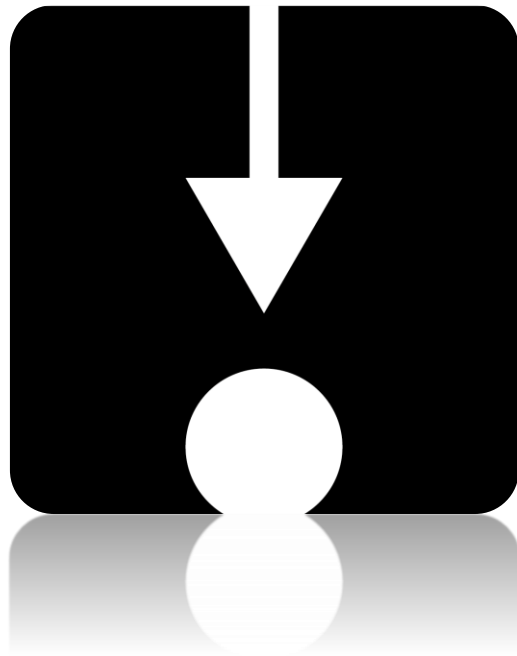


# #LancsBox 6.0 manual



Citation for #LancsBox:

Brezina, V., Weill-Tessier, P., & McEnery, A. (2021). #LancsBox v. 6.x. [software package]

Brezina, V., Weill-Tessier, P., & McEnery, A. (2020). #LancsBox v. 5.x. [software package]

Brezina, V., Timperley, M., & McEnery, A. (2018). #LancsBox v. 4.x. [software package]

Brezina, V., McEnery, T. & Wattam, S. (2015). Collocations in context: A new perspective on collocation networks. *International Journal of Corpus Linguistics*, 20(2), 139-173

.innovation in corpus linguistics

**#LancsBox**

@Lancaster University

## Contents

|      |  |    |
|------|--|----|
| 1    | Downloading and running #LancsBox version 6.0..... | 5  |
| 2    | Loading and importing data .....                   | 9  |
| 2.1  | Visual summary of Corpora tab.....                 | 9  |
| 2.2  | Load your corpora and wordlists.....               | 9  |
| 2.3  | Supported file formats .....                       | 10 |
| 2.4  | Download #LancsBox corpora and wordlists .....     | 10 |
| 2.5  | Create a new corpus from the web                   | 10 |
| 2.6  | Working with corpora and wordlists                 | 12 |
| 2.7  | Saving corpora.....                                | 12 |
| 2.8  | Pre-processing of corpora (Advanced users).        | 12 |
| 3    | Key functionalities .....                          | 15 |
| 3.1  | Mouse clicks .....                                 | 15 |
| 3.2  | Shortcut Keys.....                                 | 15 |
| 3.3  | Tools and Tabs.....                                | 16 |
| 3.4  | Split screen .....                                 | 16 |
| 3.5  | Saving results.....                                | 17 |
| 3.6  | Copy/pasting selected results .....                | 17 |
| 4    | KWIC tool (key word in context) .....              | 18 |
| 4.1  | Visual summary of KWIC tab .....                   | 18 |
| 4.2  | Searching and displaying results.....              | 19 |
| 4.3  | Settings and full text pop-up .....                | 19 |
| 4.4  | Sorting, randomising and filtering..               | 20 |
| 4.5  | Statistical analysis.....                          | 20 |
| 5    | Whelk tool .....                                   | 22 |
| 5.1  | Visual summary of Whelk tab .....                  | 22 |
| 5.2  | Top panel: KWIC.....                               | 22 |
| 5.3  | Bottom panel: Frequency distribution               | 22 |
| 5.4  | Statistical analysis.....                          | 23 |
| 6    | GraphColl .....                                    | 24 |
| 6.1  | Visual summary of GraphColl tab ...                | 24 |
| 6.2  | Producing a collocation graph .....                | 24 |
| 6.3  | Reading collocation table.....                     | 25 |
| 6.4  | Reading collocation graph .....                    | 26 |
| 6.5  | Extending graph to a collocation network.....      | 27 |
| 6.6  | Shared collocates.....                             | 28 |
| 6.7  | Problems with graphs: overpopulated graphs         | 29 |
| 6.8  | Reporting collocates: CPN.....                     | 30 |
| 7    | Words tool .....                                   | 31 |
| 7.1  | Visual summary.....                                | 31 |
| 7.2  | Producing frequency list .....                     | 32 |
| 7.3  | Visualizing frequency and dispersion               | 32 |
| 7.4  | Producing keywords.....                            | 33 |
| 7.5  | Producing corpus statistics .....                  | 33 |
| 8    | Ngram tool.....                                    | 35 |
| 8.1  | Visual summary.....                                | 35 |
| 9    | Text .....   | 37 |
| 9.1  | Visual summary.....                                | 37 |
| 9.2  | Searching in Text.....                             | 37 |
| 9.3  | Settings .....                                     | 38 |
| 10   | Wizard .....                                       | 39 |
| 10.1 | Visual summary.....                                | 39 |
| 10.2 | Selecting settings and running Wizard              | 40 |
| 10.3 | Data analysis .....                                | 40 |
| 10.4 | Research report .....                              | 41 |
| 11   | Searching in #LancsBox.....                        | 42 |
| 12   | Statistics in #LancsBox .....                      | 48 |
| 12.1 | Frequency measures.....                            | 48 |
| 12.2 | Dispersion measures.....                           | 48 |

|      |                            |    |
|------|----------------------------|----|
| 12.3 | Keyword measures .....     | 48 |
| 12.4 | Collocation measures ..... | 49 |
| 13   | Glossary .....             | 50 |
| 14   | Messages.Properties .....  | 53 |

## #LancsBox v.6.0: License

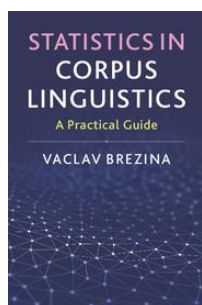
#LancsBox is licensed under BY-NC-ND Creative commons license. #LancsBox is free for non-commercial use. The full license is available from: <http://creativecommons.org/licenses/by-nc-nd/4.0/legalcode>

#LancsBox uses the following third-party tools and libraries: Apache Tika, Gluegen, Groovy, JOGL, minlog, QuestDB, RSyntaxTextArea, smallseg, TreeTagger. Full credits are available <http://corpora.lancs.ac.uk/lancsbox/credits.php>

When you report research carried out using #LancsBox, please cite the following:

- Brezina, V., McEnery, T. & Wattam, S. (2015). [Collocations in context: A new perspective on collocation networks](#). *International Journal of Corpus Linguistics*, 20(2), 139-173.
- Brezina, V., Weill-Tessier, P., & McEnery, A. (2021). #LancsBox v. 6.x. [software package]
- Brezina, V., Timperley, M., & McEnery, A. (2018). #LancsBox v. 4.x. [software package].

## Statistical help



Brezina, V. (2018). *Statistics for corpus linguistics: A practical guide*. Cambridge: Cambridge University Press.

If you are interested in finding out details about statistical procedures used in corpus linguistics, refer to Brezina (2018); visit also Lancaster Stats Tools online at <http://corpora.lancs.ac.uk/stats>

## Further reading and materials

Brezina, V. (2016). Collocation Networks. In Baker, P. & Egbert, J. (eds.) *Triangulating Methodological Approaches in Corpus Linguistic Research*. Routledge: London.

Brezina, V. (2018). Statistical choices in corpus-based discourse analysis. In Taylor, Ch. & Marchi, A. (eds.) *Corpus approaches to discourse: a critical review*. Routledge: London.

Brezina, V. & Gablasova, D. (2017). The corpus method. In: Culpeper, J, Kerswill, P., Wodak, R., McEnery, T. & Katamba, F. (eds). *English Language (2nd edition)*. Palgrave.

Brezina, V., McEnery, T. & Wattam, S. (2015). Collocations in context: A new perspective on collocation networks. *International Journal of Corpus Linguistics*, 20(2), 139-173.

Brezina, V., & Meyerhoff, M. (2014). Significant or random. *A critical review of sociolinguistic generalisations based on large corpora*. *International Journal of Corpus Linguistics*, 19(1), 1-28.

Gablasova, D., Brezina, V., & McEnery, T. (2017). Collocations in corpus-based language learning research: Identifying, comparing, and interpreting the evidence. *Language Learning*, 67 (S1), 155–179.

Gablasova, D., Brezina, V., & McEnery, T. (2017). Exploring learner language through corpora: comparing and interpreting corpus frequency information. *Language Learning*, 67 (S1), 130–154.

- More materials (video lectures, exercises, slides etc.) are available: on the #LancsBox website: <http://corpora.lancs.ac.uk/lancsbox/materials.php>

# 1 Downloading and running #LancsBox version 6.0

#LancsBox is a new-generation corpus analysis tool. Version 6 has been designed primarily for 64-bit operating systems (Windows 64-bit, Mac and Linux) that allow the tool’s best performance. #LancsBox also operates on older 32-bit systems, but its performance is somewhat limited. Version 6 of #LancsBox comes with an installer, which makes installation of #LancsBox even easier.

---

① **Select and download:** Select the version suitable for your operating system and download installer to your computer.



② **Run installer**

Agree to security warnings on your machine – #LancsBox is safe to run – and follow the steps in the installer. Always install #LancsBox to a folder, where the tool has ‘read and write’ privileges such as the User folder or Desktop; On Windows, never install #LancsBox to Program Files.

---

## Important note: System privileges

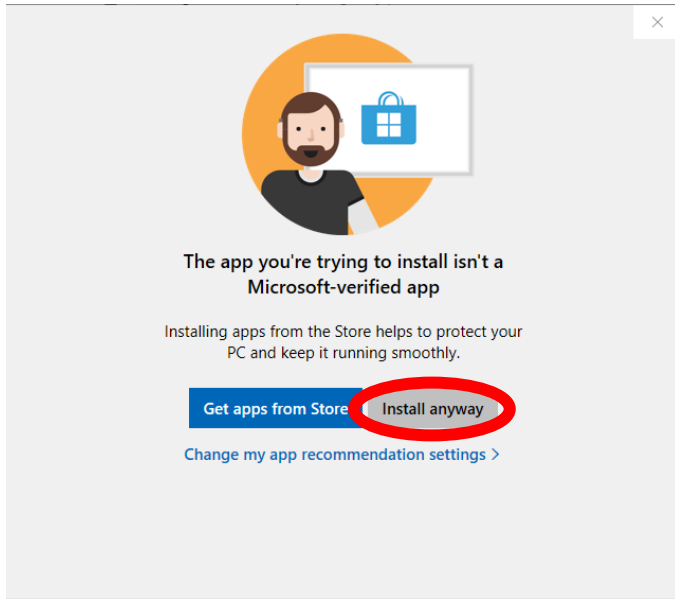
Please follow the instruction below for your specific operating system.

### Windows 10

Windows 10 will display either of these two messages.

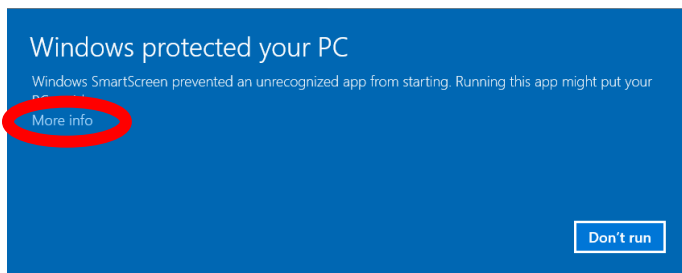
>Newer builds

“The app you are trying to install isn’t a Microsoft-verified app.”. If this warning message appears, click on ‘Install anyway’.

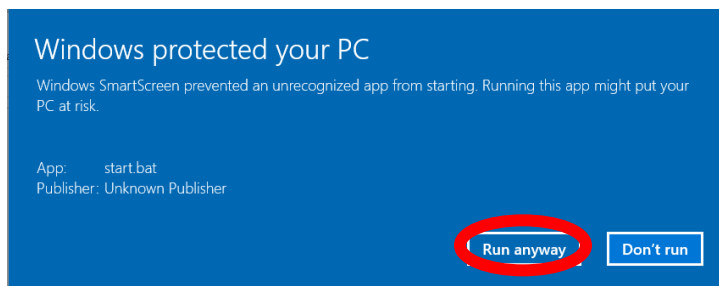


>Older builds

“Windows protected your PC”. If this warning message appears, click on ‘More info’

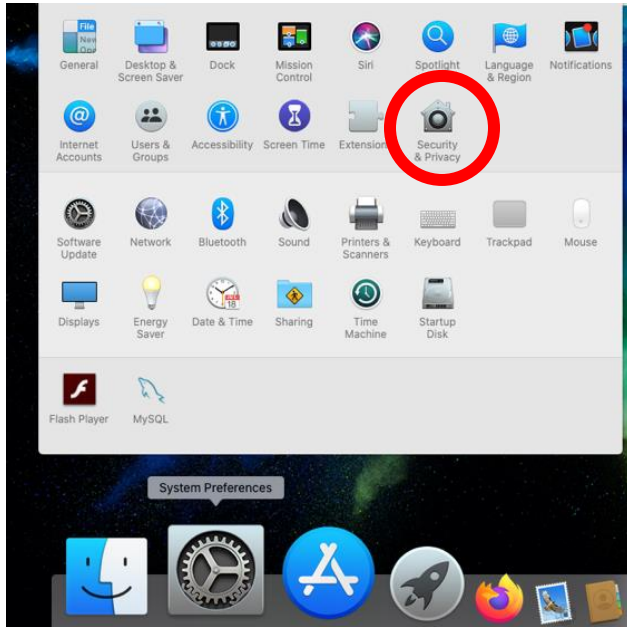


Then click on ‘Run anyway’.

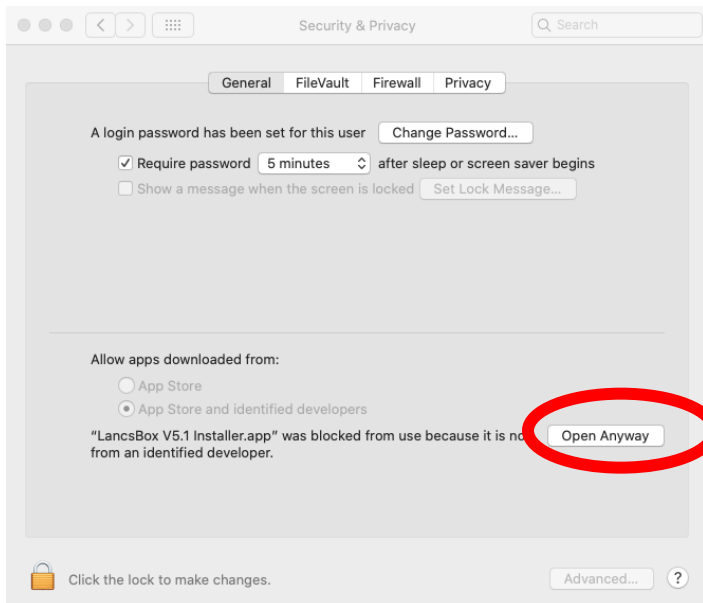


## MAC

Open “System Preferences” in the dock, click on “Security & Privacy”.

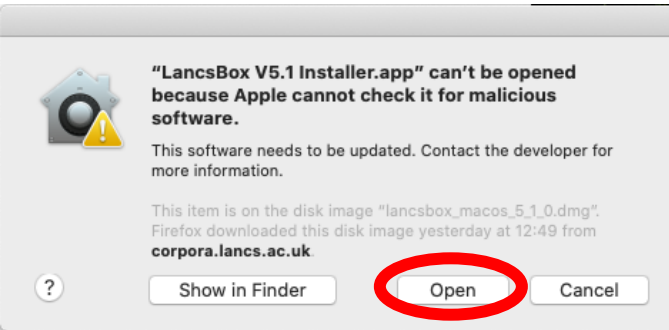


Click on “Open Anyway” next to the message “LancsBox V6.0 Installer was blocked because it is not from an identified developer”.



Click on “open” when the message “LancsBox V6.0 Installer.app” can’t be opened because Apple cannot check it for malicious software” is displayed in a new window.

