# Cambridge Sketch Engine

## Advanced Help (v.2.0)

## Contents:
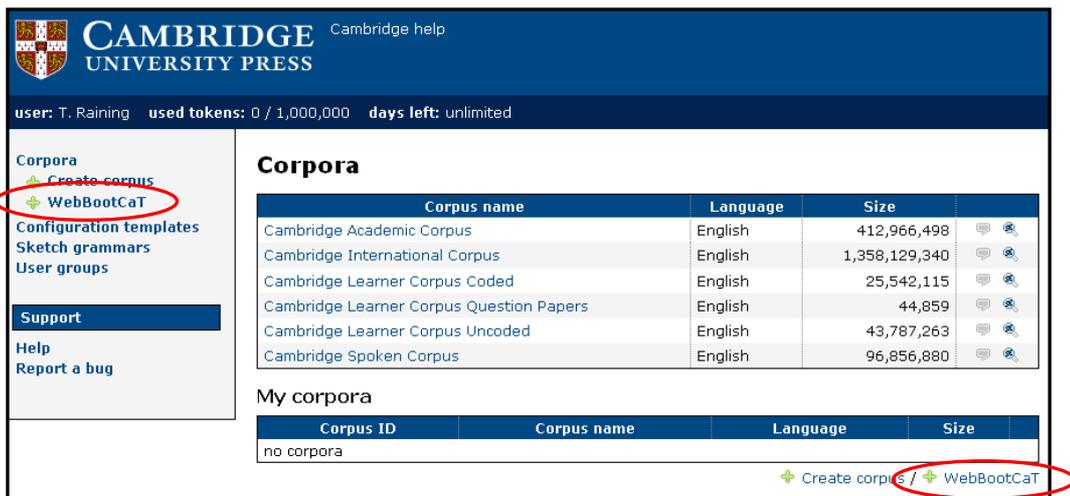
# 1. Creating corpora from the web - WebBootCaT

The WebBootCaT function allows you to build corpora from the web. WebBootCaT can be used to complement the existing Cambridge Corpus resources – it can be used to compile a corpus on a particular topic or subject not currently found within the Corpus (for example, on *new technology*, *international law* or *retail*.)

WebBootCaT requires only a list of seed words (terms that are expected to be typical of the domain of interest) as an input. For example, to generate a *new technology* corpus, you might use seed words such as *phone, wi-fi, email, wireless, Internet*, etc. BootCaT then generates a corpus based on searches for these seed words.

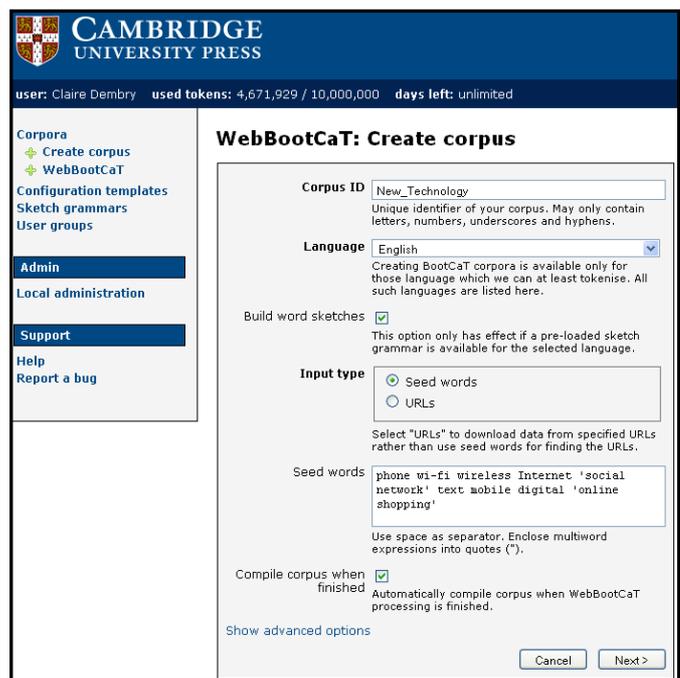To build your own corpus, click on WebBootCaT (shown in red below) from the Sketch homepage:



This takes you to a screen where you can enter the name of your corpus, the language of your corpus (English) and your seed words.

