </pre>
5	CorpusInstaller	Install corpus with TreeTagger (Brazilian Portuguese)	StandardToolInstaller	<pre> tool => "TreeTagger" language => "portuguese-br" annotation_template_id => "2" xml_template_id => "2" </pre>
6	CorpusInstaller	Install corpus with TreeTagger (German)	StandardToolInstaller	<pre> tool => "TreeTagger" language => "german" annotation_template_id => "2" xml_template_id => "2" </pre>

Annotator configuration

Corpus setup configuration

Interface

Install a new corpus

You have access to the following Corpus Installer systems. Please select the installer you want to use.

Nexis Corpus Importer

Install a corpus up to **100,000,000** tokens.

[Details] [Select]

Install corpus with CLAWS/USAS

Install a corpus up to **1,000,000** tokens.

[Details] [Select]

Install corpus with TreeTagger (English)

Install a corpus up to **1,000,000** tokens.

[Details] [Select]

Install corpus with TreeTagger (Portuguese)

Install a corpus up to **1,000,000** tokens.

[Details] [Select]

Install corpus with TreeTagger (Brazilian Portuguese)

Install a corpus up to **1,000,000** tokens.

[Details] [Select]

Install corpus with TreeTagger (German)

Install a corpus up to **5,000,000** tokens.

[Details] [Select]

Select files

Select files from your upload area to include in the corpus. If the files contain more data than you are allowed to install, the overflow will be left out of the corpus.

Include?	Filename	Size	Date modified
<input type="checkbox"/>	example.txt	2.8 KB	2020-Apr-27 06:15
<input type="checkbox"/>	sde.txt	3.0 KB	2019-Jul-31 13:27
<input type="checkbox"/>	Times March 2019.zip	8.9 MB	2019-May-02 11:46

Set options

Specify a title for your corpus

(Optional) Select the language used in your corpus ▼

Tick here if the main script in the corpus is right-to-left (*Arabic, Hebrew, Aramaic, etc.*)

Select a colour scheme for the interface to your corpus (*will be overridden if you have the "monochrome" option enabled*) ▼

Tick here to receive an alert by email when your corpus is ready to use.

Run corpus installer



Once indexed...

Menu	CQPweb User Page						
Your account	Your installed corpora						
Overview	Corpus	Indexing date	Size			Disk space	Actions
Settings			Tokens	Types	Texts		
Macros							
Corpus permissions	dsaaedx (_00137)	2019-Jun-10 06:12	4,032,258	87,021	5,758	138 MB	[Delete corpus]
Your files and corpora							
View your corpora	A test! 9 (_0016e)	2020-Jun-25 07:20	494	273	1	380 KB	[Delete corpus]
Install a new corpus	Your corpus data storage allowance						
View install queue	Amount of disk space your installed corpora currently take up:					139 MB	
Manage your files	Your disk space allocation for installed corpora:					10,240 MB	

How do plugins work?

- Each type has certain rules about what it *receives* and what it *produces*
 - Downloader
 - *Receives* a concordance from command-line CQP
 - *Produces* the file that the user gets as download
 - Annotator
 - *Receives* input text (file uploaded by user)
 - *Produces* vertical format text for import
 - Installer
 - *Receives* list of input files
 - *Produces* settings for corpus installation plus annotated versions of the files.

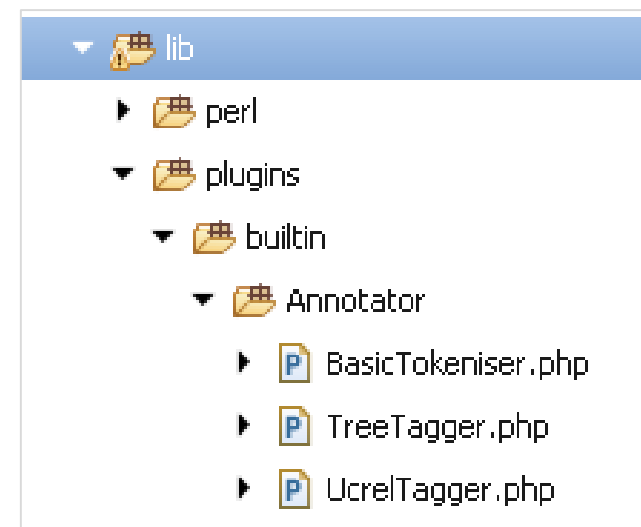
How do plugins work?