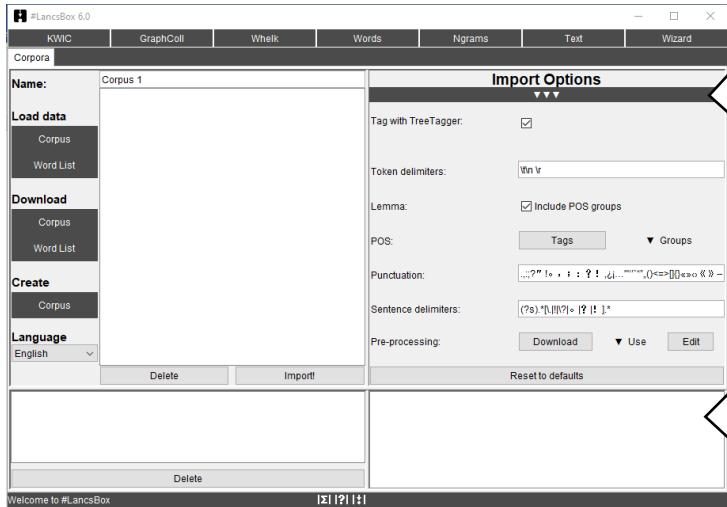




## 2 Loading and importing data

Data can be loaded and imported into #LancsBox on the 'Corpora' tab. This tab opens automatically when you run #LancsBox. #LancsBox works with corpora in different formats (.txt, .xml, .doc, .docx, .pdf, .odt, .xls, .xlsx, .zip etc.) and with wordlists (.csv). There are three options for loading corpora and wordlists: i) load (your own) data, ii) download corpora and wordlists that are distributed with #LancsBox and iii) create your own corpora from the web.

### 2.1 Visual summary of Corpora tab



**Top panel: Importing corpora and wordlists**

**You can:**

- Select your corpus or wordlist to load.
- Download a corpora and wordlists distributed with #LancsBox.
- Create a corpus
- Select language.
- Review POS tags.
- Review punctuation marks and sentence delimiters.

**Bottom panel: Working with corpora and**

**You can:**

- Activate or delete imported corpora or wordlists.
- Review corpus and text size (tokens, types, lemmas).
- Preview texts.
- Save processed corpora with pos-tags etc.

### 2.2 Load your corpora and wordlists

#LancsBox allows you to work easily with your own corpora and wordlists. These corpora are those stored on your computer or at a location accessible from your computer (memory stick, shared drive, dropbox, cloud etc.).

1. In the Corpora tab, left-click on 'Corpus' or 'Word List' under 'Load data', depending on whether you want to load a corpus or a wordlist.
2. This will open a window where you can navigate to the location (folder) where your corpus or wordlist is stored.
3. You can select a specific file, select multiple files by holding down Ctrl and left-clicking on your chosen files, or select all files in the folder by holding down Ctrl + A.
4. Left-click 'Open' to load your files.
5. Select the language of your corpus or wordlist. #LancsBox supports automatic lemmatisation and POS tagging in multiple languages. This is done using Tree Tagger. If your language is not listed, select 'Other'; in this case, automatic lemmatisation and POS tagging will be disabled.
6. [Optional: You can review/change the import options by left-clicking on a bar with three triangles (▲▲▲). In most cases, you can use the default options.]

7. Left-click 'Import!' to import your corpus into #LancsBox. By default, #LancsBox automatically adds POS tags to the corpus.

### 2.3 Supported file formats

---

#LancsBox supports different file formats (.txt, .xml, .doc, .docx, .pdf, .odt, .xls, .xlsx, .zip and many others) of corpus files. #LancsBox automatically extracts and processes text available in corpus files. For wordlists, #LancsBox assumes the comma-delimited file format (.csv).

---

1. Corpus formats: .txt, .xml, .doc, .docx, .pdf, .odt, .xls, .xlsx, .zip – full list: [Apache Tika](#).
2. Wordlist format: csv (see example below).

```
Corpus: BNC| Language: English| 4055 files| 96996843 tokens| 662414 types| 716618 lemmas|
"Type", "Frequency: 01 - Freq", "Dispersion: 01_CV"
"the", "6054524.000000", "0.286889"
"of", "3049295.000000", "0.400166"
"and", "2622080.000000", "0.263099"
"to", "2599355.000000", "0.223254"
"a", "2168976.000000", "0.221813"
"in", "1945319.000000", "0.333547"
```

### 2.4 Download #LancsBox corpora and wordlists

---

#LancsBox allows you to work with existing corpora that are freely distributed with #LancsBox under a specific license. Two modes for corpus sharing are available: i) open access and ii) restricted access. We are constantly adding more corpora to this list.

---

1. In the corpora tab, left-click on 'Corpus' or 'Word List' under 'Download'.
2. This will open a window where you can select corpora or wordlists distributed with #LancsBox. By left-clicking on a corpus, you will be shown additional information about the corpus or wordlist, including the language, date, text type, license etc.
3. Review and agree with the corpus license.
4. Left-click 'Download' to download the selected corpus or wordlist.
5. Left-click 'Import!' to import your corpus into #LancsBox. By default, #LancsBox automatically adds POS tags to the corpus.

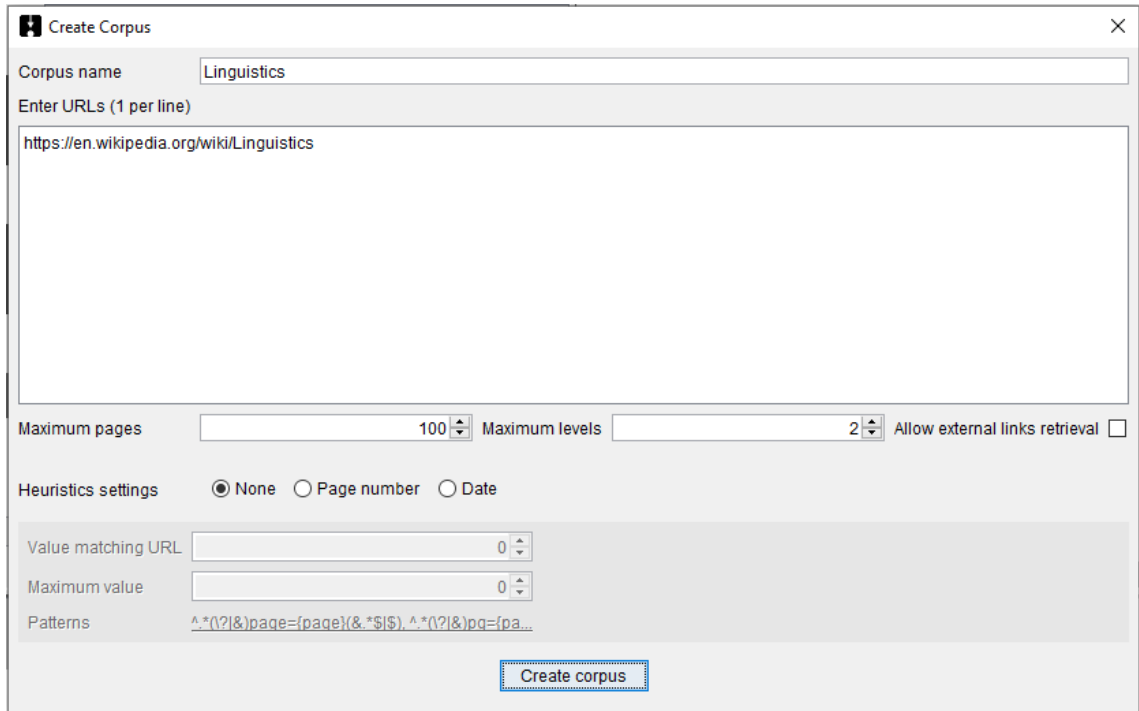
### 2.5 Create a new corpus from the web

---

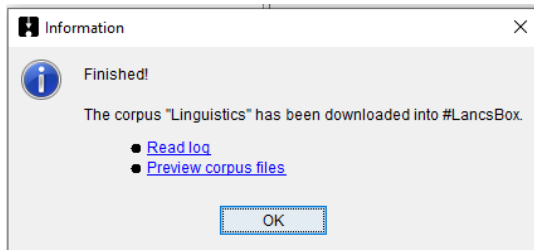
#LancsBox allows you to create a new corpus based on URLs provided. It downloads multiple webpages and extracts text from them.

---

1. In the corpora tab, left-click on 'Corpus' under 'Create'.
2. This will open a window where you can enter the corpus name, URLs and specify select extraction criteria.



3. By default, #LancsBox extracts 100 webpages at two levels of embedding.
4. Left-click on 'Create corpus' to start the extraction process.
5. When the process is finished, #LancsBox will notify you. You can then open the log with details about the websites identified and extracted and view the folder with the data (txt).



6. Click 'OK' and then 'Import' to load data into #LancsBox directly.

Note: #LancsBox allows heuristic searches for websites that include multiple pages or dates e.g. [https://www.mumsnet.com/Talk/am\\_i\\_being\\_unreasonable/4170918-to-even-consider-giving-husband-second-chance?pg=1](https://www.mumsnet.com/Talk/am_i_being_unreasonable/4170918-to-even-consider-giving-husband-second-chance?pg=1)

Switch on the heuristics setting and specify i) heuristic type (page number or date) ii) the value matching the URL (1 in the example above) and iii) the maximum value. -1 Indicates unlimited number of pages until all pages are downloaded, while a date in the future indicates an unlimited number of dates until all dates are exhausted.

## 2.6 Working with corpora and wordlists

---

All corpora and wordlists that have been imported into #LancsBox are displayed in the bottom panel on the 'Corpora' tab. This panel allows reviewing corpora, previewing files and fast reloading of corpora and wordlists when #LancsBox is closed and re-opened.

---

1. If you have imported a corpus (📁) or wordlist (📄) it will appear in the bottom panel, alongside any other corpora or wordlist you have already imported. These can be removed by left-clicking 'delete'. In the bottom-right section, you can view the corpus structure: the individual text files that the corpus is composed of.
2. In the bottom panel (bottom left window), the default corpus can also be specified. The default corpus is a corpus that #LancsBox offers as a default choice in the individual modules. The default corpus can be specified by left-double-clicking on the name of the corpus; a filled rectangle (■) will appear next to the name of the default corpus.
3. If #LancsBox is closed, the corpora and wordlists will remain imported but will be unloaded. To activate (reload) the corpora or wordlists for use, left-double-click on the corpora or wordlists.
4. You can also preview the files by right-clicking on them. They will appear in the Text tool (see Section 8). The list of files (including the info about their size) can also be copied (Ctrl/Command+C) and pasted (Ctrl/Command+V) into a spreadsheet or text document.
5. Corpora are now ready to be analysed using five modules: KWIC, Whelk, GraphColl, Words and Text. Wordlists can be used in the Words tool.

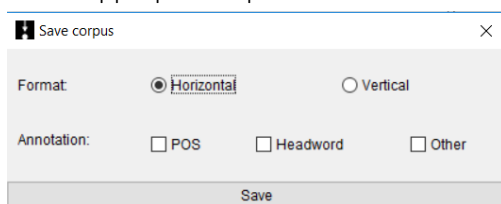
## 2.7 Saving corpora

---

#LancsBox saves corpora in the horizontal or the vertical format.

---

1. Right-click on the corpus which you wish to save.
2. Select appropriate options.



3. Click 'Save'.

## 2.8 Pre-processing of corpora (Advanced users).

---

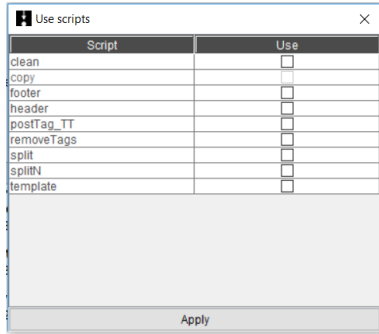
#LancsBox allows pre-processing data as part of the import procedure. This is set up in the 'Import options' under 'Pre-processing'. Data can be modified in different ways using a variety of Groovy scripts, which are fully customisable.

---

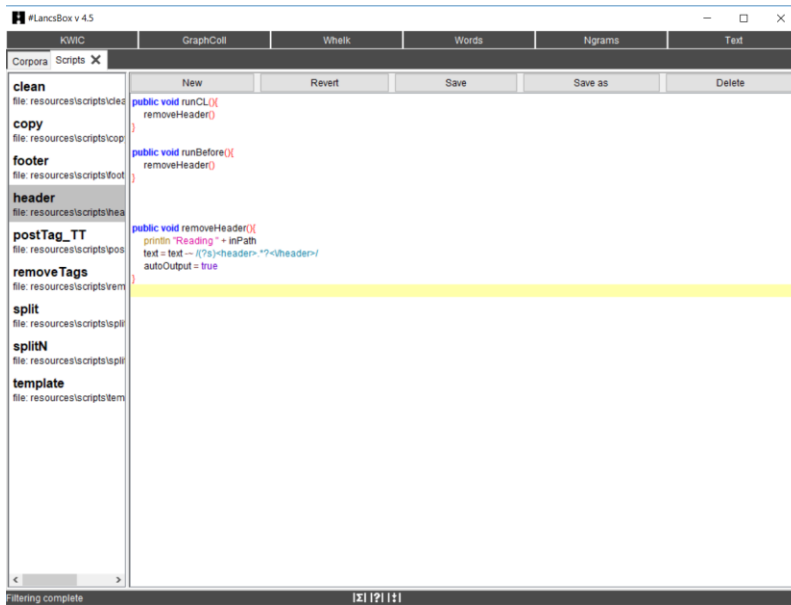
1. Under 'Pre-processing' three options are available:



2. 'Download' allows the user to download scripts and their newest versions available from the #LancsBox website.
3. 'Use' displays a list of currently available scripts and a checkbox next to each script for the user to indicate, which scripts will be used in the pre-processing stage.



4. 'Edit' displays the scripts in a script editor, which allows modifying existing scripts and creating new scripts.



5. The structure of a script is as follows. More information about the Groovy scripting language can be found at <http://groovy-lang.org>.

| Script   | Comments  |
|--|---|
| <pre>public void runCL(){     println "Ran on the command line." }</pre>         | Scripts run via the command line.   |
| <pre>public void runBefore(){     println "Ran as a pre-process script." }</pre> | Scripts run when the files are being loaded. This allows splitting files, deleting or changing texts or structuring elements e.g. xml tags. |

```
public void runAfter(Token token){
    println "Ran after the tagging step."
}
```

```
public void removeHeader(){
    println "Reading " + inPath
    text = text ~/((?s)<header>.*?</header>/
    autoOutput = true
}
```

Scripts run after part-of-speech tagging. This allows modifying the output of the Tree Tagger, e.g. correcting tagging errors.

An example of a simple script deleting header indicated by <header></header> tags in text.

### ► Did you know?

The Brown corpus and the LOB (Lancaster-Oslo/Bergen) corpus are one of the first modern corpora stored and processed on computers. Each consists of one million running words (tokens), a size that was very ambitious at the time of their compilation. Brown was compiled in the 1960s by Henry Kučera and W. Nelson Francis at Brown University (US). It was originally stored and processed on IBM punch cards. In the early 1970s, a British counterpart to the Brown corpus was compiled as a collaboration between Lancaster University (UK) and two Norwegian universities: Oslo and Bergen. The project was initiated by Geoffrey Leech from Lancaster University.

### 3 Key functionalities

This section reviews key functionalities of #LancsBox that are common to multiple #LancsBox modules.