For example, if we wanted to make a *new technology* corpus, our screen might look like this:
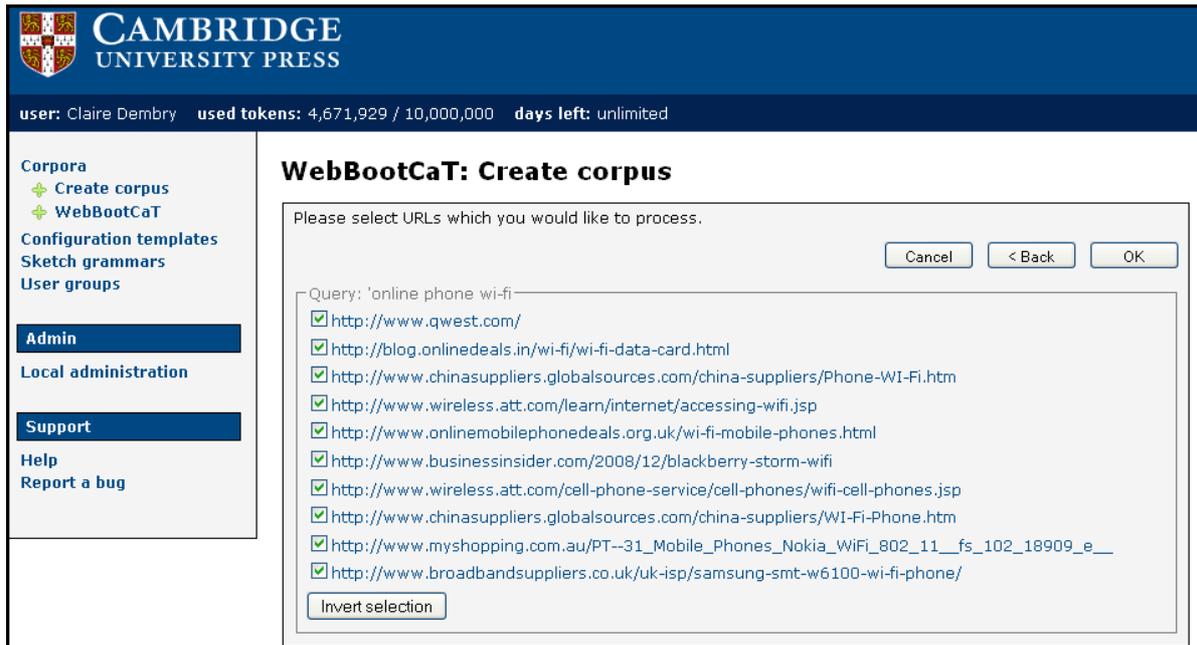
Seed words should be separated from each other by a space. Seed words that consist of more than one word (such as *social network* in the example on the right) should be enclosed in quotation marks.

The more seed words you enter, the larger your corpus will be.

Once you have entered all of your seed words click on *Next* at the bottom of the page.
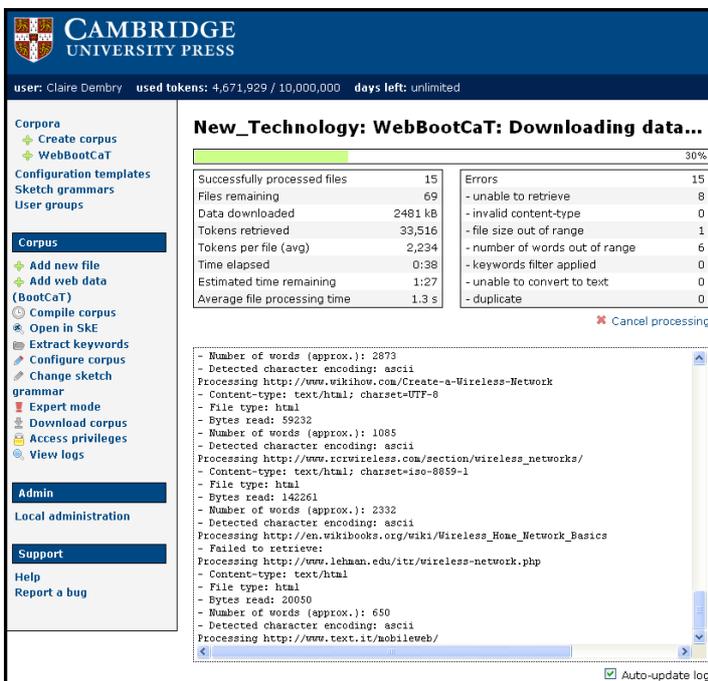