- Write a PHP file containing a Class for the plugin
- Put the file in the right place in CQPweb's code

```

23
24 /**
25  * Annotator plugin for the UCREL toolchain with CLAWS/USAS.
26  *
27  * Note this assumes the UCREL toolchain is installed as it is
28  * on our systems at Lancaster.
29  *
30  * Elsewhere, it probably won't work without tweaks.
31  */
32 class UcrelTagger extends AnnotatorBase implements Annotator
33 {
34     private $script = 'ucrel-tagger-toolchain';
35     private $script_opts = [];
36
37
38     // TODO in future there will be 4 modes, depending on the Claws resources (maybe also USAS resoruces? )
39
40     const MODE_CLASSIC = 1; /* classic lexicon used for all those years in Wmatrix etc., based on Written BNC1994 */
41     const MODE_WRITTEN = 2; /* including new Writ BNC2014 vocab (added to classic lexicon) */
42     const MODE_SPOKEN = 3; /* including new Spok BNC2014 vocab (added to Old BNC Spoken lexicon) */
43     const MODE_EMODENG = 4; /* including new Shakespearean vocab (added to classic lexicon) */
44     // all depends on the toolchain knowing about them, of course.
45
46
47     private $mode = self::MODE_CLASSIC;
48
49     public function __construct(array $extra_config = [])
50     {
51         /* the only extra params we accept are (a) options for ucrel-tagger-toolchain, or (b) a few specials. */
52         foreach($extra_config as $param => $val)
53         {
54             switch($param)
55             {
56                 case '_special_use_this_script':
57                     $this->script = realpath($val);
58                     break;

```





```
199- /**
200  * Interface for CorpusInstaller Plugins.
201  *
202  * A CorpusInstaller Plugin is a driver for corpus setup by users.
203  *
204  * It carries out all steps in the procedure. This may include running an Annotator,
205  * checking file formats, doing cleanup, etc. Files may be generated anew,
206  * or created from operations on the user's existing files.
207  *
208  * The plugin does not actually run cwb-encode and friends.
209  * But it does supply options for setup based on what it's done.
210  */
211- interface CorpusInstaller extends CQPwebPlugin
212- {
213-     /**
214      * Make the plugin restrict the amount of data by tokens.
215      * If 0 or a negative limit is set, no restriction at all applies.
216      * @param int $max
217      */
218     public function set_max_input_tokens(int $max) : void;
219
220
221-     /**
222      * Tell the plugin what the CQPweb handle ofr the corpus will be.
223      * (Necessary for it to generate the correct SQL statements.)
224      * @param string $name  A lowercase CQP-corpus name.
225      */
226     public function set_corpus_name(string $name) : void;
227
228
229-     /**
230      * Ascertain whether this CorpusInstaller needs input files
231      * (some plugins generate data externally). This controls
232      * whether the file selector is enabled/disabled, and whether
233      * any input files are passed in.
234      *
235      * @return bool  True if input files are needed.
236      */
237     public function needs_input_files() : bool;
238
239-     /**
240      * Add a given path to the list of files to be used as input to the installation.
241      * @param string|array $path  A path, or an array of paths.
242      */
243     public function add_input_file($path) : void;
244
245-     /**
246      * Run setup - that is, anything that needs to be done to get files
247      * ready to be encoded. This might include tagging, or even building a corpus.
248      *
249      * @return bool  True if setup worked OK; false if not.
250      */
251     public function do_setup() : bool;
252
```

INTERFACES

*Define functions
that a plugin Class
must provide*

Only PHP?

- Right now, yes
- RFace
- PyFace

RFace in action: Log Ratio

```

1276 @/**
1277 * Calculates a Z-unit to be used as the offset for the LR confidence interval.
1278 *
1279 * @param float $alpha The alpha to use. Caller should adjust for familywise if necessary. Defaults to 0.05.
1280 * @param RFace $r An R Face object to use. Defaults to null, in which case the function starts an R slave itself.
1281 * @return float The Z unit (to use, for instance, embedded into an SQL query).
1282 */
1283 function calculate_Z_for_LR_confinterval(float $alpha = 0.05 , ?RFace $r = null) : float
1284 {
1285     switch ($alpha)
1286     {
1287         /* optimise for some frequently-used cases
1288          * (see calculate_LL_threshold() for more notes on this!)
1289          *
1290          * alpha here changes to the CI width (0.05 = 2.5% each way...)
1291          */
1292         case (float)'0.05':          $Z_unit = 1.959964;    break; /* 95% CI; in R: qnorm(0.025, lower.tail=FALSE) evals to 1.959964 */
1293         case (float)'0.01':          $Z_unit = 2.575829;    break; /* 99% CI; etc. */
1294         case (float)'0.001':         $Z_unit = 3.290527;    break;
1295         case (float)'0.0001':        $Z_unit = 3.890592;    break;
1296         case (float)'0.00001':       $Z_unit = 4.417173;    break;
1297         case (float)'0.000001':      $Z_unit = 4.891638;    break;
1298         case (float)'0.0000001':     $Z_unit = 5.326724;    break;
1299
1300     default:
1301         global $Config;
1302         $internal_r = is_null($r);
1303         if ($internal_r)
1304             $r = new RFace($Config->path_to_r);
1305
1306         list($Z_unit) = $r->read_execute(sprintf("qnorm(%E, lower.tail=FALSE)", $alpha/2.0));
1307
1308         if ($internal_r)
1309             unset($r);
1310
1311         break;

```

PyFace... coming soon honest

- Plugin writing in R / Python
 - Write a very small PHP class which fulfils the Plugin Interface
 - Make the PHP class pass commands to an RFace or PyFace
 - Write the main program in [language of choice]
 - Retrieve the results from the RFace / PyFace
 - Make the PHP class return the result to CQPweb

Thank you for your attention