---

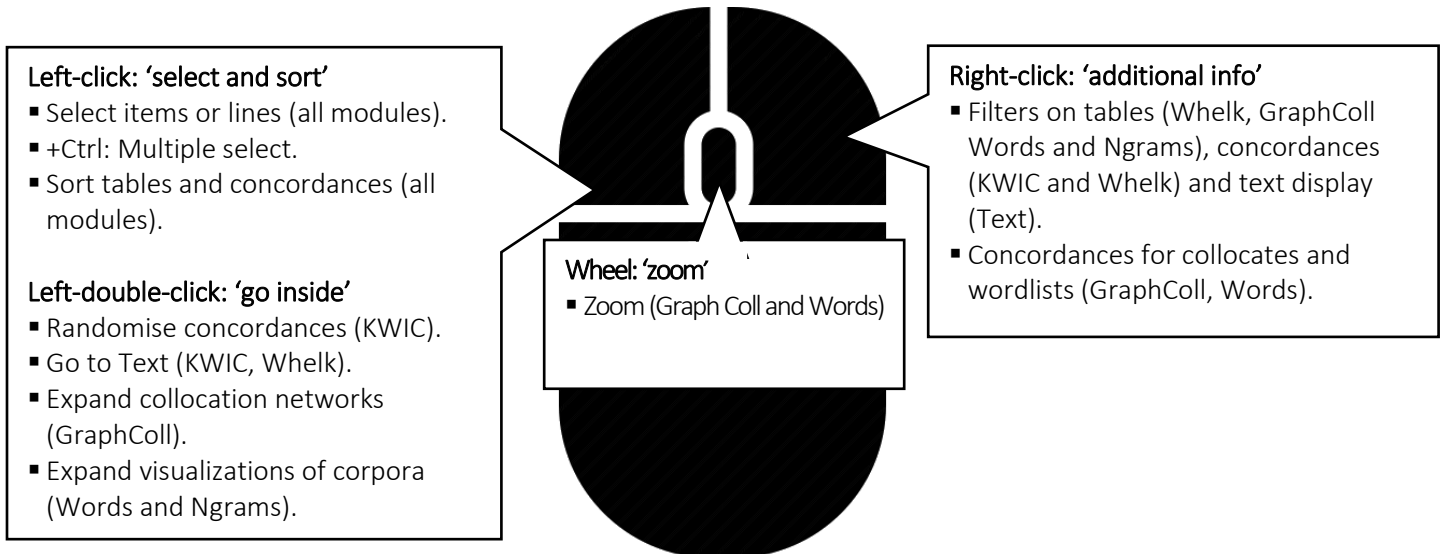
#### 3.1 Mouse clicks

#LancsBox doesn't use drop-down menus. Instead, all commands are literally just one mouse click away.

---



Hover with the mouse pointer for tooltips (brief contextual explanation of key functionalities/terms) to appear.



► **Note:** Mac users need to review their specific setup of the mouse clicks. By default, right-click is defined as Control + click. Alternatively, a standard two-button mouse with a wheel can be connected to a Mac machine.

#### 3.2 Shortcut Keys

#LancsBox allows changing the size of the text for easy readability. This works both in graphs and tables.

---

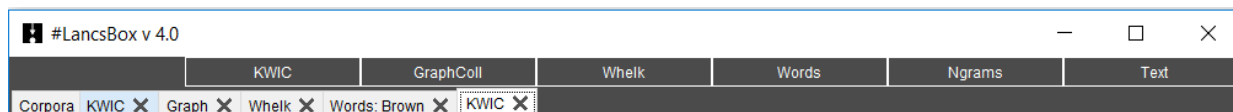
|                       |            |
|-----------------------|------------|
| Make all text bigger  | Ctrl and + |
| Make all text smaller | Ctrl and - |



### 3.3 Tools and Tabs

#LancsBox supports multiple simultaneous analyses and multiple corpora. #LancsBox has five main modules (tools): KWIC, Whelk, GraphColl, Words and Text. Each tool can be called multiple times on separate tabs. The modules in #LancsBox are interconnected: they can be launched as pop-ups inside a module.

1. The figure below show the top bar in #LancsBox with buttons for individual modules and multiple tabs open.



2. The modules in #LancsBox have the following functionalities:

KWIC produces concordances.

Whelk shows distribution of the search term in corpus files.

GraphColl identifies and visualizes collocations.

Words produces wordlists and identifies and visualizes keywords.

Ngrams produces lists of ngrams and identifies and visualizes key ngrams.

Text displays a full context of a search term.

### 3.4 Split screen

#LancsBox supports split-screen comparisons that allow displaying two separate analyses, one in the top and one in the bottom panel.


3. To use split screen, left-click on a bar with three triangles: ▲▲▲. This brings up the bottom panel.
4. To activate the bottom (or the top) panel in the split-screen view, left-click on the panel. An active panel is indicated by a light blue border (□).
5. To close the split-screen view, left-click on the bar with three triangles: ▼▼▼. This will hide the bottom panel but will not clear the results, so the bottom panel can be brought back later, if needed.

### 3.5 Saving results

---

#LancsBox supports easy saving of results. It saves concordances, wordlists, tables and graphics.

---

1. To save the results that #LancsBox produces, left-click on the save icon () in the top right-hand corner.
2. Select the location where you wish to save the results.
3. Click 'Save'.

### 3.6 Copy/pasting selected results

---

#LancsBox supports easy copy/pasting of selected results.

---

1. Select results which you wish to copy/paste by left-clicking on them; the results will be highlighted. To select discontinuous results, hold down Ctrl while selecting. To select all results, press Ctrl + A [Mac: Command + A].

| Index | File        | Left  | Node   | Right  |
|-------|-------------|---|--------|--|
| 1     | A_Press_rep | and Juliet" was the irresponsibility of young | love   | pushed into tragedy by Shakespeare." Othello" is |
| 2     | A_Press_rep | a cultivated, brave man who comes to          | love   | too late, and does not know what                 |
| 3     | A_Press_rep | not to know what to do with                   | love." | Zeffirelli does not mention the colour of        |
| 4     | A_Press_rep | Logue writes fierce, noisy poems about war,   | love,  | and Logue. Son of a Southampton civil            |
| 5     | A_Press_rep | go up in flames one day. In                   | love,  | he wrote:—" I can not see Smiles                 |
| 6     | C_Press_rev | Byron's leaving him, the scandal of his       | love   | affair with his half-sister, Augusta Leigh, the  |
| 7     | C_Press_rev | him to a point that looks like                | love,  | had fanned the enthusiasm which had sent         |

2. Press Ctrl + C [Mac: Command + C].
3. In the new location (e.g. text file, spreadsheet) press Ctrl + V [Mac: Command + V].

## 4 KWIC tool (key word in context)

The KWIC tool generates a list of all instances of a search term in a corpus in the form of a concordance. It can be used, for example, to:

- Find the frequency of a word or phrase in a corpus.
- Find frequencies of different word classes such as nouns, verbs, adjectives.
- Find complex linguistic structures such as the passives, split infinitives etc. using ‘smart searches’.
- Sort, filter and randomise concordance lines.
- Perform statistical analysis comparing the use of a search term in two corpora.

### 4.1 Visual summary of KWIC tab

The screenshot shows the KWIC tool interface with the search term 'love'. The interface includes a search bar, a list of concordance lines, and a bottom panel for filtering. Callouts provide the following instructions:

- Save results**: Located at the top right of the interface.
- Statistical analysis**: Located at the top right of the interface.
- Left-double-click 'Index' to randomise concordance lines.**: Points to the 'Index' column header.
- Left-click concordance header to sort.**: Points to the 'File' column header.
- Right-click concordance header to use advanced filter.**: Points to the 'Context' column header.
- Left-double-click concordance display to see text.**: Points to a concordance line.
- Right-click inside to apply filter.**: Points to the text within a concordance line.
- Pull up the bottom panel.**: Points to the bottom panel.

#### Main search box

#### You can:

- Search for a word or phrase.
- Search for number ranges, e.g. >1930&<=1945
- Use \* wildcards, e.g. new\*
- Use case sensitive regular expressions, e.g. /[abc].\*/
- Use case insensitive regular expressions, e.g. /dog|cat/i
- Search for punctuation, e.g. /.\*\./p
- Use ‘smart searches’, e.g. PASSIVE, NOUN
- Use Corpus Query Language (CQL)

#### Additional search boxes

**Headword**  
 **POS**

#### You can:

- Search at different levels of annotation.
- Combine search terms at various levels.
- Use regular expressions, e.g. /N.\*/
- Define batch searches.

## 4.2 Searching and displaying results

#LancsBox supports powerful searching of corpora. The search box can be used for simple as well as advanced searches at different levels of annotation.

1. Simple searches: type in the word or phrase of interest in the search box in the top left-hand corner and left-click 'Search'.
2. Advanced searches: click on the triangle inside the search box (▼) to activate advanced searches at different levels of corpus annotation. You can type search terms as separate constraints into one or more advanced search boxes. For example, the following advanced search is a search for the lemma 'go'.

|    |          |
|----|----------|
|    | Search   |
| go | Headword |
| V* | POS      |

Text level empty → no constraint.

Headword is go.

AND

POS is any verbal use.

3. A concordance is generated. The search term, called the 'node', is positioned in the centre and highlighted (orange colour), with words displayed to the left and right of it.
4. KWIC displays basic information about the frequency of the search term and its distribution in texts; the second example shows an application of a filter (see Section 4.4):

|        |          |             |            |       |       |
|--------|----------|-------------|------------|-------|-------|
| Search | research | Occurrences | 158 (1.57) | Texts | 13/15 |
|--------|----------|-------------|------------|-------|-------|

Read: The search term 'research' occurs 158 times in the corpus with the relative frequency 1.57 per 10k words in 13 out of 15 texts.

|        |          |             |              |       |      |
|--------|----------|-------------|--------------|-------|------|
| Search | research | Occurrences | 7/158 (0.07) | Texts | 3/15 |
|--------|----------|-------------|--------------|-------|------|

Read: When a filter is applied (indicated by blue colour), the search term 'research' occurs 7 times out of 158 in the corpus with the relative frequency 0.07 per 10k words in 3 out of 15 texts.

## 4.3 Settings and full text pop-up

KWIC settings include Corpus, Context and Display options. KWIC also allows full-text pop-ups.

1. Corpus: this setting changes the corpus which is being searched. Note that different corpora can be searched in the top and bottom panel in split-screen view.
2. Context: this setting changes the number of words that are displayed in the concordance to the left and to the right of the node.
3. Display: this setting changes the display type. The 'Plain text' default can be changed to 'Text with POS', 'Lemmatized text' and 'All annotation'. The example below demonstrates these four display formats:

**Plain text:** The new life looks promising for Mr. Noyce.

**Text with POS:** The\_DT new\_JJ life\_NN looks\_VVZ promising\_JJ for\_IN Mr.\_NP Noyce.\_NP

**Lemmatized text:** the\_DT new\_JJ life\_NN look\_VVZ promising\_JJ for\_IN Mr\_NP Noyce\_NP

**All annotation:** [The{the}\_DT] [new{new}\_JJ] [life{life}\_NN] [looks{look}\_VVZ] [promising{promising}\_JJ]  
[for{for}\_IN] [Mr.{Mr}\_NP] [Noyce.{Noyce}\_NP]

4. Full text pop-up: Double left-click on a concordance line to display the entire text with the appropriate line highlighted.

## 4.4 Sorting, randomising and filtering

---

KWIC concordance can be sorted alphabetically, randomised and filtered.

---

1. Alphabetical sorting: Left-click the concordance header (any column) to sort the column alphabetically in the A-Z (ascending) order; click again to re-sort alphabetically in the Z-A (descending) order. The sorting is indicated by arrows: A-Z (▲) and Z-A (▼).
2. Randomising: Left-double-click the header of the 'Index' column to randomise the concordance lines. Randomisation is indicated by the tilde sign (~).
3. Simple filtering: Right-click anywhere inside the concordance to activate the simple filter on that column. Input a word or phrase or a regular expression enclosed in forward slashes (/ /) and click 'Apply'. Filtering is indicated by light blue colour of the filtered text. The filter also updates the results (Occurrences and Texts) in the top display panel (see Section 4.2, point 4).
4. Advanced filtering: Right-click any part of the concordance header to activate the advanced filter. Select an exact column or position for filtering (see below), enter value and click 'Add' and 'Apply'. Filtering is indicated by light blue colour on text and in the results display panel (Occurrences and Texts).

An example of positions for advanced filtering:

L5 L4 L3 L2 L1 Node R1 R2 R3 R4 R5  
is Mr. Robert Weaver of New York. One of his tasks

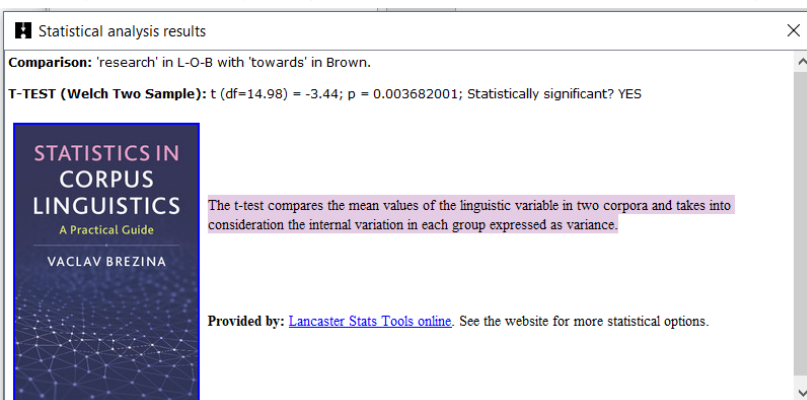
## 4.5 Statistical analysis

---

KWIC connects to Lancaster Stats Tools online to perform statistical analysis of the data in split panels.

---

When search results appear in both the top and the bottom panel in split-screen, these can be compared by clicking on the statistical analysis button (📊). The tool automatically connects to Lancaster Stats Tools online (Brezina 2018) and performs the t-test. The results are reported as follows:



Statistical analysis results

Comparison: 'research' in L-O-B with 'towards' in Brown.

T-TEST (Welch Two Sample): t (df=14.98) = -3.44; p = 0.003682001; Statistically significant? YES

STATISTICS IN CORPUS LINGUISTICS  
A Practical Guide  
VACLAV BREZINA

The t-test compares the mean values of the linguistic variable in two corpora and takes into consideration the internal variation in each group expressed as variance.

Provided by: [Lancaster Stats Tools online](#). See the website for more statistical options.

### ► Did you know?

In 1992, when reviewing the state of the art in corpus linguistics, Leech (1992) considers a concordance program “[t]he simplest and the most widely-used tool for corpus-based research” (p. 114). 25 years later, a concordance program such as KWIC still belongs to the essential toolkit of a corpus linguist. The simple and direct access to data that a concordance program facilitates combined with more sophisticated functions such as sorting, filtering and randomising provides a powerful analytical technique.

Leech, G. (1992). Corpora and theories of linguistic performance. In: *Directions in corpus linguistics*, 105-122.

## 5 Whelk tool

The Whelk tool provides information about how the search term is distributed across corpus files.

It can be used, for example, to:

- Find absolute and relative frequencies of the search term in corpus files.
- Filter the results according to different criteria.
- Sort files according to absolute and relative frequencies of the search term.

### 5.1 Visual summary of Whelk tab

The screenshot shows the Whelk tool interface. The top panel displays search results for the term 'love' across various corpora. The bottom panel shows a table of the search term's distribution across individual files.

| File                   | Tokens | Frequency | Relative frequency per 10k |
|------------------------|--------|-----------|----------------------------|
| P_Romance.txt          | 58197  | 75        | 12.887262                  |
| C_Press_review.txt     | 34289  | 39        | 11.37391                   |
| A_Fiction_gen.txt      | 58515  | 69        | 10.253783                  |
| L_Fiction_myst.txt     | 48259  | 15        | 3.1022284                  |
| F_Pop_lore.txt         | 88742  | 26        | 2.9298415                  |
| N_Adventure.txt        | 58322  | 16        | 2.7433903                  |
| G_Belle_lett_inogr.txt | 155271 | 35        | 2.2543234                  |
| E_Skills.txt           | 76613  | 16        | 2.0884185                  |
| M_Science_lect.txt     | 12037  | 2         | 1.6615435                  |
| D_Religion.txt         | 34257  | 4         | 1.1675448                  |
| J_Acad_writing.txt     | 161289 | 10        | 0.6200051                  |
| A_Press_report.txt     | 88805  | 5         | 0.5630314                  |
| R_Humour.txt           | 18087  | 1         | 0.55298327                 |
| B_Press_crit.txt       | 54397  | 0         | 0.0                        |
| H_Misc_non_lect.txt    | 60627  | 0         | 0.0                        |

Top panel: Searching corpora

You can:

- Search, sort and filter.
- Use simple and advanced searching functionality.
- Use 'smart' searches.

Bottom panel: Displaying distribution

You can:

- View the distribution of the search term in individual files.
- Sort, filter and copy/paste.

### 5.2 Top panel: KWIC

The top panel in Whelk has the same powerful search, sort and filter functionalities as the KWIC tool (see Section 4). It is directly connected to the bottom panel: any update in the top panel is immediately reflected in the bottom panel.