This takes you to a screen that shows the websites that BootCaT has found that relate to your seed words. Here you can browse this list of URLs and untick any sources that you do not wish to use - all of the boxes are ticked by default.

Once you've selected the URLs you wish to use, click *OK* (at the top, and also at the bottom of the page) to download data from these sources.

Sketch Engine will now process the text so that it can be used with the same functionality as the Cambridge Corpus resources - the toolbar at the top of the screen keeps you updated on BootCaT's progress.



This can take a few moments – particularly if you've entered lots of seed words.

You can navigate away from this screen (e.g. by clicking on the Cambridge logo at the top left of the page) and leave BootCaT running in the background while you complete other corpus searches.

You can return to check on the progress of BootCaT at any time by clicking on the name of your corpus from the opening menu page (navigate there by clicking on the Cambridge logo.)

Once the data has downloaded, the following screen is shown:

In order to use all of the functions available in Sketch with your new corpus, (e.g. Word Sketch, Sketch Diff) you now need to <u>compile</u> your corpus. To do this, click on *Compile corpus*, and then on *Compile*. This adds the information needed to use all of the functions.

Once the processing is complete, click on *OK*. Your corpus is now ready to use in Sketch Engine – to do this click on *Open in SkE*.



Your corpus will now also appear on the opening page, underneath the other Cambridge corpora that are available to use.

Sketch users can create their own corpora up to a maximum of 1m words (it may be possible to increase this amount, on request – contact cdembry@cambridge.org for more info).

It is also possible to **share** corpora that you have built using WebBootCaT with other Sketch Engine users:
- click on the name of your new corpus from the main screen (shown above)
- click on *access privileges*
- select the level of access you would like the other person to have
- type their name into the box then select it from the list that appears and click *ok.*
- this user can now access the corpus you have built.

**A note of caution!** WebBootCaT gathers data from the web in an indiscriminate way. Because of this, occasionally unusual words or characters may appear in your concordance lines - you should bear this in mind when drawing conclusions from any results you find.

## 2. View options – advanced

*View options* allows you to change the type and amount of information that is displayed in your concordance lines. For a discussion of changing the references column and page size, please see Sketch Engine – Getting Started, available here: http://www.cambridge.org/sketch/help/

From time to time you may want to display additional information in your concordance lines.

For example, you may wish to see the Part of Speech tag that a particular word has been assigned.

Under *Attributes* a selection of information that can be displayed is shown. You can choose to show as many of these as you wish (although adding in more than one at a time can make your results difficult to read!)



Once selected, the attributes appear in the concordance separated from the words by a forward slash, as shown in the examples below:

- <u>Word</u> displays the words in the concordance.

- <u>Tag</u> displays the Part of Speech that has been assigned to each word.
  For example: **Girls**/*NNS* **are**/*VBP*

- <u>Lempos</u> displays both the lemma and the Part of Speech that a word belongs to.
  For example: **Girls/***girl-n*  **are**/*be-v*

- <u>Lemma</u> displays the lemma (or headword) that each word belongs to.
  For example: **Girls/***girl* **are**/*be*

- <u>Lc</u> indicates the lowercase forms.
  For example: **Girls/***girls* **are/***are*

- <u>Lemma_lc</u> displays the lowercase form of the lemma that each word belongs to.
  For example: **Girls/***girl* **are/***be*

A full list of the PoS tags can be found in the Cambridge Help: http://www.cambridge.org/sketch/help/

You can choose to display attributes for <u>each token</u> in your results (by selecting *display attributes for each token*) or for your <u>search term only</u> (by selecting *display attributes KWIC tokens only*).