### 5.3 Bottom panel: Frequency distribution

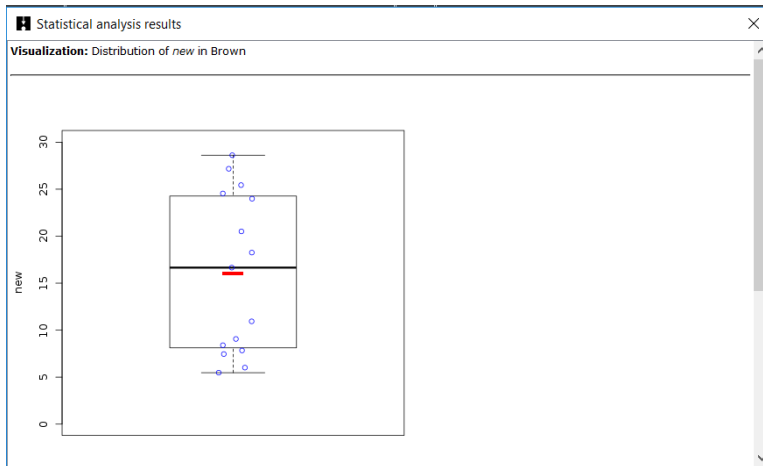
The bottom panel in Whelk provides detailed information about the distribution of the search term.

1. 'File' column lists the name of the individual files in the corpus.
2. 'Tokens' column provides the information about the size of each file in running words (tokens).
3. 'Frequency' column provides absolute frequencies of the search term i.e. refers to how many instances of the search term there are in each file.
4. 'Relative frequency per 10k' provides relative frequency normalised to the basis of 10,000 tokens; this value is comparable across files and corpora.

## 5.4 Statistical analysis

Whelk connects to Lancaster Stats Tools online to perform statistical analysis of the data.

When search results appear, these can be visualised using a boxplot by clicking on the statistical analysis button (📊). The tool automatically connects to Lancaster Stats Tools online (Brezina 2018) and displays the result:



### ► Did you know?

The Whelk tool (both the name and the functionality) is inspired by Kilgarriff's (1997: 138ff) notion of the 'whelks problem'. Imagine, says Kilgarriff, that you have a corpus which includes one text (a book) about whelks – small snail-like sea creatures (🐌). In this text, the word *whelks* will appear many times and hence will appear as a frequent word in the entire corpus, although its use is limited to one specific context. To overcome the problem and present more accurate information about word distribution, the Whelk tool shows the frequency distribution of search terms in individual corpus files.



## 6 GraphColl

The GraphColl tool identifies collocations and displays them in a table and as a collocation graph or network.

It can be used, for example, to:

- Find the collocates of a word or phrase.
- Find colligations (co-occurrence of grammatical categories).
- Visualise collocations and colligations.
- Identify shared collocates of words or phrases.
- Summarise discourse in terms of its 'aboutness'.

### 6.1 Visual summary of GraphColl tab

The screenshot shows the GraphColl interface with a table of collocates on the left and a network graph on the right. Callout boxes provide instructions for various features:

- Save results.
- View options.
- Change collocation settings.
- Display collocates in a table.
- Display collocation graphs and networks.
- Pull up the bottom panel.

| Index | Status | Position |       |           |    |      |
|-------|--------|----------|-------|-----------|----|------|
| 1     | O      | L        |       |           |    |      |
| 2     | O      | R        |       |           |    |      |
| 3     | O      | R        |       |           |    |      |
| 4     | O      | R        |       |           |    |      |
| 5     | O      | L        |       |           |    |      |
| 6     | O      | M        | much  | 5.1538576 | 10 | 1072 |
| 7     | O      | L        | life  | 6.0951286 | 6  | 684  |
| 8     | O      | R        | years | 5.9931353 | 9  | 1086 |
| 9     | O      | R        | most  | 5.6568897 | 7  | 1059 |
| 10    | O      | L        | where | 5.4703566 | 6  | 1033 |
| 11    | O      | L        | how   | 5.4451597 | 5  | 876  |
| 12    | O      | R        | first | 5.9725773 | 7  | 1237 |
| 13    | O      | L        | on    | 5.0722905 | 11 | 7332 |
| 14    | O      | L        | my    | 5.0675035 | 8  | 1821 |
| 15    | O      | L        | their | 5.0278540 | 12 | 2808 |

### 6.2 Producing a collocation graph

GraphColl produces collocations graphs on the fly. After selecting the appropriate settings you can start searching for the node and its collocates.

1. Select the appropriate settings for the collocation search:
  - i) **Span**: how many words to the left (L) and to the right (R) of the node (search term) are being considered when searching for collocates [default: 5L, 5R].
  - ii) **Statistics**: the association measure used to compute the strength of collocation [default: frequency – no association measure is preferred because the choice depends on the research question].

- iii) Threshold: The minimum frequency and statistics cut-off values for an item (word, lemma, POS) to be considered a collocate.
  - iv) Corpus: The corpus that is being searched.
  - v) Unit: The unit (type, lemma, part of speech [POS] tag) used for collocates.
2. Type the search term into the search box (top left) and left-click 'Search'.
  3. This will produce a collocation table (left) and a collocation graph (right).

### 6.3 Reading collocation table

A collocation table is a traditional way of displaying collocates. In GraphColl, the table shows the following pieces of information for each collocate: i) status, ii) position, iii) stat, iv) collocation frequency and v) frequency of the collocate anywhere in the corpus. By default, the table is sorted according to the selected collocation statistic (largest-smallest).

1. The following is a visual description of the collocation table.

Right-click header: filter

Left-click header: sort

node (expanded)

Left-double-click: expand collocation network

Right-click: show concordance

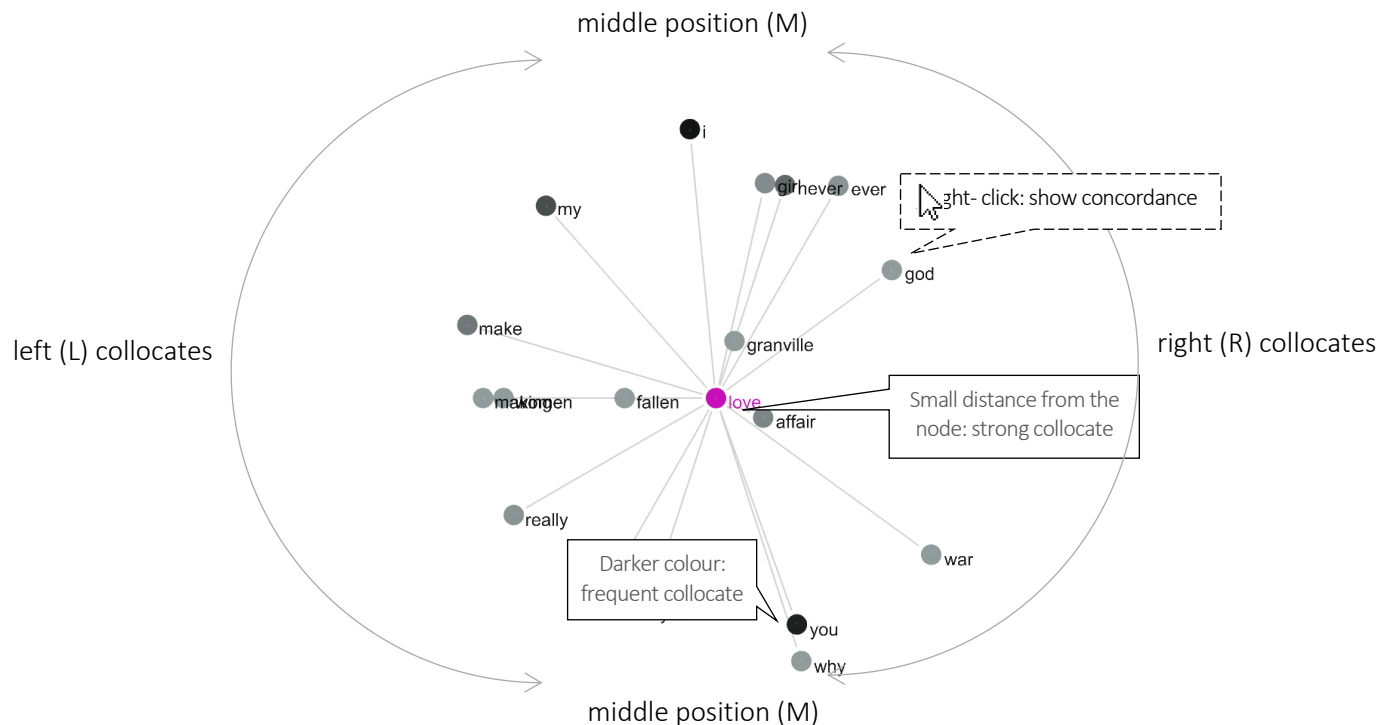
| Status | Position | Collocate | Stat           | Freq (coll.) | Freq (corpus) |
|--------|----------|-----------|----------------|--------------|---------------|
| ○      | R        | granville | 9.208810483... | 5            | 28            |
| ○      | R        | affair    | 9.179667251... | 7            | 40            |
| ○      | L        | fallen    | 8.623849712... | 5            | 42            |
| ●      |          | love      | 6.768239667... | 10           | 304           |
| ○      | L        | making    | 5.999358660... | 5            | 259           |
| ○      | L        | really    | 5.989182871... | 6            | 313           |
| ○      | R        | girl      | 5.754072064... | 5            | 307           |
| ○      | R        | you       | 5.392047146... | 46           | 3630          |
| ○      | R        | war       | 5.386810358... | 5            | 396           |
| ○      | L        | my        | 5.256041267... | 21           | 1821          |
| ○      | L        | make      | 5.223919105... | 9            | 798           |
| ○      | L        | never     | 5.075197...    | 7            | 688           |

2. The meaning of the individual columns is:
  - i) Status: shows whether the collocate has been expanded; ○ indicates a non-expanded collocate, while ● indicates expanded collocate (node) in a collocation network.
  - ii) Position: shows textual position of the collocate, which can be either left (L) of the node, right (R) of the node or middle (M), i.e. with equal frequency L and R.
  - iii) Collocate: shows the collocate in question.
  - iv) Stat: displays the value of the selected association measure.
  - v) Freq (coll): displays the frequency of the collocation (combination of node + collocate).
  - vi) Freq (corpus): displays the frequency of the collocate anywhere in the corpus.

## 6.4 Reading collocation graph

The graph displays three dimensions: i) strength of collocation, ii) collocation frequency and iii) position of collocates. To find out more about a collocate, right-click on it to obtain concordance lines (KWIC), in which the collocates co-occurs with the node.

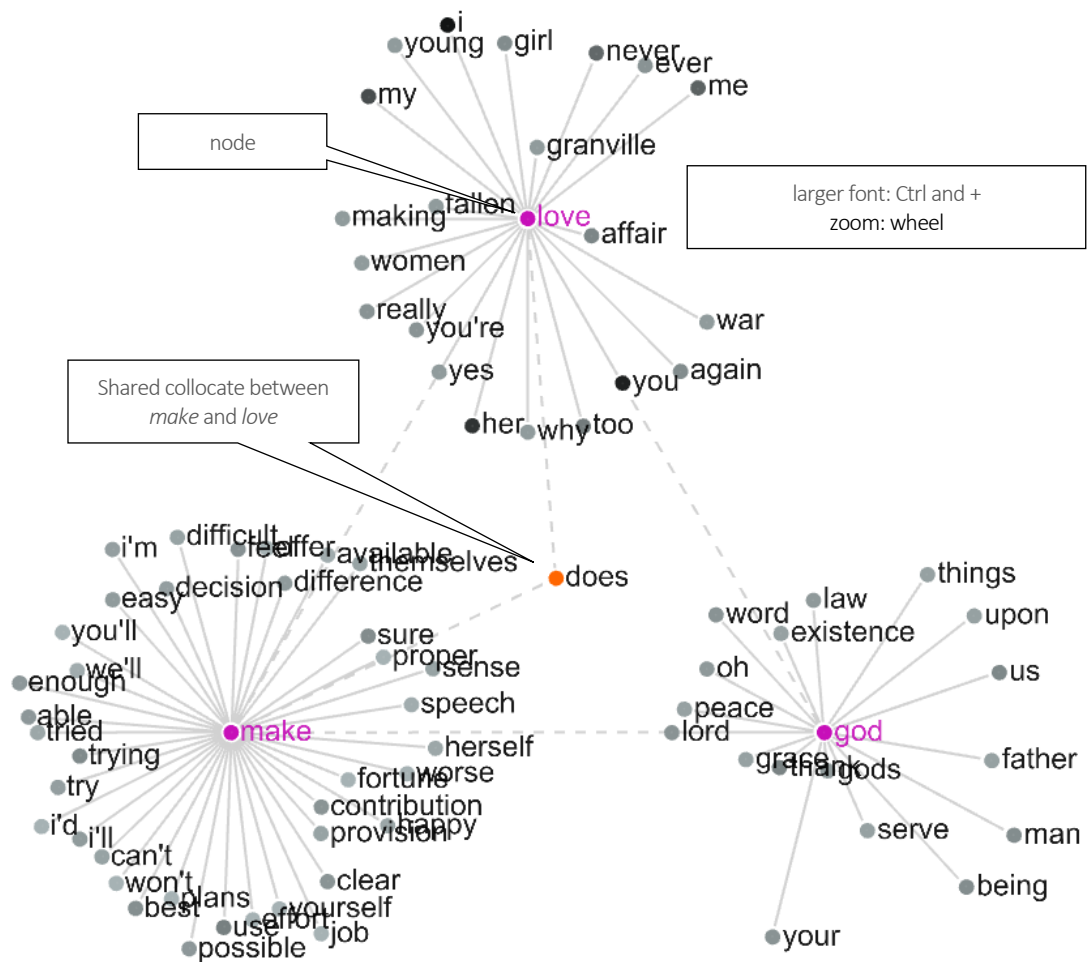
1. **Strength:** The strength of collocation as measured by the association measure is indicated by the distance (length of line) between the node and the collocates. The closer the collocate is to the node, the stronger the association between the node and the collocate ('magnet effect').
2. **Frequency:** Collocation frequency is indicated by the intensity of the colour of the collocate. The darker the shade of colour, the more frequent the collocation is.
3. **Position:** The position of collocates around the node in the graph reflects the exact position of the collocates in text: some collocates appear (predominantly) to the left of the node, others to the right; others still appear sometimes left and sometimes right (middle position in the graph). For the ease of display (if multiple collocates appear in a similar position and hence overlap), the tool allows 'spreading out' collocates evenly around the node. This is done by clicking on the 'Spread out' button (top right). When this is done, the collocates are dispersed evenly around the node with a 'L' or 'R' index displayed above the collocate circle indicating their original position to the left and to the right respectively.



## 6.5 Extending graph to a collocation network

A collocation network is an extended collocation graph that shows i) shared collocates and ii) cross-associations between several nodes.

1. To expand a simple collocation graph (see above) into a collocation network, either search for more nodes or left-double-click on a collocate in either the table or the graph.
2. A collocation network displays nodes with unique collocates (outer rim of the graph) and shared collocates (middle of the graph). The links between nodes and shared collocates are indicated by a dash-dot line (— · — · —).



## 6.6 Shared collocates

Shared collocates are collocates shared by at least two nodes in the graph. Shared collocates are displayed in the middle of the graph with links to the relevant nodes.

1. A full list of shared collocates can be obtained by clicking on the text 'Shared collocates'.

**old**

Freq: 660 - Collocates: 139 - [Shared collocates: 50](#)

2. The list of shared collocates is displayed in a tabular form.

| ID | ▲ Collocate | Freq (corpus) | No of nodes | Nodes    |
|----|-------------|---------------|-------------|----------|
| 1  | a           | 23204         | 2           | new, old |
| 2  | about       | 1815          | 2           | new, old |
| 3  | after       | 1070          | 2           | new, old |
| 4  | and         | 28852         | 2           | new, old |
| 5  | are         | 4395          | 2           | new, old |
| 6  | as          | 7251          | 2           | new, old |
| 7  | at          | 5377          | 2           | new, old |
| 8  | because     | 883           | 2           | new, old |
| 9  | being       | 712           | 2           | new, old |
| 10 | called      | 401           | 2           | new, old |



## 6.8 Reporting collocates: CPN

It is important to realise that there is no one definite sets of collocates: different statistical procedures and threshold values highlight different sets of collocates. We therefore need to report the statistical choices involved in the identification of collocations using standard notation called Collocation Parameters Notation (CPN). When saving the results, GraphColl saves the settings in the form of CPN.