Tags for the KWIC tokens only in a simple search for *chase* will return results like this:



Last updated: February 2012

It is also possible to display markers that indicate boundaries of different sorts using the *Structures* column. For example, sentence boundaries, <s>, can be seen in green an excerpt below:



## 3. Simple vs. multilevel sort

Sorting allows you to change the order in which your concordance is displayed. You might want to use *sort* to look at patterns within your results.

To sort the words in the concordance alphabetically to the **left**, **right**, on the **node,** or by **references** choose these option from the concordance menu.

To run more complex sorts, click on the *sort* button itself, situated on the lower left menu. This takes you to the screen shown below:



There are two types of searches that you can run:

- **Simple Sort** allows you to sort your concordance results by a feature that you select (e.g. by word, tag, genre or variety etc).

- **Multilevel Sort** allows you to sort in the same way, but to sort for more than one parameter, e.g. you can sort by word then by PoS tag within the same results.

To run a **Simple Sort,** select options from the top section of the sort screen. Use the dropdown box next to *Attribute* to select the feature you would like to sort by. For example, you can sort by *word, tag, source, or genre*.

You can then select whether you'd like to sort the concordance to the *left, right* or on the *node* (i.e. your initial search term). Choose the option you'd like to select from the buttons next to *Sort key*

You can select how many tokens are sorted e.g. choose 3, if you wanted to sort 3 words to the left. You can select to ignore case and you can also select to sort backwards (reverse the sort)
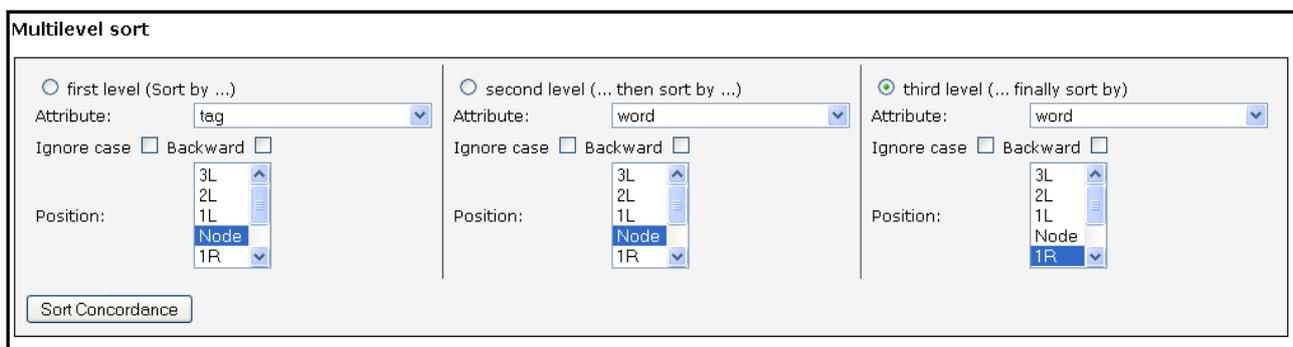
Note! If you choose to sort by e.g. tag, in order to see the actual tags in the results you must choose this option from the view options menu (look back to Section 2. for more information on this).

To run a **Multilevel Sort** you need to fill in the boxes in the bottom half of the sort screen. The choices are the same as the simple sort, but with a multilevel sort you can sort on 1, 2 or 3 levels – this means that you can sort within a sort.

For example, you can run a simple search for *fire* (i.e. all inflected forms, all PoS). Then, by using Multilevel Sort you can then sort on 3 levels. For example you could sort by:

- First level - part-of-speech of the node (so, sorting *fire* as a noun from *fire* as a verb).
- Second level - the word (i.e. the actual word form of the node, so sorting the verbs *fire* from *fired*, etc).
- Third level – word to the left of the node (1L).
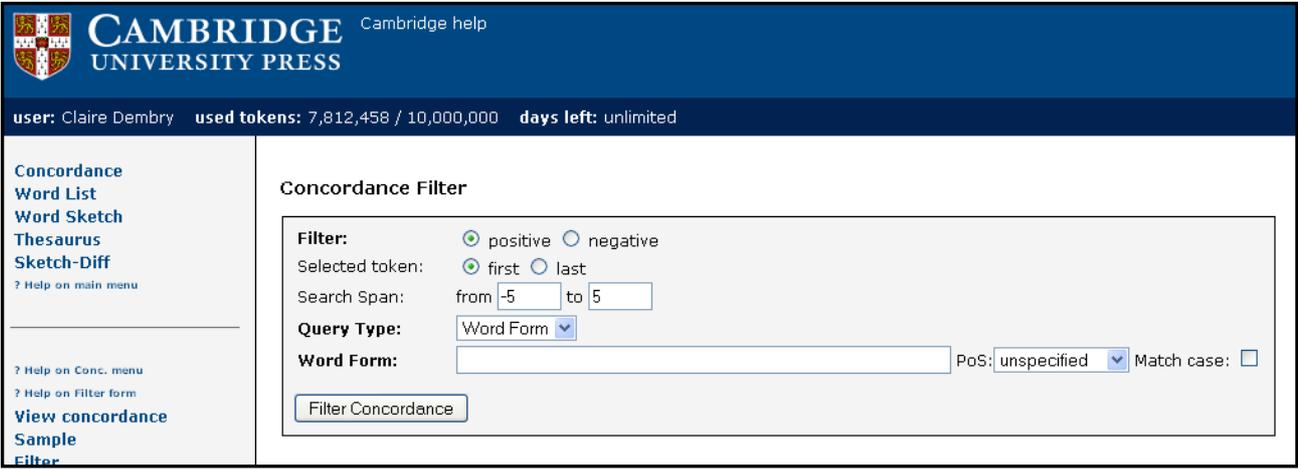
These options would look like this:



If you're sorting by more than one level, you should click the button at the highest level, i.e. to search by first and second level, click the second level button.

Note! If you choose to sort by e.g. *tag*, or *lempos* you must choose this option from the view options menu in order to see the actual tags or lemma-pos displayed in the results.

## 4. Filter

The **filter** function allows you to narrow down your search to include or exclude a particular word or phrase in your analysis.

To run a filter, first make a concordance of your search term. From the lower left hand side menu, click on the *filter* button. This takes you to the screen shown below:



In the main panel you have various options for filtering your concordance.

You have the following choices:

- **Positive:** includes a word or phrase in the subsequent results.
- **Negative:** excludes a word or phrase from the subsequent results.
- **First** and **Last:** use this option if your original query was for more than one word - you can determine which token (first or last) is used when calculating distance using the Search Span.
- **Search Span:** set how many tokens before or after your node word that the filter should apply to.

The *query type* box works in the same way as a 'normal' concordance query – you are able to run simple, lemma, word form, phrase and CQL searches. Here you should enter the word or phrase that you would like to include or exclude in your filter. You can also match the PoS and case of the word you enter.

## 5. Frequency

*Frequency* allows you to investigate the most common features of your search (e.g. most frequent words, tags, genres etc).

To look at frequency results for the **node tags, node forms, doc IDs** and **Text Types,** choose these options from the concordance menu.

To obtain more complex frequency results, click on the *frequency* button itself, situated on the lower left menu. This takes you to the screen shown below:

There are two types of frequency distributions that you can run:

- **Text Type frequency distribution** is straightforward to use – it allows you to show how your search term is distributed through the texts in the Corpus. You may find, for example, that a word like *police* appears significantly more often in newspaper texts than in other text types. To display results for more than one category, hold ctrl and clicking.