Brezina et al. (2015) propose CPN as a specific notation to be used for accurate description of collocation procedure and replication of the results. The following parameters are reported.

| Statistic ID                                     | Statistic name | Statistic cut-off value | L and R span | Minimum collocate freq. (C) | Minimum collocation freq. (NC) | Filter                 |
|--|----------------|-------------------------|--------------|-----------------------------|--------------------------------|------------------------|
| 4b   | MI2            | 3                       | L5-R5        | 5                           | 1                              | function words removed |
| 4b-MI2(3), L5-R5, C5-NC1; function words removed |                |                         |              |                             |                                |                        |

### ► Did you know?

The name GraphColl is an acronym for *graphical collocations* tool. GraphColl was the first module in #LancsBox (v.1.0) with the other tools being added at a later stage. Graphical display of collocations and collocation networks is inspired by the work of Phillips (1985), who demonstrated the concept of lexical networks (Phillip's term for 'collocation networks') with small specialised corpora. GraphColl takes this notion further, offering different statistical choices and producing collocation networks on the fly with both small and large corpora.

Phillips, M. (1985). *Aspects of text structure: An investigation of the lexical organisation of text*. Amsterdam: North-Holland.

## 7 Words tool

The Words tool allows in-depth analysis of frequencies of types, lemmas and POS categories as well as comparison of corpora using the keywords technique.

It can be used, for example, to:

- Compute frequency and dispersion measures for types, lemmas and POS tags.
- Visualize frequency and dispersion in corpora.
- Compare corpora using the keyword technique.
- Visualize keywords.

### 7.1 Visual summary

The screenshot shows the #LancsBox v 4.0 interface. The top menu bar includes KWIC, GraphColl, Whelk, Words, Ngrams, and Text. The main window is titled 'Corpora Words: L-O-B, Brown'. It features two data tables and a visualization area.

**Table 1: L-O-B Corpus**

| Type | Frequency: 01 - Freq | Dispersion: 01_CV |
|------|----------------------|-------------------|
| the  | 21197.000000         | 0.162650          |
| of   | 11196.000000         | 0.291942          |
| and  | 11003.000000         | 0.078241          |
| to   | 10516.000000         | 0.081485          |
| a    | 10034.000000         | 0.114962          |
| in   | 9304.000000          | 0.208365          |
| that | 8796.000000          | 0.175456          |
| is   | 7339.000000          | 0.565229          |
| was  | 7201.000000          | 0.461677          |
| it   | 7188.000000          | 0.225496          |
| for  | 6997.000000          | 0.187196          |
| he   | 6353.000000          | 0.635331          |
| as   | 5652.000000          | 0.097503          |
| with | 4907.000000          | 0.111463          |
| be   | 4907.000000          | 0.277399          |
| on   | 4907.000000          | 0.160148          |

**Table 2: Brown Corpus**

| Type | Frequency: 01 - Freq | Dispersion: 01_CV |
|------|----------------------|-------------------|
| the  | 69970.000000         | 0.108141          |
| of   | 36408.000000         | 0.275038          |
| and  | 28852.000000         | 0.091673          |
| to   | 26148.000000         | 0.085560          |
| a    | 23204.000000         | 0.133061          |
| in   | 21340.000000         | 0.180183          |
| that | 17568.000000         | 0.175684          |
| is   | 15546.000000         | 0.554664          |
| was  | 14907.000000         | 0.490791          |
| he   | 13306.000000         | 0.655043          |
| for  | 9488.000000          | 0.202990          |
| it   | 8762.000000          | 0.265833          |
| with | 7289.000000          | 0.113924          |
| as   | 7251.000000          | 0.161254          |
| his  | 6996.000000          | 0.497809          |
| on   | 6747.000000          | 0.160148          |

**Visualization:** Two circles represent corpora: a pink circle labeled 'L-O-B' and a black circle labeled 'Brown'. Callouts explain: 'Drag corpora together to produce keywords.', 'Left-double-click on the corpus to see its internal structure.', and 'Right-click on the corpus to see corpus statistics.'

**Left:** Creating frequency lists, computing dispersion and keywords.

**Right:** Visualizing frequencies, dispersions and keywords.



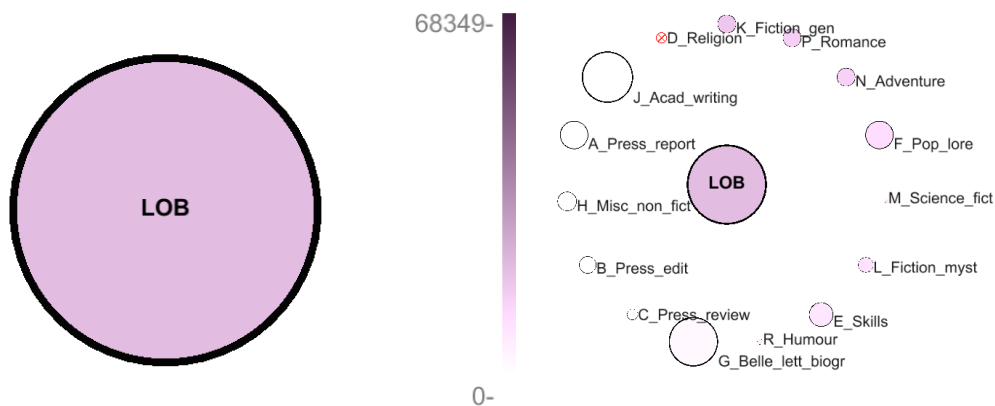
## 7.2 Producing frequency list

On start, Words produces a frequency list (table) based on the default corpus (see Section 2.6, point 2) and default settings. These settings can be changed and a different frequency list is produced.

1. The following are the settings for frequency lists:
  - i) Corpus: The corpus that is being used.
  - ii) Frequency: Absolute or relative frequency [default: absolute frequency].
  - iii) Dispersion: The dispersion statistic [default: coefficient of variation (CV)].
  - iv) Unit: The unit used in the frequency list (type, lemma or part of speech tag).
2. Changing any of these settings triggers re-computing of the frequency list.
3. Frequency lists can be searched using the search box (top left).
4. Frequency lists can be sorted by left-clicking on the header.
5. Frequency lists can be filtered by right-clicking on the header and applying a filter.
6. Two different frequency lists can be computed in the split-screen view, which is triggered by left-clicking on a bar with three triangles: ▲▲▲. This brings up the bottom panel.

## 7.3 Visualizing frequency and dispersion

The Words module displays corpora and corpus files (when a corpus is left-double-clicked). It visualises frequency and dispersion of words using intensity of colour and position of individual files displayed as circles; the size of the circle indicates the relative size of the corpus/file.



Display of frequency in the whole corpus on the scale of 0 - 68,349 (most frequent item).      Display of frequency per file (when corpus is left-double-clicked).

1. To visualize frequency of an item in the table, left-click on the item in the frequency table. The shade of the colour of the corpus will change according to the frequency value of this item. The scale on the right offers a reference point for interpretation.
2. To visualize dispersion of an item in the table, left-double-click on the corpus (large circle). The corpus will expand to display individual files (small circles) of which the corpus consists. The size of each circles is proportional to the size of the corpus subpart. The shade of the colour of the


small circles will change according to the frequency value of the item in the frequency list. Crossed-out (⊗) circles indicate that the item does not occur in the given corpus file. In addition, the corpus files are ordered according to the relative frequency of the item with the file with the largest relative frequency of the item appearing at the 12-o'clock position ( ) and the other files ordered clockwise according to decreasing relative frequency of the item ( ).

## 7.4 Producing keywords

---

The Words module computes a comparison of frequencies between two corpora/wordlists using a selected statistical measure. It identifies and visualizes positive keywords, negative keywords and lockwords.

---

1. Left-click on ▲▲▲ to bring up the bottom panel.
2. In the bottom panel, select a comparison (reference) corpus, while in the top panel keep your corpus of interest.
3. In the visualisation panel (right), drag the circles that represent the two corpora together . Alternatively, press the space bar.
4. The resulting table will display frequency and dispersion info about the two corpora as well as the keyword statistic; the graphics will identify top 10 positive keywords, top 10 negative keywords and top 10 lockwords.
5. In the settings, you can change the i) keyword statistic and ii) threshold.  
Keyword statistic: This is a measure that compares two frequency lists [default: simple maths with constant  $k = 100$ ].  
Threshold: Threshold values for the identification of positive keywords, negative keywords and (by implication) lockwords.

## 7.5 Producing corpus statistics

---

The Words module computes essential corpus statistics: i) Complexity stats and ii) Lexical stats

---

1. Right-click on corpus .
2. In the pop-up table toggle between Complexity stats and Lexical stats.

## Mean sentence length and Standard deviation (SD)

| ▼ Complexity Stats    |                        | ▼ Lexical Stats      |                    |                  |  |
|-----------------------|------------------------|----------------------|--------------------|------------------|--|
| File                  | Sentence Length (mean) | Sentence Length (SD) | Word Length (mean) | Word Length (SD) |  |
| A_Press_report.bt     | 19.159855              | 11.671002            | 4.745014           | 2.592452         |  |
| B_Press_edit.bt       | 20.061825              | 12.509228            | 4.734839           | 2.6490588        |  |
| C_Press_review.bt     | 22.179173              | 14.621512            | 4.77955            | 2.7150402        |  |
| D_Religion.bt         | 19.105968              | 13.838464            | 4.5256734          | 2.5267594        |  |
| E_Skills.bt           | 20.938234              | 13.569921            | 4.603331           | 2.51522          |  |
| F_Pop_lore.bt         | 21.013971              | 12.89571             | 4.6807714          | 2.5748186        |  |
| G_Belle_left_biogr.bt | 24.429043              | 15.205565            | 4.714493           | 2.6827366        |  |
| H_Misc_non_fict.bt    | 25.527159              | 20.760244            | 4.882379           | 2.7997973        |  |
| J_Acad_writing.bt     | 26.358719              | 16.505852            | 4.851614           | 2.8534663        |  |
| K_Fiction_gen.bt      | 14.338397              | 12.206561            | 4.3068104          | 2.27138          |  |
| L_Fiction_myst.bt     | 12.934602              | 9.881333             | 4.30815            | 2.259926         |  |
| M_Science_fict.bt     | 12.371017              | 11.275793            | 4.5213094          | 2.4187284        |  |
| N_Adventure.bt        | 11.963488              | 9.186616             | 4.262817           | 2.1768787        |  |
| P_Romance.bt          | 12.555987              | 9.679649             | 4.236387           | 2.1641104        |  |
| R_Humour.bt           | 17.87253               | 14.513976            | 4.5027366          | 2.506564         |  |

## Type-token ratio (TTR), Standardised type-token ratio (STTR), Moving average type-token ratio (MATTR)

| ▼ Complexity Stats    |       | ▼ Lexical Stats |             |            |            |
|-----------------------|-------|-----------------|-------------|------------|------------|
| File                  | Types | Tokens          | TTR         | STTR       | MATTR      |
| A_Press_report.bt     | 12079 | 88805           | 0.13601711  | 0.7342071  | 0.7342669  |
| B_Press_edit.bt       | 7909  | 54367           | 0.14547427  | 0.73095614 | 0.7306529  |
| C_Press_review.bt     | 7703  | 34289           | 0.22464931  | 0.74618065 | 0.74707484 |
| D_Religion.bt         | 5399  | 34257           | 0.15760283  | 0.69137025 | 0.6896752  |
| E_Skills.bt           | 10808 | 76613           | 0.14107266  | 0.72006595 | 0.7209448  |
| F_Pop_lore.bt         | 12274 | 88742           | 0.13831106  | 0.72313124 | 0.72300106 |
| G_Belle_left_biogr.bt | 17485 | 155271          | 0.112609565 | 0.7196904  | 0.7203814  |
| H_Misc_non_fict.bt    | 6717  | 60627           | 0.11079222  | 0.6818785  | 0.6824127  |
| J_Acad_writing.bt     | 15743 | 161289          | 0.097607404 | 0.685145   | 0.685025   |
| K_Fiction_gen.bt      | 7841  | 58515           | 0.13399982  | 0.7243858  | 0.72329557 |
| L_Fiction_myst.bt     | 6632  | 48259           | 0.13742514  | 0.7332717  | 0.7323574  |
| M_Science_fict.bt     | 3187  | 12037           | 0.26476696  | 0.7563636  | 0.7587221  |
| N_Adventure.bt        | 7638  | 58322           | 0.13096258  | 0.73029095 | 0.7307091  |
| P_Romance.bt          | 6525  | 58197           | 0.11211918  | 0.7355844  | 0.73544407 |
| R_Humour.bt           | 4452  | 18087           | 0.24614364  | 0.7351933  | 0.73470604 |

### ▶ Did you know?

The statistical technique of keyword analysis was originally developed by Mike Scott (1997) and it was implemented in WordSmith Tools. It relied on corpus comparison using the chi-squared test or the log-likelihood test. As Kilgarriff pointed out, the chi-squared test and the log-likelihood test are not entirely appropriate for this type of comparison. Kilgarriff's solution implemented in Sketch Engine was to compare corpora using a 'simple maths' procedure, a simple ratio between relative frequencies of words in the two corpora we compare. In addition to 'simple maths', #LancsBox offers also other types of solutions for corpus comparison.

Scott, M. (1997). PC analysis of key words—and key key words. *System*, 25(2), 233-245.

Kilgarriff, A. (2009, July). Simple maths for keywords. In *Proceedings of the Corpus Linguistics Conference*. Liverpool, UK.

## 8 Ngram tool

The Ngram tool allows in-depth analysis of frequencies of n-grams (bigrams, trigrams etc.), which could be defined as contiguous combinations types, lemmas and POS. The tool also produces key ngrams by comparing two corpora using a technique similar to keywords.

It can be used, for example, to:

- Identify n-grams, lexical bundles and p-frames (also skip grams)
- Compute frequency and dispersion measures for ngram types, lemmas and POS tags.
- Visualize frequency and dispersion of ngrams in corpora.
- Compare ngrams in two corpora using the keyword technique.
- Visualize key ngrams.

### 8.1 Visual summary

Right-click on the table header to activate filter.

| Type     | Frequency: 01 - Freq | Dispersion: 01_CV |
|----------|----------------------|-------------------|
| of the   | 18.000000            | 0.381724          |
| in the   | 00                   | 0.224633          |
| to the   | 00                   | 0.149529          |
| on th    | 00                   | 0.140736          |
| and t    | 00                   | 0.270452          |
| it is    | 1985.000000          | 0.652103          |
| for the  | 1977.000000          | 0.343772          |
| to be    | 1912.000000          | 0.224275          |
| at the   | 1745.000000          | 0.211144          |
| that the | 1651.000000          | 0.551571          |
| it was   | 1555.000000          | 0.553916          |
| with the | 1525.000000          | 0.258497          |
| from the | 1509.000000          | 0.159117          |
| of a     | 1501.000000          | 0.254168          |
| by the   | 1486.000000          | 0.503977          |
| in a     | 1259.000000          | 0.247329          |

Drag corpora together to produce key ngrams.

Left-double-click on the corpus to see its internal structure.

Right-click on the corpus to see corpus statistics.

Right-click inside the table to activate a Whelk pop-up.

| Type                 | Frequency: 01 - Freq | Dispersion: 01_CV |
|----------------------|----------------------|-------------------|
| your expenses        | 1.000000             | 3.741657          |
| owe additional       | 1.000000             | 3.741657          |
| foundation during    | 1.000000             | 3.741657          |
| surprise he          | 3.000000             | 2.176043          |
| health hazard        | 1.00                 | 3.741657          |
| parables being       |                      | 3.741657          |
| with lipstick        |                      | 3.741657          |
| sullam that          |                      | 3.741657          |
| drank slowly         |                      | 3.741657          |
| horsemanship classes | 1.000000             | 3.741657          |
| have fashioned       | 1.000000             | 3.741657          |
| for lunch            | 3.000000             | 2.055493          |
| themselves from      | 3.000000             | 2.102494          |
| unlikely synonyms    | 1.000000             | 3.741657          |
| noble or             | 1.000000             | 3.741657          |

**Left:** Creating frequency lists, computing dispersion and key ngrams.

**Right:** Visualizing frequencies, dispersions and key ngrams.

### ► Did you know?

Multi-word expressions are extremely important when describing language. There are different terms to describe multi-word expressions such as collocations (Brezina et al. 2015; Gablasova et al. 2017), n-grams, lexical bundles and p-frames. While collocations, which are identified in the GraphColl module, typically represent non-contiguous expressions, the n-gram type multi-word expressions represent contiguous lexico-grammatical patterns. They are defined as follows.

- n-gram: a sequence of n types, lemmas, POS from a text or corpus.
- lexical bundle: an ngram with certain frequency and distributional (dispersion) properties, e.g. relative freq. 10 per million and range > 5.
- p-frame (also skip gram): an n-gram that allows for variability at one or more positions such as *it would be \* to*.

All these types of multi-word expressions can be identified using the Ngram tool in #LancsBox.

Brezina, V., McEnery, T. & Wattam, S. (2015). Collocations in context: A new perspective on collocation networks. *International Journal of Corpus Linguistics*, 20(2), 139-173.

Gablasova, D., Brezina, V., & McEnery, T. (2017). Collocations in corpus-based language learning research: Identifying, comparing, and interpreting the evidence. *Language Learning*, 67 (S1), 155–179.

## 9 Text

The Text tool enables an in-depth insight into the context in which a word or phrase is used.

It can be used, for example, to:

- View a search term in full context.
- Preview a text.
- Preview a corpus as a run-on text.
- Check different levels of annotation of a text/corpus.

### 9.1 Visual summary

The screenshot displays the Text tool interface. At the top, the search term 'new' is entered, and the results show 181 occurrences (20.30 per 10k). The interface is divided into several sections: 'Search Term', 'Occurrences', 'Corpus', 'Text', and 'Text'. The main area shows a list of lines with the search term highlighted in red. A tooltip for the search term 'new' shows 'Absolute and relative frequency (per 10k)'. Two callout boxes provide additional information: one states 'All instances of a search term are highlighted in text.' and another states 'Up (↑) and down (↓) arrow to move between the occurrences.'

### 9.2 Searching in Text

Texts and corpora can be searched easily using a simple search box.

1. Type the search term into the search box (top left). Left-click 'Search'.
2. This will highlight all lines in the text where the search term appears in dark grey with the search term itself in red. To move between the highlighted lines up (↑) and down (↓) arrows can be used.
3. Frequency information (both an absolute and relative frequency per 10,000 tokens) will appear under 'Occurrences'.
4. A single line can be highlighted by left-clicking on the line. To highlight multiple lines, Ctrl (Command) + Left-click the desired lines.
5. Highlighted lines can be copied (Ctrl/Command+C) and pasted (Ctrl/Command+V) into a text editor.

### 9.3 Settings

---

The following settings are used in Text: i) Corpus, ii) Text and iii) Display.

---

1. Corpus: this setting allows changing the corpus which is being displayed and searched. Note that different corpora can be searched in the top and the bottom panel in the split-screen view.
2. Text: this setting allows changing the text that is being displayed and searched.
3. Display: this setting allows changing the display format. The 'Plain text' default can be changed to 'Text with POS', 'Lemmatized text' and 'All annotation'.

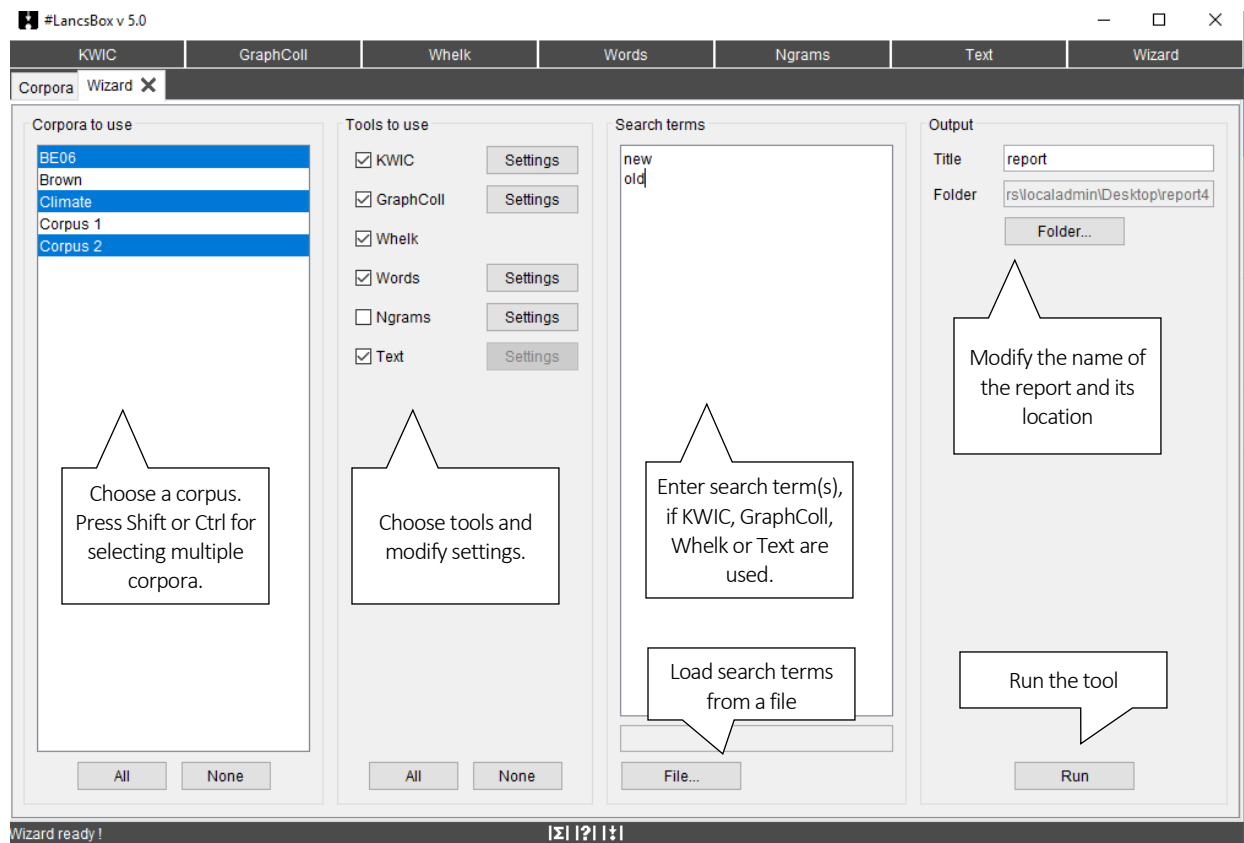
# 10 Wizard

The Wizard tool combines the power of all tools in #LancsBox, searches corpora and produces research reports for print (docx) and web (html).

It can be used, for example, to:

- Carry out simple or complex research.
- Produce a draft report.
- Download all relevant data.

## 10.1 Visual summary





## 10.2 Selecting settings and running Wizard

---

Wizard produces research reports automatically. All you need to do is select the corpus/corpora and procedures to use.

---





1. Select corpora in the 'Corpora to use' panel (left). For multiple adjacent corpora, hold Shift on the keyboard while selecting; for multiple corpora that are not next to each other, hold Ctrl (or Control on mac).
2. If the corpus you want does not appear on the list, go to 'Corpora' tab and add it.
3. Choose the tools you want to employ. The choice of the tools depends on the type of analysis you want to perform.
  - KWIC: analysis of concordance lines.
  - GraphColl: collocation analysis.
  - Whelk: analysis of frequencies in individual texts.
  - Words: analysis of frequencies and dispersions of individual lexical items and keyword analysis.
  - Ngrams: analysis of frequencies and dispersions of ngrams.
  - Text: analysis of broader contexts.
4. Adjust the tool settings by clicking on the 'Settings' button next to the tool.
5. Enter search term(s), if, KWIC, GraphColl, Whelk or Text are used. Alternatively, load search terms from a text (txt) file.
6. Choose the location where the report and extracted data will be saved or leave default (Desktop).
7. Press 'Run'

## 10.3 Data analysis

---

Wizard produces the data analysis in the background and informs the user about the progress. The complete data set is saved together with a report.

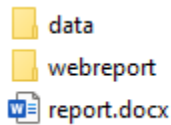
---

1. The data set includes the following folders.
  -  csv
  -  images
  -  tsv
  -  xml
2. csv (comma separated files; open in Excel or Calc) include data for Lancaster Stats Tools online.
3. images (png; open in a standard graphics app) include graphs and other graphical output.
4. tsv (tab separated files; open in Excel or Calc) include complete data output from the individual tools.
5. xml (extensible markup language; opens in a text editor) includes a complete Wizard data set, which can be used for further computational processing.

## 10.4 Research report

Wizard produces a structured data report in two formats: docx and html

1. The .docx report can be easily edited in Word, Writer or a similar word processor.



2. The length and contents of the report depend on the number of corpora and tools that were selected.
3. The report follows the structure of an academic research report.

Created by #LancsBox Wizard

June 9, 2020 - 22:07

# Comparison of British and American English

## 1 Introduction

This research report was automatically produced by #LancsBox (Brezina et al. 2020), a corpus analysis tool developed at Lancaster University. It uses cutting-edge technology and statistical sophistication (Brezina 2018) to analyze and visualize corpus data. For more information and tips on research report writing see the [Research Report Guide](#).

## 2 Method

### 2.1 Data

The study analyzed the following corpora:

Table 1. Corpora used

| Name  | Language | Texts | Tokens    | Additional information          |
|-------|----------|-------|-----------|---------------------------------|
| Brown | English  | 15    | 1,014,361 | Types: 49,686<br>Lemmas: 44,622 |
| L-O-B | English  | 15    | 1,007,677 | Types: 48,349<br>Lemmas: 43,920 |

In the study, 2 corpora were used of the total size of 2,022,038 running words (tokens) in 30 texts. A full description of the corpora is available in [data\tsv\corpora](#).

### 2.2 Procedure

#LancsBox (Brezina et al. 2020) software package was employed to analyse the data. The following tool from the package was used: KWIC. The KWIC tool generates a list of all instances of a search term in a corpus in the form of a concordance. The following search terms were used: "new", "old" and "some".

## 3 Results

## 11 Searching in #LancsBox

Throughout the tool, #LancsBox offers powerful searches at different levels of corpus annotation using i) simple searches, ii) wildcard searches, iii) smart searches, iv) regex searches and v) batch searches. In addition, from v. 5.1, complex searches are available using multiple conventions including the Corpus Query Language (CQL).

1. Simple searches are literal searches for a particular word (*new*) or phrase (*New York Times*). Simple searches are case insensitive; this means that *new*, *New*, *NEW*, *NeW* etc. will return the same set of results.
2. Wildcard searches are searches including one of three special characters \*, <, > and =.

| Special character | Meaning                                       | Example of use   |
|-------------------|---|--|
| *                 | 0 or more characters<br>any word [with space] | <i>new*</i> [ <i>new, news, newly, newspaper...</i> ]<br><i>new *</i> [ <i>new car, New York, new ideas...</i> ] |
| >                 | larger than                                   |  |
| <                 | smaller than                                  |  |
| =                 | equals [combined with < and >]                |  |

3. Smart searches are searches predefined in the tool to offer users easy access to complex searches; smart searches are unique to #LancsBox. These searches are used for searching for word classes (NOUN, VERB etc.), complex grammatical patterns (PASSIVE, SPLIT INFINITIVE etc.) and semantic categories (PLACE ADVERB).

Smart searches are defined specifically for a particular language inside the tool. Currently, a group of features is pre-defined in the resources folder: `resources\languages\[name of language]\Searches.txt`. The user can edit this file by adding or deleting items.

The following smart searches are available for English:

|                     |                    |
|---------------------|--------------------|
| !                   | DETERMINER         |
| ,                   | DO                 |
| .                   | DOWNTONER          |
| ?                   | EXISTENTIAL THERE  |
| ADJECTIVE           | GERUND             |
| ADVERB              | HAVE               |
| BE                  | INFINITIVE         |
| BOOSTER             | HYPHENATED WORD    |
| COLLECTIVE NOUN     | INDEFINITE PRONOUN |
| COMPARATIVE         | INFINITIVE         |
| COMPLEX NOUN PHRASE | INTERJECTION       |
| CONDITIONAL         | LINKING ADVERB     |
| CONNECTOR           | LONG WORD          |
| CONTRACTION         | MODAL              |
| DEGREE ADVERB       | NEGATION           |

NOMINALIZATION  
 NOUN  
 NUMBER  
 PARTICLE  
 PASSIVE  
 PAST TENSE  
 PAST PARTICIPLE  
 PERFECT INFINITIVE  
 PHRASAL VERB  
 PLACE ADVERB  
 PREPOSITIONAL PHRASE

PRESENT TENSE  
 PRONOUN  
 PROPER NOUN  
 REFLEXIVE PRONOUN  
 REPETITION  
 SHORT WORD  
 SPLIT INFINITIVE  
 SUPERLATIVE  
 SWEARWORDS  
 TIME ADVERB  
 VERB

4. Regex searches are advanced searches that allow to search for any combination of characters. Any expression enclosed in forward slashes (/) is interpreted as regular expression. #LancsBox supports perl-compatible regular expressions.

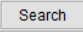
| Regex            | Explanation   | Regex         | Explanation                                     |
|------------------|---|---------------|---|
| <b>Word</b>      | A string of characters (case sensitive)   | <b>a{3}</b>   | Exactly 3 of a                                  |
| <b>/word/i</b>   | A string of characters (case insensitive)   | <b>a{3,}</b>  | 3 or more of a                                  |
| <b>/word\./p</b> | Punctuation search: A string of characters followed by full stop (case sensitive) | <b>a{3,6}</b> | Between 3 and 6 of a                            |
| <b>[abc]</b>     | A single character either a, b or c.  | <b>\d</b>     | Any digit                                       |
| <b>[^abc]</b>    | Any single character except: a, b, or c   | <b>\D</b>     | Any non-digit                                   |
| <b>[a-z]</b>     | Any single character in the range a-z   | <b>\w</b>     | Any word character (letter, number, underscore) |
| <b>[a-zA-Z]</b>  | Any single character in the range a-z or A-Z                                      | <b>\W</b>     | Any non-word character                          |
| <b>[0-9]</b>     | A single number in the range 0-9  |               |   |
| <b>.</b>         | Any single character  |               |   |
| <b>(a b)</b>     | a or b  |               |   |
| <b>a?</b>        | Zero or one of a  |               |   |
| <b>a*</b>        | Zero or more of a   |               |   |
| <b>a+</b>        | One or more of a  |               |   |

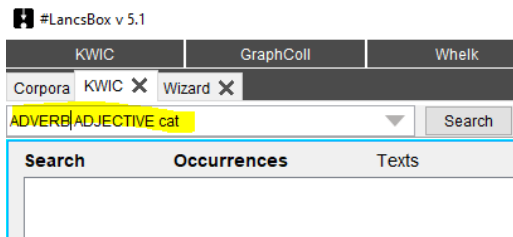
5. Batch searches allow to search for multiple search terms recursively and saving the results automatically; #LancsBox supports both simple and complex batch searches. Batch searches can be used in KWIC, GraphColl and Whelk modules when the corpora are tagged. Here is how batch searches work.
- a) Click on the down arrow in the search box to activate Advanced search options. The last option is a batch search. Click on 'Batch'.

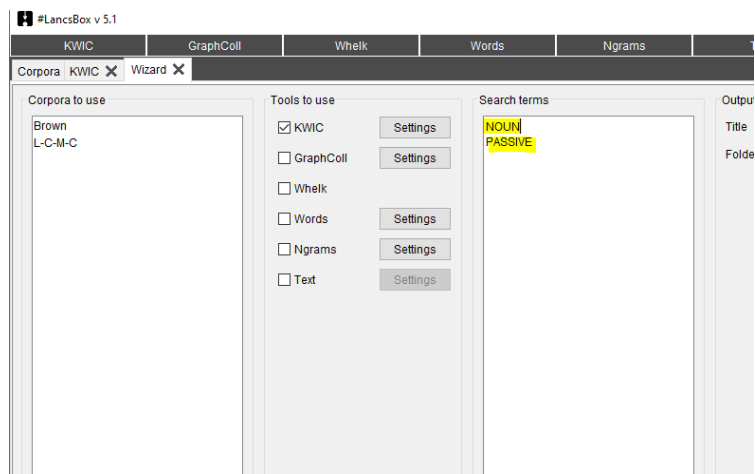


- b) Navigate to and load a text file with the appropriate search terms, one per line. Simple search terms include a list of word forms to be searched; complex search terms are defined via a combination of criteria such as word form, pos tag, headword etc... Consecutive criteria need to be present on the same line separated by tab (\t) in the following order: label – wordform – headword – pos – user tag. This is best achieved by creating the file with advanced batch search terms in Excel or Calc. Examples of simple and complex searches can be seen below.

| Simple batch search: each search term on a separate line | Complex batch search: label – wordform – headword – pos – user tag (tab separated)  |                 |    |       |       |   |   |   |                 |                 |  |  |  |   |          |  |    |  |  |   |     |  |  |    |  |   |      |  |  |  |       |   |              |              |    |       |       |
|--|---|-----------------|----|-------|-------|---|---|---|-----------------|-----------------|--|--|--|---|----------|--|----|--|--|---|-----|--|--|----|--|---|------|--|--|--|-------|---|--------------|--------------|----|-------|-------|
| my   |   |                 |    |       |       |   |   |   |                 |                 |  |  |  |   |          |  |    |  |  |   |     |  |  |    |  |   |      |  |  |  |       |   |              |              |    |       |       |
| cat  |   |                 |    |       |       |   |   |   |                 |                 |  |  |  |   |          |  |    |  |  |   |     |  |  |    |  |   |      |  |  |  |       |   |              |              |    |       |       |
| go   |   |                 |    |       |       |   |   |   |                 |                 |  |  |  |   |          |  |    |  |  |   |     |  |  |    |  |   |      |  |  |  |       |   |              |              |    |       |       |
| went   |   |                 |    |       |       |   |   |   |                 |                 |  |  |  |   |          |  |    |  |  |   |     |  |  |    |  |   |      |  |  |  |       |   |              |              |    |       |       |
|  | <table border="1"> <thead> <tr> <th></th> <th>A</th> <th>B</th> <th>C</th> <th>D</th> <th>E</th> </tr> </thead> <tbody> <tr> <td>1</td> <td>multiple-simple</td> <td>/cat dog mouse/</td> <td></td> <td></td> <td></td> </tr> <tr> <td>2</td> <td>headword</td> <td></td> <td>go</td> <td></td> <td></td> </tr> <tr> <td>3</td> <td>pos</td> <td></td> <td></td> <td>N*</td> <td></td> </tr> <tr> <td>4</td> <td>user</td> <td></td> <td></td> <td></td> <td>Mytag</td> </tr> <tr> <td>5</td> <td>combibnation</td> <td>/going went/</td> <td>go</td> <td>/V.*/</td> <td>Mytag</td> </tr> </tbody> </table> |                 | A  | B     | C     | D | E | 1 | multiple-simple | /cat dog mouse/ |  |  |  | 2 | headword |  | go |  |  | 3 | pos |  |  | N* |  | 4 | user |  |  |  | Mytag | 5 | combibnation | /going went/ | go | /V.*/ | Mytag |
|  | A   | B               | C  | D     | E     |   |   |   |                 |                 |  |  |  |   |          |  |    |  |  |   |     |  |  |    |  |   |      |  |  |  |       |   |              |              |    |       |       |
| 1  | multiple-simple   | /cat dog mouse/ |    |       |       |   |   |   |                 |                 |  |  |  |   |          |  |    |  |  |   |     |  |  |    |  |   |      |  |  |  |       |   |              |              |    |       |       |
| 2  | headword  |                 | go |       |       |   |   |   |                 |                 |  |  |  |   |          |  |    |  |  |   |     |  |  |    |  |   |      |  |  |  |       |   |              |              |    |       |       |
| 3  | pos   |                 |    | N*    |       |   |   |   |                 |                 |  |  |  |   |          |  |    |  |  |   |     |  |  |    |  |   |      |  |  |  |       |   |              |              |    |       |       |
| 4  | user  |                 |    |       | Mytag |   |   |   |                 |                 |  |  |  |   |          |  |    |  |  |   |     |  |  |    |  |   |      |  |  |  |       |   |              |              |    |       |       |
| 5  | combibnation  | /going went/    | go | /V.*/ | Mytag |   |   |   |                 |                 |  |  |  |   |          |  |    |  |  |   |     |  |  |    |  |   |      |  |  |  |       |   |              |              |    |       |       |

- c) Once the file with search terms is loaded, click on the ‘Search’ button (  ) and navigate to the location where the results will be saved.
6. Complex searches. From #LancsBox v. 5.1, extended search functionality is available. #LancsBox automatically identifies different search conventions and performs the desired search automatically. The user just needs to type the pattern in the main search box (KWIC, Whelk, GraphColl) or Wizard search box.





#LancsBox also automatically corrects errors in CQL (Corpus Query Language) – one of the conventions explained below.

For example, the erroneous search term

[headvor="cat

is interpreted correctly as [headword="cat"]

The following conventions are available in complex searches.

- a) Multiple smart searches can be used in the same query; smart searches can also be combined with simple searches.
  - 1) ADVERB ADJECTIVE
  - 2) PRONOUN PASSIVE
  - 3) ADVERB ADJECTIVE NOUN was
  
- b) The OR operator can be used in simple searches to indicate alternatives; it can be combined with parentheses to indicate which words belong together as shown in 3) and 4)
  - 1) cat OR dog
  - 2) car OR dog OR mouse
  - 3) my (cat OR dog)
  - 4) (my cat) OR (my dog)

Note: It is not possible to use the OR operator for combining expressions of different length – with a different number of words, e.g. (my cat) OR dog

c) The NOT operator can be used in simple searches to negate a search term (meaning 'anything but X'); it can be combined with parentheses to indicate which words belong together as shown in 3), 4) and 5).

- 1) NOT my
- 2) NOT my friend
- 3) NOT (my friend)
- 4) NOT (a good) idea
- 5) NOT (a good or bad) idea NOT me

d) #LancsBox also supports CQL (Corpus Query Language). It can be used for defining complex searches at different levels of annotation (1-4) or their combinations. All queries in CQL inside double quotes are interpreted as case insensitive regular expressions; for case sensitivity double equals sign (==) is required, e.g. 5).

CQL allows searching at the following levels of annotation: i) word, ii) headword (hw, lemma), iii) pos and iv) tag. While i)-iii) are supplied automatically for languages with full grammatical support, iv) represents an optional level of a user-defined tag. For example, a single item can be defined in CQL as

```
[word="goes" & headword="go" & pos="V.* "]
```

This is interpreted as a form of the word *goes* with the headword *go* and part-of speech tag *V.\** (verb). Note that the ampersand (&) is used to separate different levels of annotation inside square brackets. If a level of annotation is not specified, no restriction is applied at that level.

In CQL, square brackets [] separate slots in a phrase. Thus, for instance, the following CQL expression

```
[pos="VB.*"] [[]{0,3} [pos="V.N"]]
```

is interpreted as a verb to be (*VB.\**) followed by between 0 and 3 words without any restriction (*[[]{0,3}*) and followed by the past participle (*V.N*).

- 1) [word="cat"]
- 2) [headword="go"]
- 3) [pos="V.\*"]
- 4) [tag="XX"]
- 5) [word=="Cat"]

- 6) [word="go" & headword="go" & pos="N.\* "]
- 7) [headword="go" & pos="V.\*"] [word="to"]
- 8) [headword="very" & pos="R.\*"]{2} [pos="J.\*"]



## 12 Statistics in #LancsBox

#LancsBox uses statistics for calculating measures of i) frequency, ii) dispersion, iii) keywords and iv) collocation. The equations of these measures can be reviewed and modified on the 'Stats' tab, which is called by clicking on the  $\Sigma$  button.

---

### 12.1 Frequency measures

1. absolute frequency =  $o_{11}$
2. relative frequency =  $(o_{11}/r_1) \times 10,000$

### 12.2 Dispersion measures

1.  $CV = SD/\text{mean}$
2.  $SD = \sqrt{\frac{\sum(x-\text{mean})^2}{n}}$
3. Range = no of files where the search term occurs at least once
4.  $\text{Range \%} = \frac{\text{Range}}{\text{number of files}} \times 100$
5.  $D = 1 - \frac{CV}{\sqrt{\text{number of files}-1}}$
6.  $DP = \frac{\text{Sum of absolute values of (observed-expected proportions)}}{2}$

### 12.3 Keyword measures

1. simple maths parameter =  $\frac{\text{relative frequency of } w \text{ in } C + k}{\text{relative frequency of } w \text{ in } R + k}$
2.  $\log \text{ likelihood}_{\text{short}} = 2 \times \left( O_{11} \times \log \frac{O_{11}}{E_{11}} + O_{21} \times \log \frac{O_{21}}{E_{21}} \right)$
3.  $\% \text{ DIFF} = \frac{(\text{relative freq. in } C - \text{relative freq. in } R) \times 100}{\text{relative freq. in } R}$
4.  $\text{Log Ratio} = \log_2 \left( \frac{\text{relative freq. in } C}{\text{relative freq. in } R} \right)$
5.  $\text{Cohen's } d = \frac{\text{Mean}_{\text{in } C} - \text{Mean}_{\text{in } R}}{\text{pooled SD}}$

## 12.4 Collocation measures

| ID | Statistic               | Equation  | ID | Statistic                | Equation  |
|----|-------------------------|---|----|--------------------------|---|
| 1  | Freq. of co-occurrence  | $O_{11}$  | 8  | T-score                  | $\frac{O_{11} - E_{11}}{\sqrt{O_{11}}}$   |
| 2  | MU                      | $\frac{O_{11}}{E_{11}}$   | 9  | DICE                     | $\frac{2 \times O_{11}}{R_1 + C_1}$   |
| 3  | MI (Mutual information) | $\log_2 \frac{O_{11}}{E_{11}}$  | 10 | LOG DICE                 | $14 + \log_2 \frac{2 \times O_{11}}{R_1 + C_1}$                                       |
| 4  | MI2                     | $\log_2 \frac{O_{11}^2}{E_{11}}$  | 11 | LOG RATIO                | $\log_2 \frac{O_{11} \times R_2}{O_{21} \times R_1}$                                  |
| 5  | MI3                     | $\log_2 \frac{O_{11}^3}{E_{11}}$  | 12 | MS (Minimum sensitivity) | $\min\left(\frac{O_{11}}{C_1}, \frac{O_{11}}{R_1}\right)$                             |
| 6  | LL (Log likelihood)     | $2 \times \left( O_{11} \times \log \frac{O_{11}}{E_{11}} + O_{21} \times \log \frac{O_{21}}{E_{21}} + O_{12} \times \log \frac{O_{12}}{E_{12}} + O_{22} \times \log \frac{O_{22}}{E_{22}} \right)$ | 13 | DELTA P                  | $\frac{O_{11}}{R_1} - \frac{O_{21}}{R_2}; \frac{O_{11}}{C_1} - \frac{O_{12}}{C_2}$    |
| 7  | Z-score <sub>1</sub>    | $\frac{O_{11} - E_{11}}{\sqrt{E_{11}}}$   | 14 | Cohen's <i>d</i>         | $\frac{\text{Mean}_{in\ window} - \text{Mean}_{outside\ window}}{\text{pooled\ } SD}$ |

## 13 Glossary

**Absolute (or raw) frequency** – The simple frequency with which a search term occurs in a corpus or its part(s); a number of hits of a search term in a corpus.

**Batch search** – A batch search enables searching for multiple search terms recursively and saving the results automatically; #LancsBox supports both simple and complex (i.e. defined via a combination of criteria such as wordform, pos tag, headword etc.) searches.

**Colligation** – Systematic co-occurrence of grammatical categories (e.g. POS tags) in text identified statistically.

**Collocate** – A word that systematically occurs with the node (word or phrase of interest, search term).

**Collocation** – Systematic co-occurrence of words in text identified statistically.

**Collocation graph** is a visual display of the association between a node and its collocates. See GraphColl.

**Collocation network** is a visual display of complex associations (collocations) in language and discourse. It consists of multiple inter-connected collocation graphs. See GraphColl.

**Concordance line** – A single line in the KWIC display representing a node (search term) with the words before and after it (the right and left context).

**Concordance** is a typical form of display of examples of language use found in a corpus with the node (search term) centred in the middle and several words of context displayed left and right of the node. Concordance is sometimes also called a 'KWIC (display)'.

**Corpus** (pl. corpora) – A collection of language data that can be searched by a computer.

**Dispersion** – is the spread of values of a variable (e.g. relative frequencies of a search term) in a dataset (corpus). Dispersion is measured statistically using metrics such as standard deviation (*SD*), coefficient of variation (*CV*), range, Juillard's *D*, *DP* etc. See Words.

**Frequency** – The number of times a search term occurs in the corpus. A distinction is made between absolute (absolute number of hits) and relative frequency (proportional frequency per *X* number of tokens).

**Frequency distribution** – frequency distribution provides information about the frequencies of a word or phrase in different parts of the corpus. See Whelk.

**GraphColl** is a module in #LancsBox, which identifies collocations and builds collocation networks on the fly.

**Import** – In #LancsBox, processing of corpus data and making it available to all modules in the package.

**KWIC** is an abbreviation for 'keyword in context'. This is a typical form of display of examples found in a corpus with the node (word or phrase of interest) centred in the middle and several words of context displayed left and right of the node. KWIC is sometimes also called a 'concordance'. KWIC is also the name of a module in #LancsBox.

**Left context** – The words preceding a particular search term (node). Individual positions in the left-context are referred to as L1 (position immediately preceding), L2, L3 etc.

**Lemma** – All inflected forms belonging to one stem; in #LancsBox by default, a combination of a headword and a grammatical category (e.g. go + VERB). For example, a lemma 'go' includes the following word forms (types): 'go', 'goes', 'went', 'going' and 'gone'.

**Lexical bundle** – an n-gram with certain frequency and distributional (dispersion) properties, e.g. relative freq. 10 per million and range > 5.

**Loaded** – In #LancsBox, when a corpus is loaded it is available to be analysed. To re-load a corpus, double-left-click on the name of the corpus.

**Module** – A specific tool within #LancsBox offering particular analytical functionalities. #LancsBox includes five different modules: KWIC, Whelk, GraphColl, Words and Text.

**N-gram** – a sequence of n types, lemmas, POS from a text or corpus.

**Node** – The word, phrase or grammatical structure of interest. See Search term.

**Part of speech (POS)** – A grammatical category, a word class. Part-of-speech is usually assigned automatically using a process called part-of-speech tagging (see below). #LancsBox includes TreeTagger, which performs part-of-speech tagging for a range of languages.

**Part-of-speech tagging (POS tagging)** – A process of adding information about the grammatical category of each word in a text or corpus. For example, the following sentence was POS-tagged: Automatically\_RB annotates\_VBZ data\_NNS for\_IN part-of-speech\_NN.

**P-frame (also skip gram)** – an n-gram that allows for variability at one or more positions such as it would be \* to.

**Regular expressions (regex)** – A special meta-language that allows advanced users to search for any combination of strings. In #LancsBox, regex searches are enclosed in forward slashes e.g. /. \*ions?/

**Relative (or normalized) frequency (RF)** is calculated as the proportion of the absolute frequency of a word we are interested in divided by the total number of words (tokens) in the corpus. This number is usually multiplied by an appropriate basis for normalization (e.g. 10,000).

**Right context** – The words following a particular search term (node). Individual positions in the right-context are referred to as R1 (position immediately following), R2, R3 etc.

**Split screen** – A comparison option in #LancsBox where the screen can be split into two panels; each panel can display a different type of analysis. #LancsBox allows second panel to be opened and minimised via left-clicking on three small triangles (▲▲▲/▼▼▼).

**Tab** – A further ‘page’ that can be opened in #LancsBox to run multiple analytical procedures simultaneously. Each module in #LancsBox can run on an unlimited number of tabs.

**Tagging** – The process of adding linguistic information to the words in a text or corpus, automatically or semi-automatically. See Part-of-speech tagging.

**Text** – A basic unit of a corpus; a corpus is a collection multiple texts. Text is also the name of a module in #LancsBox that displays and searches texts in corpora.

**Threshold** – Setting options in GraphColl and Words to display only relevant collocates or keywords respectively.

**Token** is a single occurrence of a word form in a text or corpus.

**TreeTagger** is a part-of-speech tagger developed by Helmut Schmid, which performs part-of-speech tagging for a range of languages.

**Type** is a unique word form in a text or corpus.

**Whelk** is a module in #LancsBox which provides information about how the search term is distributed across corpus files.

**Words** is a module in #LancsBox which allows in-depth analysis of frequencies of types, lemmas and POS categories as well as comparison of corpora using the keywords technique.

## 14 Messages.Properties

How to configure #LancsBox for advanced users

The Messages.properties file lets you customise #LancsBox. The things you change in here will change how the program operates and looks.

### 15 Making Changes

To change a setting in Messages.Properties: First, look for the setting you want. This will look something like: an.interesting.setting = value. We will call the first part (before the =) the key and the second part (after the =) the value.

Changing the value will change the setting. Each of the values has it's own type. This just means that when you change a colour it should be for another colour, not for a word. We will now introduce the value types used in Messages.properties.

- Path – This tells LancsBox where on your computer to look for something. This could be where to find an Icon, or where to find other information that LancsBox relies on. It will look like this: **/resources/path/to/a/file**. Any path starting with /resources refers to something within the resources folder.
- Integer (This just means a whole number)
  - As a number – Some settings want an actual number, like the default span for KWIC searches.
  - As a selector – Sometimes we use an integer to pick between several options. The options will be described in comments so you will know what your choices are.
  - As a true/false – Some settings turn things on or off in #LancsBox. If you want the setting on then the value should be **1**. **0** means turn it off.
- Colour – You can change many of the colours used in #LancsBox. These follow a particular format. We recommend using an online colour picker to find out what the value should be. It will begin with a **#**.
- Regular expression – These are like the regular expressions that you use within #LancsBox with two changes. Firstly, the //s are omitted, Secondly the options (like i) are in the Java format. These are in round brackets, have a question mark and precede the expression itself. like: **(?i) regular expression**
- Literal text – Whatever you type as the value here will be used directly. However, please note that the UI for #LancsBox uses an English font, this may limit the utility of these types of settings.
- Number Format – Java has it's own way of defining number formats. Changing these in Messages.Properties will change how #LancsBox displays numbers and can be useful in altering the number of decimal places in tables.

| Message.properties file   | Explanation   |
|---|---|
| <b># This is the Messages.properties file</b><br><b># tagger</b><br>tagger.dir = resources/tagger<br>tagger.langs = resources/tagger/models | <b>Tagger</b> .dir and .langs both have path values and refer to parts of a Tree-Tagger installation. The root directory of the Tree-Tagger is defined by.dir and .langs is a folder containing the language specific files (.par extension). |

### # database

```
# 0 - RAM database, basic persistence
# 1 - QuestDB (recommended) (64bit)
```

```
database.use = 1
database.dir = resources/corpora
database.cache.size = 2000000
```

### #download locations

```
downloads.corpora      = resources/downloads/corpora
downloads.wordlists    = resources/downloads/wordlists
```

### # language settings

```
langs.dir = resources/languages
```

### # default tokenizer settings

```
defaults.punctuation
    =.,:;!\\u00e3\\u0080\\u0082\\u00ef\\u00bc\\u008c\\u00ef\\u00bc\\u009b\\u00ef
\\u00bc\\u009a\\u00ef\\u00bc\\u009f\\u00ef\\u00bc\\u0081\\u00e2\\u0080\\u009a\\u00c
2\\u00bf\\u00c2\\u00a1\\u00e2\\u0080\\u00a6'\"\\u00e2\\u0080\\u0098\\u00e2\\u0080\\
u0099`\\u00e2\\u0080\\u009c\\u00e2\\u0080\\u009d\\u00e2\\u0080\\u009e(<=>[{}\\u0
0e2\\u0080\\u00b9\\u00e2\\u0080\\u00ba\\u00e3\\u0080\\u008a\\u00e3\\u0080\\u008b-
\\u00e2\\u0080\\u0093\\u00e2\\u0080\\u0094\\u00e4\\u00b8\\u0080*
defaults.segmentation      =\\t\\n \\r
defaults.sentence_boundary =(?s).*[\\.|!|\\?|.| | ? | ! ].*
```

### # Script directory

```
stats.dir      = resources/stats
```

### Database

.dir lets you change where corpora will be stored and is a path value. The number of tokens held in RAM can be limited or expanded by changing .size, which is an integer value. A number of different databases can be used within #LancsBox and .use lets you change which one is being used. The integer value can be one of a number of options which are part of the comment above the setting. Note that you can't load a corpus using a database if it wasn't created by the same database.

### Download locations

When you download wordlists and corpora they are saved in the downloads folder in resources prior to being imported into a corpus. Changing the path values of these settings lets you change where they get saved.

### Language settings

The language-specific settings are stored in the languages folder in resources. Changing the .dir setting lets you change this location.

### Default tokenizer settings

The tokenizer can be configured on the corpora panel in #LancsBox. The default values that appear on those boxes come from here. Making the change in Messages.Properties means that you only have to make the change once. Please note that .punctuation and .segmentation are literal values (which include an additional escape character - \) whereas .sentence\_boundary is a regular expression. The sentence boundary is used for calculating average sentence length and similar metrics.

```
stats.threshold = resources/groovy/default_threshold.groovy
stats.dir.collocate = resources/stats/collocate
stats.dir.keyword.frequency = resources/stats/keyword/frequency
stats.dir.keyword.dispersion = resources/stats/keyword/dispersion
stats.dir.keyword.statistic = resources/stats/keyword/statistic
shaders.dir = resources/shaders
```

### # Tool logo

```
icons.logo = resources/images/logo.png
```

### # Fonts

```
fonts.all.size = 12
fonts.table.size = 12
fonts.2d.scale = 0.25
fonts.graph.size = 84
fonts.graph.size.scale = 0.125
fonts.keyword.size = 84
fonts.keyword.size.scale = 0.4
```

# Select the fonts to use when there is no custom font installed.

# The custom font is the first .ttf file found in the resources/fonts folder.

# Java logical font options:

# 1 - Dialog

# 2 - DialogInput

# 3 - Monospaced

# 4 - Serif

# 5 - SansSerif

```
fonts.default.ui = 5
```

```
fonts.default.3d = 1
```

### Script directory

#LancsBox calculates statistics using a number of external scripts, which you can also edit. Each of the groups of scrips lives in a different folder. These path values let you change where they are read from.

### Tool logo

The path value can be changed to change the #LancsBox logo for another image.

### Fonts

The font sizes used in #LancsBox can be altered here. The graph and keyword fonts have large sizes which use the .scale values to shrink them. This gives high resolution text at a good size. To increase the beauty of text in graphs and words tools make the appropriate .size values larger and the .scale values smaller.

Java uses what it calls logical fonts. The #LancsBox UI uses different logical fonts by default. You can change the default font options by changin the .ui and .3d integer values. A comment precedes the settings to inform you of the available options. The data font can be overwritten from this by placing a single .ttf file in the resouces/fonts folder.



## # Misc

```
window.size.width = 1024
window.size.height = 768
slider.lock = 0
tokeniser.allowRtoL = 0
display.default.RtoL = 0
numbers.format.integer = ###,###,###,###,###
numbers.format.real = #####0.000000
numbers.format.real_short = #####0.00
```

## # General program colours

```
colours.bar = #4B4B4B
colours.highlight = #00A4FF
colours.text_highlight = #ff6600
colours.text = #4B4B4B
colours.advanced_arrow = #B3B1B0
```

## # General UI paths and settings

```
icons.frame = resources/images/icon.png
icons.tabs.attach = resources/images/pin1.png
icons.tabs.close = resources/images/cross.png
icons.generic.right_arrow = resources/images/right-arrow.gif
icons.corpora = resources/images/corpora.png
icons.save = resources/images/save.png
icons.stats = resources/images/stats.png
icons.about = resources/images/about.png
icons.help = resources/images/help.png
icons.kwic = resources/images/kwic.png
icons.graph = resources/images/graph.png
icons.compare = resources/images/compare.png
```

## Misc

.lock allows you to either lock or unlock (true/false) the slider which lets you resize the tables in graph and words.

Most right to left corpora are actually in left to right format in the files, but are displayed in reverse. If the actual data is stored as right to left (very unlikely) then .allowRtoL can be enabled and a new checkbox will appear in import options on the corpora tab. In the much more likely event that the data is left to right but should be displayed right to left then the default display direction of #LancsBox using the .RtoL setting. Both of these settings are also true/false values. The format of numbers in tables can be changed using the .integer, .real and .real\_short settings.

## General program colours

Some of the more widely used colours in #LancsBox can be changed here.

## General UI paths and settings

The path values can be changed to load custom icons in the #LancsBox UI. The default message on the status bar can be changed by altering the .welcome\_message value, which is a literal string.

icons.compare.disabled = resources/images/compareDisabled.png  
statusbar.welcome\_message = Welcome to #LancsBox

### # table icons

icons.sort.ascending = resources/images/upArrow.png  
icons.sort.descending = resources/images/downArrow.png  
icons.sort.ascending.filtered = resources/images/upArrowSquare.png  
icons.sort.descending.filtered = resources/images/downArrowSquare.png  
icons.sort.filter = resources/images/square.png  
icons.sort.random = resources/images/random.png

### # The tooltips for various buttons

buttons.tooltip.corpora = Corpora  
buttons.tooltip.save = Save  
buttons.tooltip.graph = Collocation graphs and networks tool  
buttons.tooltip.kwic = <html>Concordance tool</html>  
buttons.tooltip.whelk = Dispersion tool  
buttons.tooltip.keywords = Wordlists and keywords tool  
buttons.tooltip.ngram = N-Gram tool  
buttons.tooltip.text = Text tool  
buttons.tooltip.help = Help  
buttons.tooltip.stats = Statistics  
buttons.tooltip.about = About  
buttons.popup.close = Apply

### # Generic button labels, reused throughout

buttons.generic.browse = Load data  
buttons.generic.delete = Delete  
buttons.generic.clear = Clear  
buttons.generic.run = Run  
buttons.generic.new = New  
buttons.generic.load = Load

### Table Icons

Changing these path values lets you change the icons that appear in #LancsBox tables.

### The tooltips for various buttons

The tooltips are customisable for the main buttons. These are the tool and status bar buttons that you first see when loading #LancsBox.

### Generic button labels, reused throughout

The text of some buttons in the UI can be changed by altering these literal string values. This includes the apply button on some popups.

```
buttons.generic.save = Save
buttons.generic.close = Close
#buttons.generic.stop = Stop
```

### # Load pane

```
labels.load.prompt_name = Name:
labels.load.corpus_name = Corpus
labels.load.case = Clamp types to lowercase
labels.load.punctuation = Store punctuation
buttons.load.new = Import!
buttons.load.reset = Reset to defaults
icons.load.corpus = resources/images/corpus.png
icons.load.wordlist = resources/images/wordlist.png
```

### # Stats pane text

```
tabs.name.stats = Statistics
labels.stats.name = Name:
buttons.stats.commit = Save
buttons.stats.save = Save as...
buttons.stats.load = Open...
buttons.stats.remove = Remove
buttons.stats.revert = Revert
```

### # n-gram settings

```
defaults.ngrams = 2
```

### # keywords renderer

```
colours.keywords.corpus_name_dark = #000000
colours.keywords.corpus_name_light = #bababa
colours.keywords.text = #000000
colours.keywords.target = #c60db8
colours.keywords.target.max = #c60db8
```

### Load pane

The main corpora pane uses some literal strings that can be changed here.

### Stats pane text

The stats panel uses some literal strings that can be changed here.

### N-gram settings

The n-grams tool defaults to being a bigram tool. This can be changed by changing this integer value.

### Keywords renderer

The words / ngrams tool has a number of colours which can be changed. Those which have a corresponding .max value denote a colour range. The frequency colours will be interpolated using these ranges.

```
colours.keywords.reference      = #2e3131
colours.keywords.reference.max  = #2e3131
colours.keywords.highlight     = #ff6600
colours.keywords.table         = #d1d1d1
colours.scroll                  = #5f5f5f
colours.no_scroll               = #d1d1d1
```

### # Graph pane text

```
buttons.graph.export  = Export
buttons.graph.export.dot = .dot File
buttons.graph.export.img = .png Image
buttons.graph.labels  = Labels
buttons.graph.run     = Search
buttons.graph.layout  = Layout
buttons.graph.kwic    = KWIC
buttons.graph.threshold = Threshold
buttons.graph.stat    = Stat
```

### # Graph Renderer

```
colours.graph.node          =#c60db8
colours.graph.collocate_light =#e6f7f9
colours.graph.collocate_dark  =#000000
colours.graph.highlight     =#ff6600
colours.graph.edge          =#d1d1d1
colours.graph.text          =#000000
colours.graph.shared        =#ff6600
colours.graph.shared_background =#583e82
renderer.screenshot.width    =7680
renderer.screenshot.height   =4320
renderer.default.sphere_size =6
renderer.default.sphere_resolution =50
colours.toggle.free         = #31c831
colours.toggle.hybrid       = #ffc200
```

### Graph pane text

A number of string literals are given here for the GraphColl tool. The string literals can be changed here.

### Graph renderer

The GraphColl tool has colours and colour ranges which can be changed. These are the colour values given here. Additionally the size of screenshots can be changed here (though they also apply to words) by changing the integer values of .width and .height. The number of sides a sphere has (all 3d tools) can be changed by altering the integer value of the .sphere\_resolution setting. This can drastically speed up crowded graphs but not all numbers will work on all computers. The size of graph spheres can also be changed using .sphere\_size. This gives you even greater control than just changing the font size.

```
colours.toggle.positional    = #ff0040
colours.toggle.word_class    = #cd00cd
```

### # Whelk searches

```
whelk.window.span = 100
```

### #KWIC pane colours

```
colours.kwic.node           =#ff6600
colours.kwic.highlight      =#00a4ff
colours.kwic.highlight_not  =#5e626b
```

### # KWIC window size settings

```
kwic.left.min              =3
kwic.left.def               =5
kwic.left.max               =20
kwic.right.min              =3
kwic.right.def              =5
kwic.right.max              =20
```

### # POS group colours

```
colours.group.1            = #0080ff
colours.group.2            = #ff0080
colours.group.3            = #00cd67
colours.group.4            = #ff6500
colours.group.5            = #cd00cd
colours.group.6            = #d5ff00
colours.group.7            = #a6a6a6
colours.group.8            = #00e6e6
colours.group.9            = #ff4dff
colours.group.10           = #006200
```

### Whelk searches

STTR and MATTR searches can be performed in #LancsBox. These use a window size of a number of tokens. This number can be changed by setting the value of the .span setting to a different integer.

### KWIC pane colours

The colours used in the KWIC tool can be changed by altering the colour values of these settings.

### KWIC window size settings

The default span settings for KWIC searches can be set here. The integer values only define the defaults, you can still change them in the program.

### POS group colours

The POS groups / aliases can be defined in the import options. The first ten of them will be assigned these colours when viewing lemma graphs in the word class mode.

### # Wizard settings

```
wizard.debug.enabled=false  
wizard.reportTitle=report  
wizard.xml.truncate=-1  
wizard.enable_filter=false  
wizard.statsurl=http://corpora.lancs.ac.uk/lancsbox/stats_wizard.php  
wizard.results.extraxml.threshold=5000000  
wizard.report.tab.symbol=\uFFEB
```

### # Web corpus download settings

```
webcorpus.download.logfolder = resources/web  
webcorpus.download.minpausems = 500  
webcorpus.download.maxpausems = 2000  
webcorpus.download.querynumberforpause = 10  
webcorpus.download.heuristics.repeatedsize.stopthreshold = 3
```

### Wizard

Various default settings for Wizard including the default report title, url for statistical analysis etc...

### Create corpus

Various default settings connected with the automatic creation of corpora in #LancsBox. Importantly, the users can modify the min and max values for pauses (in milliseconds) that are included when requesting data from servers to prevent automatic blocking of the process.