- **Multilevel frequency distribution** allows you to look at the frequency of more than one parameter, e.g. you can sort by the most frequent word then by the most frequent PoS tag within the same results.
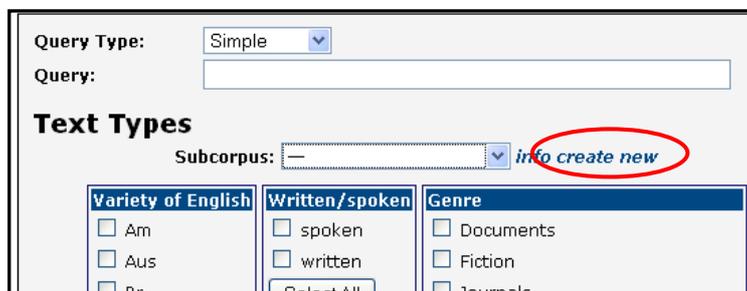
This multilevel approach works on the same basis as multilevel sorting (described in Section 3. earlier). We can explore multilevel frequency distribution further by looking at an example:

- Run a simple search for *travel*.
- From the concordance select *Frequency* from the left hand menu
- In the Multilevel frequency distribution section, choose *word* and select *node.* This will display the frequencies of e.g. *travels, travelling, travelled*.
- To search for the most frequent words <u>following</u> these forms, we need to add a second level:
- Under *second level* set the attribute at *word* and set the position at *1R* (word one position to right of node word)
- To search for the most frequent words <u>following</u> both the node and the token after the node, we need to add a third level:
- Under *third level* set the attribute at *word* and set the position at *2R*.
- Remember that if you're sorting by more than one level, you should click the button at the highest level, i.e. to search by first <u>and </u>second level, click the second level button.
- Press *Make frequency list* to show your results.

## 6. Building a subcorpus

You can use the *subcorpus* function to explore only specific parts of the corpus. To create a subcorpus:

- Click on *text types* to display the text types options

- Next to Subcorpus, click on *create new* (shown in red)

- (Subcorpora you have already built appear in the dropdown menu)

- In the next screen, choose a name for your subcorpus and enter it into the box.

- Select the text types you wish to include in your subcorpus by ticking the appropriate boxes.

- Click on *create subcorpus*

Your subcorpus is now added to your list. There is no limit to the number of subcorpora that you can build.

## 7. Working with subcorpora

To run a **Concordance** using only your subcorpus, select *text types* and choose your subcorpus from the dropdown menu. Then complete your concordance search in the usual way – the data returned will be from your subcorpus only.

To generate a **Frequency list** using only your subcorpus, click on *frequency list* and then choose your subcorpus from the dropdown list. The most frequent words in your subcorpus will then be listed.

To run a **Word sketch** using only your subcorpus, click on *Word sketch* then on *advanced options* from the lower left hand menu. Then select your subcorpus from the dropdown list and then run the Word sketch as usual.

You may also want to use your subcorpus to run a **Keyword analysis**. This allows you to compare two corpora or subcorpora in order to find out the words that are comparatively the most different. For example, a keyword analysis that compared a 'business' subcorpus to, for example, the Cambridge International Corpus would indicate the words that are most frequently found in the business subcorpora but not in the CIC (and are therefore more typical in that domain).

To run a **Keyword analysis** using your subcorpus, click on *frequency list.* Under *Keywords* select your subcorpus from the first drop down menu and choose a corpus to compare it to from the second drop down menu.

## 8. Building wordlists

It is possible to build word frequency list using Sketch Engine. You can do this by corpus or by subcorpus.
To run a frequency list, click 'word list' from the left hand side menu. The word list screen is shown below:



If you would like to make a word list from a subcorpus that you have already defined, select it from the drop down menu, or click 'create new' to create a new subcorpus to work with.

You can search for all words, or you can enter patterns, e.g. all words ending in "ing" or starting with "pre". These patterns should be entered using the format: *.\*ing* or *pre.\**

You can also change the minimum frequency level of your list, in order to exclude low frequency items.

Whitelists and Blacklists allow you to include and exclude words from your wordlist.
- White lists search only for those words in the list, and ranks them with respect to how frequent they are in the corpus.
- Black lists exclude words from the corpus frequency list. Blacklists can be used to exclude stopwords (for a list of commonly used stopwords that you may wish to use, click here: http://cup.sketchengine.co.uk/stopwords/english/)

Once you're happy with your selections, click on make wordlist. You wordlist will then be displayed. You can then choose 'save' to e.g. save it to Excel.

## 9. Comparing corpora using keywords

You can create a keyword list in order to compare one corpus to another. A keyword list displays those words in your corpus that are most different to the other corpus you compare it to (the reference corpus).

To create a keyword list:
- click on 'wordlist' from the left hand side menu options. You will then see the screen shown in 8 above.
- Choose your corpus or subcorpus, and then under 'output type' (in the lower half of the screen) select 'keywords'
- Select your reference corpus (and, if you wish, subcorpus). In order to get a picture of how your corpus is different to others, if may be useful to compare it to a very general corpus – in this case, the CIC is a good choice.
- Click *make Word List* to display your results

## 10. Word Sketch – Advanced options

Along with using Word Sketch in the usual way (as outlined in the Getting Started guide), a number of 'expert options' are available.



Open these options by selecting 'word sketch' then clicking 'expert options' from the lower left hand menu:

Here, the main difference is that you can search by a particular subcorpus, by selecting it from the dropdown.

You can also choose to see the grammatical relations you wish, by selecting/deselecting them from the lower part of the screen.

## 11. Uploading your own text files to Sketch Engine

Most typically, Sketch Engine is used in order to examine Cambridge corpus resources. However, it is possible to upload your own text files and analyse them in the same way.

Your own data will be POS tagged thus allowing you to use many of the same functions usually can in Sketch (such as e.g. Work Sketch). However, you are not able to use any metadata associated with you text files.

**Creating a Corpus:**
Before adding any text files, you need to create a new location for them in Sketch Engine.

Go to the main Sketch Engine homepage, and select Create Corpus (underneath the My Corpora heading). The following screen will appear:

Type in a name for your corpus in the Corpus ID and Corpus name boxes. The Corpus name will appear on your main Sketch Engine screen.

The Info box is optional, but it may be helpful to add more detail about your data for yourself and/or other users. Select the language of your corpus (typically English) from the drop-down box, and click Next.

On the next screen, choose TreeTagger for English and click Next.

On the next screen, choose English PennTB-TreeTagger MCD 2.1 and click Finish.

**Uploading text files:**
You can now upload text files to your new corpus. Choose Add new file (underneath the corpus name, and also in the left hand menu). The following screen will appear:

**Test2: Add new file: Step 1**

Supported file types: .txt, .html, .pdf, .doc, .vert (vertical).

| | | |
|---|---|---|
| Upload from disk ○ | [ ] | Browse... |
| Download from location ○ | http:// | |
| Use file or directory on the server ○ | _____ ▼ | |

☐ show files in subdirectories

FTP to cup.sketchengine.co.uk at port 10021 to upload files. Use the same user name and password as for logging into this web interface.

Paste text ○

Cancel    Next >

This is where you can select the files you want to upload. If you want to upload text files saved in a folder on your computer, select Upload from disk and click on Browse... to select the file.

You can also paste text directly into the box shown above, by selecting the Paste text button.

Once you have selected your option and selected/pasted your file, click Next, and when the file preview appears click Finish. This will take you back to the original screen for your corpus, where you can select Add new file again and repeat the process.

Each file has to be uploaded separately.

**Compiling your corpus:**

Once you have added all of your files, click on Compile corpus. You are unable to work with your corpus until it has been compiled.



On the next screen, make sure English PennTB-TreeTagger MCD 2.1 is selected (it is the default). Then click on Compile at the bottom of the screen.

When the message at the top of the screen changes from Corpus is busy to Processing done!, click on OK.



Your corpus is now ready to use. You can open it from the Sketch Engine homepage under My corpora; click on the small magnifying glass logo next to the corpus name:

## Sharing your corpus with other Sketch Engine users:

It is possible to share any corpus you upload with other Sketch Engine users.

To do this, from the main screen click on the name of the corpus you have uploaded.

The files that you have uploaded will then appear on the screen, along with a number of options displayed in the left hand menu, as shown below:

To share your corpus, click on Access privileges, as shown above.

There are three different levels of access that you can assign to users. These access levels are explained on the screen.

To add a user start to type their name in the appropriate box and then select them from the dropdown list. Then click ok. The user(s) you have selected will now have access to your corpus. Your corpus will appear underneath any corpus that they have added themselves, as shown below: