## The Prime Machine HD Guide
### *I want to …* build my own corpus

If you are studying English or linguistics as an undergraduate or postgraduate student, you may want to create your own corpus for an assignment or research project.  If you are teaching English, you may want to make a corpus of pedagogical texts – the specialist kinds of texts your students use in class or need to be prepared for.  This short guide will explain how to import text to make a new DIY corpus using The Prime Machine HD corpus tool.

Steps to complete:
1. Prepare the texts you want to use to build your corpus.
2. Import the documents.
3. Use the corpus tools to explore different features, compared to a readymade corpus or a second DIY corpus.

What you'll need to get started:
- The Prime Machine HD for Windows, MacOS, iPad, iPhone or Android (available free from https://www.theprimemachine.net/ )
- Your texts as plain text, RTF, Word, PDF, PPT or EPUB.
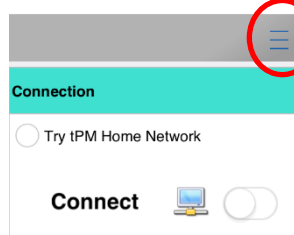- Patience, enthusiasm and an open mind!

## Getting started
The best place to get The Prime Machine HD (tPM) is from an official store.  It is free!



Windows and Android users can also download the App directly from the website: https://www.theprimemachine.net

When you first use tPM, you almost certainly will want to connect to the server to access pre-prepared corpora and resources.  There are two main views for the search screen – Simple Mode and Full Mode.  The Full Mode includes additional tabs and features for corpus research and DIY corpus work.  In this guide, you will need to use Full Mode.



The main 'hamburger' or 蒸笼 menu in the top-right corner allows you to connect and change mode.

## Step 1: Preparing files for a DIY corpus

The files you use to create your corpus could be documents on your device, or texts you download. There are free and commercial tools to help you gather texts from the internet with minimal effort.

The main point of building a corpus is usually to try to represent a specific text variety.

When viewing corpus data in tPM, the filenames and text categories can be seen on the cards display, so it is worth spending some time organizing the files and naming them in a systematic way.
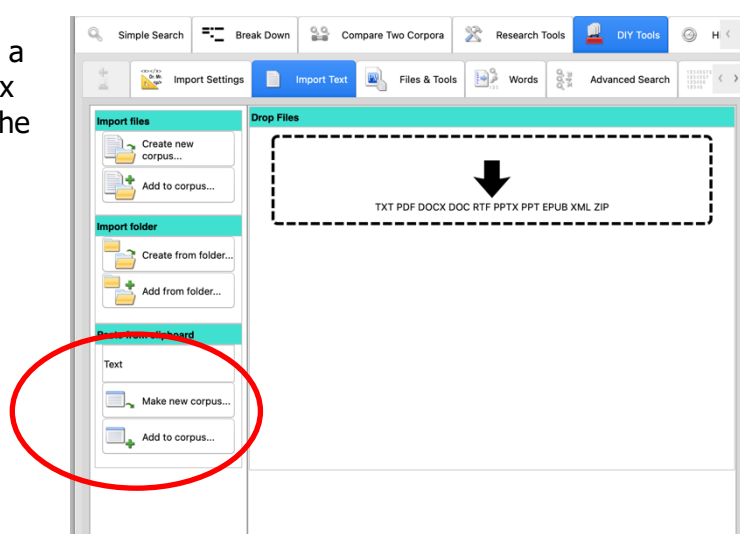
Before you start importing texts (and particularly book chapters), you should think about whether you want to split long files into smaller parts and also think about the filenames and the order of the files.

If you load your corpus as a single text, you will still be able to complete many kinds of analyses, but you won't be able to explore texts separately and it may be harder to notice if hits are limited to particular texts. If you do use separate files, make sure the filenames are neat and consistent. For example, rename book chapters "Chapter 1", "Chapter 2", etc.

## Step 2: Importing the texts

If you are using iPad, iPhone or Android, use your device's file manager to make a zip file of all the chapters first. On desktop platforms, you can load multiple files in one operation.

The "Paste from Clipboard" function is a quick way to make a small DIY corpus. Using the box above the button, you can set the name of the corpus before you paste text in.

## File formats

You can import text from a variety of file formats: PDF, DOCX, DOC, RTF, TXT, PPTX, PPT and EPUB.

However, some PDF files may not be compatible because they contain scanned images of text or have other restrictions.

Your aim should be to load a good selection of texts from suitable sources; you don't need to include everything.
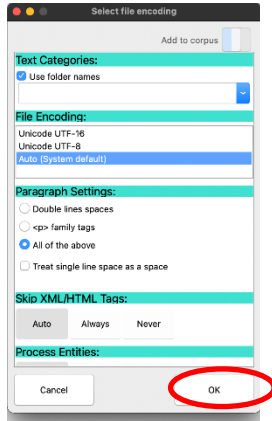
If there is a source you really want to include, but it doesn't load properly in tPM, you can consider using the free desktop app tPMCrafty (also available from The Prime Machine website) or use other tools to convert your PDF to plain text first.

tPMCrafty can help add spaces to the ends of lines, alter the spacing between paragraphs and split a text into smaller parts.



For DIY corpora of more than one text, screenshots of the procedure on all platforms are on pages 3-6. Remember to SAVE THE DIY CORPUS after you have imported all the texts.
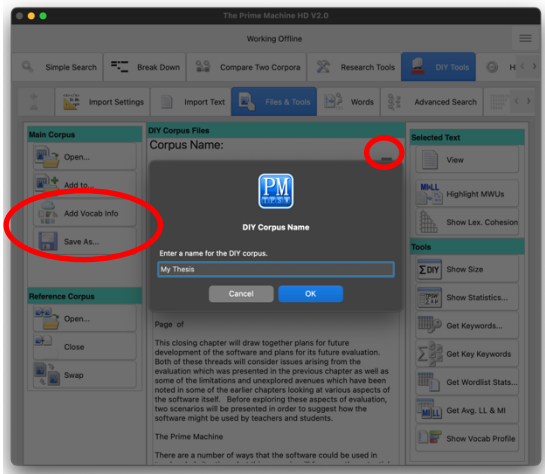
## MacOS

When the App starts, make sure Full Mode is selected and go to the DIY Tools – Import Text tab.
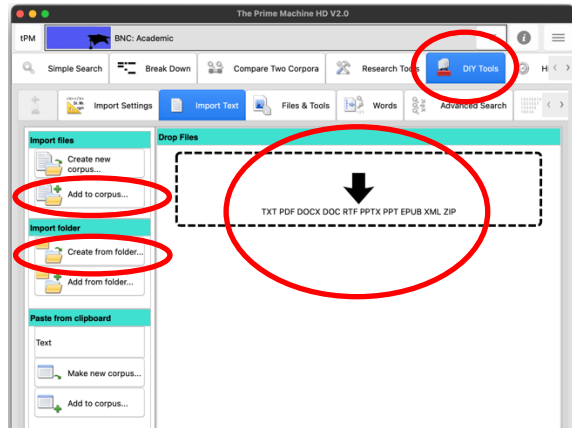
Most default settings will work well.  Alter the file encoding or paragraph settings if working with plain text files and you get unexpected results.
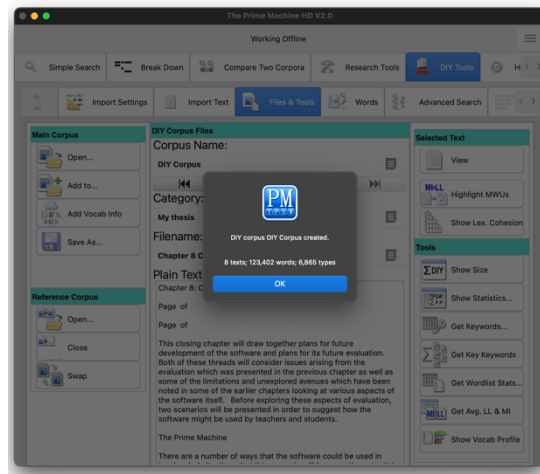
Wait patiently while the text is extracted, the sentences and paragraphs are organised and the words and combinations of words are indexed.
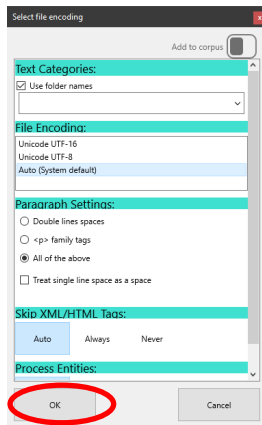
Drag and drop the files into the drop zone; or use the "Create new corpus…" button to select one or more files inside a folder; or choose "Create from folder…" to import an entire folder of files.

When it is finished, you will see the size in texts, words and types.

From the Files & Tools Tab you can rename the corpus and then save it using the "Save as…" button.
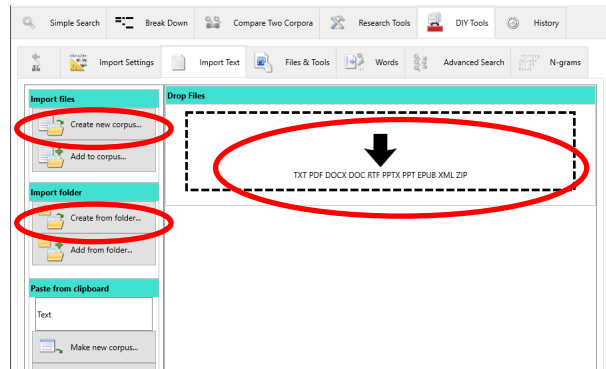
## Windows

When the App starts, make sure Full Mode is selected and go to the DIY Tools – Import Text tab.
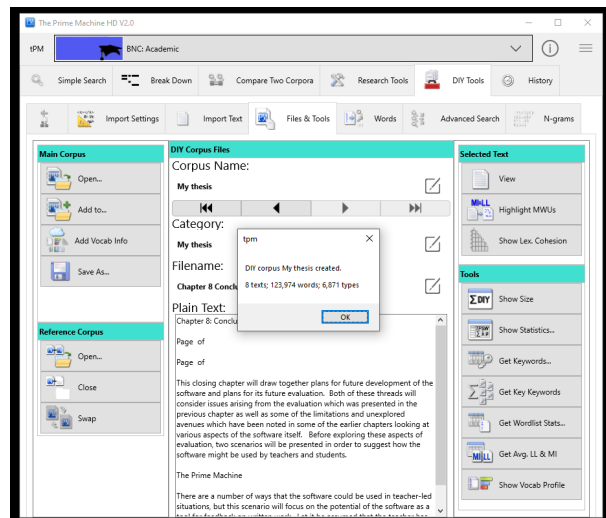
Most default settings will work well. Alter the file encoding or paragraph settings if working with plain text files and you get unexpected results.
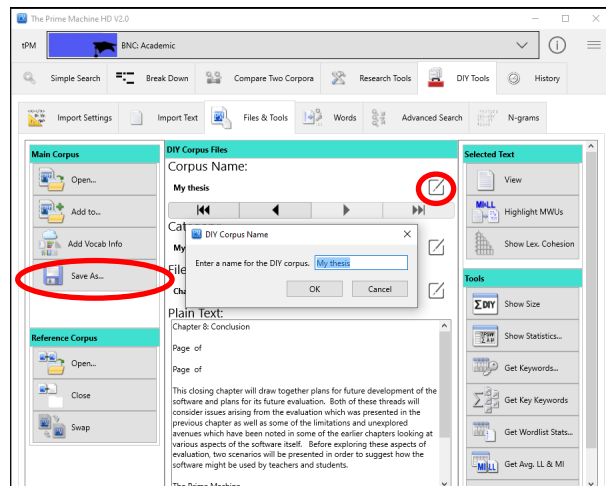
Wait patiently while the text is extracted, the sentences and paragraphs are organised and the words and combinations of words are indexed.

Drag and drop the files into the drop zone; or use the "Create new corpus…" button to select one or more files inside a folder; or choose "Create from folder…" to import an entire folder of files.
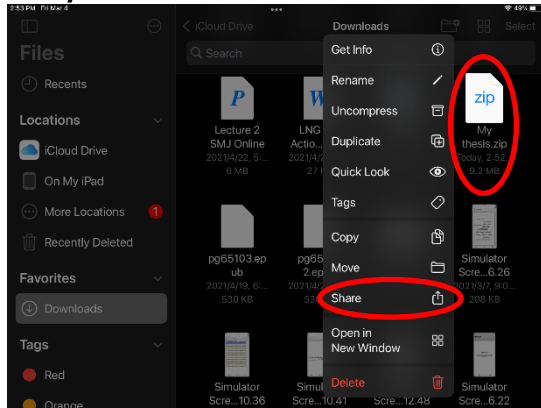
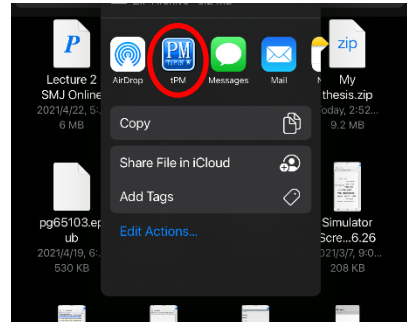When it is finished, you will see the size in texts, words and types.

From the Files & Tools Tab you can rename the corpus and then save it using the "Save as…" button.
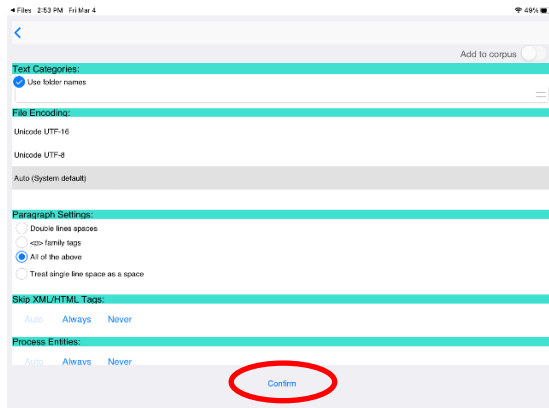
## iPad / iPhone



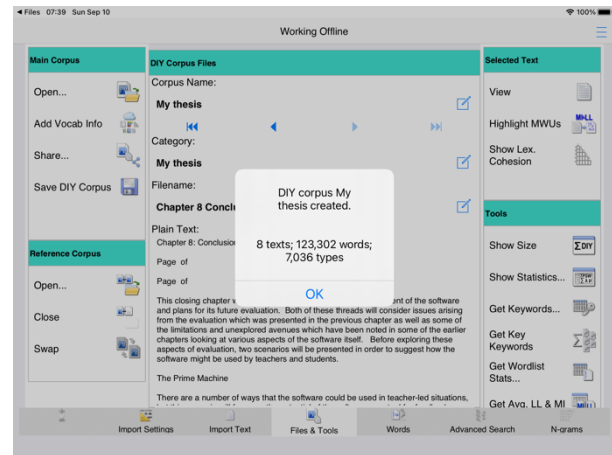Use the "Files" App to find the zip file (or single document) you want to import into tPM.  Long tap on it and then choose "Share".

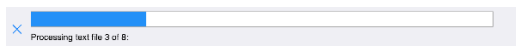**On mobile platforms, tPM only has access to your files when you share them from another app.**



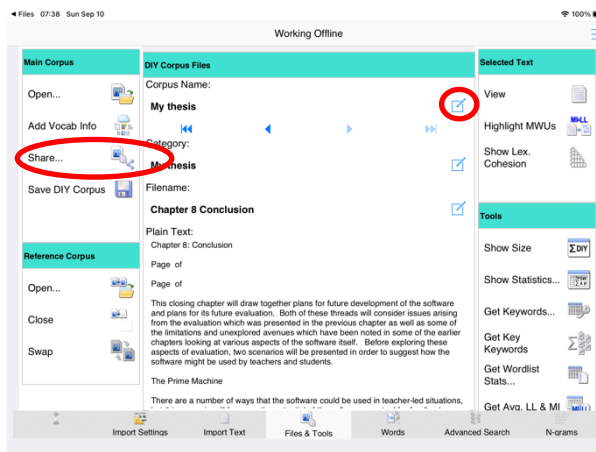Find tPM on the list of apps and tap on it.



Most default settings will work well.  Alter the file encoding or paragraph settings if working with plain text files and you get unexpected results.



Wait patiently while the text is extracted, the sentences and paragraphs are organised and the words and combinations of words are indexed.



When it is finished, you will see the size in words and types.



From the Files & Tools Tab you can rename the corpus and then save it using the "Save" button.

**Android**

Use the "Files" App to find the zip file (or single document) you want to import into tPM. Long tap on it and then choose "More" and "Open in app"

**On mobile platforms, tPM only has access to your files when you share them from another app.**

Find tPM on the list of apps and tap on it.

Most default settings will work well. Alter the file encoding or paragraph settings if working with plain text files and you get unexpected results.
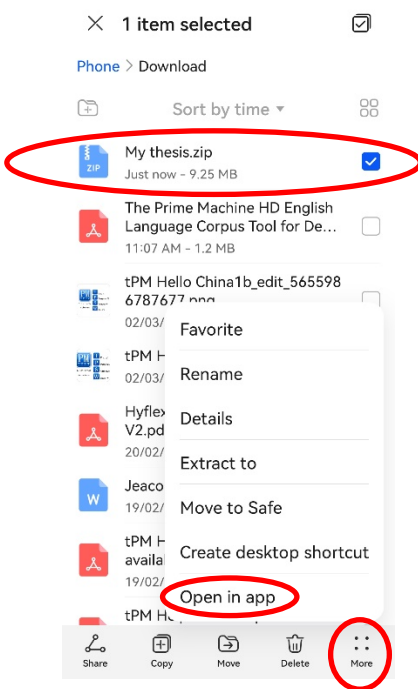
Wait patiently while the text is extracted, the sentences and paragraphs are organised and the words and combinations of words are indexed.
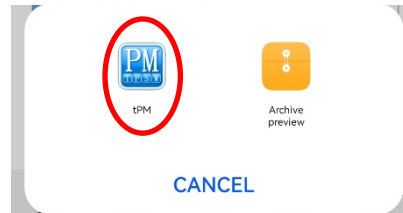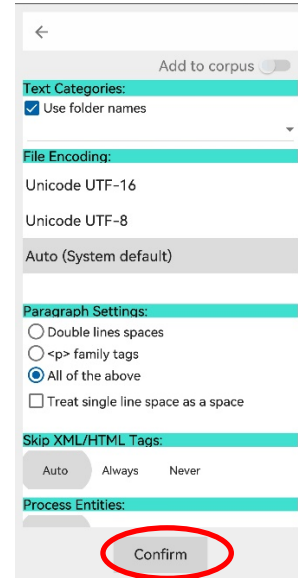
When it is finished, you will see the size in texts, words and types.

From the Files & Tools Tab you can rename the corpus and then save it using the "Save" button.

**Building a corpus out of concordance lines**

It is also possible to build a special kind of DIY corpus called a probe corpus – a corpus made up from the sentences of the concordance lines of one or more queries.

If you keep the Wider Context setting on the Lines display before you build the probe corpus, it will include up to one sentence before and after the sentence containing the node. If you change the settings on the Line display to Sentence Only before you build the probe corpus, it will only include the sentences containing the node.

If you filter the results before you build the probe corpus, it will only include the selected lines.

You can generate complicated concordance line searches to include several node words, or you can add one set of concordance lines to an existing probe corpus one set at a time.





On narrower screens (like this Android phone), the pop-up menu looks like this.

## Step 3: Analysing your corpus

Now you have your texts loaded as DIY corpora, you can explore them using data driven methods as well as specific searches. For many of these operations, a reference corpus is used as a baseline.  You should select a suitable readymade corpus from the menu at the top of the screen, or for operations which can be completed using a second DIY corpus as a reference corpus, you should load that first using the buttons on the Files & Tools tab.

| Function | Readymade | 2nd DIY |
|---|---|---|
| **Files & Tools Tab** | | |
| Show File List | ✓ | ✓ |
| Key Words | ✓ | ✓ |
| Key Keywords | ✓ | ✗ |
| Highlight MWUs | ✓ | ✓ |
| Wordlist Statistics | ✓ | ✓ |
| Average LL & MI | ✓ | ✗ |
| Vocab. Profile | ✓ | ✗ |
| **Words Tab** | | |
| Concordances | ✓ | ✓ |
| Associates | ✓ | ✓ |
| **Advanced Search Tab** | | |
| Concordances | ✓ | ✓ |
| **N grams Tab** | | |
| Match N grams | ✓ | ✗ |



**or**

The compare buttons on the Words tab and the Advanced Search tab will display concordance lines, frequency graphs and collocations for the main DIY corpus and a reference corpus side by side.  Frequency and collocation results are only available for some kinds of queries made using the Advanced Search tab.

### Files & Tools Tab

- Load and Save DIY Corpus Files; Load a DIY corpus file to use as a reference corpus;
- View the File List with statistics for Type-Token Ratios, Paragraph, Sentence and Word lengths;
- Get Keywords, Key Keywords, Wordlist Statistics and Average LL & MI Scores for the whole DIY corpus, the selected category, or the selected text;
- Highlight all the sentences in one text according to collocational strength.



These buttons are to load and save DIY corpora; the buttons under Reference Corpus are to load a previously saved DIY corpus to use as a reference corpus.

You can edit the corpus, category and filenames here.

You can view one text from the DIY corpus using this button.

**Words tab**

- Type in a word or find it on the wordlist and get concordance lines (with frequency and collocation data);
- Copy, save or share the word frequency list by right-clicking, double clicking or using a long tap on the table;
- Show tables of collocations, n grams or key associates for the word selected in the wordlist table;

You can enter more than one word in the search box, but more complicated searches can be executed using the Advanced Search tab.



The Selected Word is the one used for the top right buttons. *the* is unlikely to be interesting so choose a word further down the list first.



All the functions are available on all platforms; on narrower screens, you can hide or show a group of buttons using the corner button.

The DIY corpus wordlist with frequencies can be viewed using an extra button for narrow screens under All Words.

**Advanced Search Tab**

- Search for more than one word form at a time, entering all the different words in one search box separated by space;
- Use * to represent zero or more letters; use _ between two words to search for combinations of single and multiple words;
- Use some of the readymade wordlists from the Research Tools tab for queries;
- Search for combinations of words with multiple choices for each slot with the combinations occurring within a span of 5 running words;
- Get frequency tables, plots or concordance lines.





The Wordlists used in the Wordlist Statistics functions on DIY corpora can be viewed here.

Some wordlists can be copied for DIY searches using this button

Once it has been copied, the list can be moved up or down to other slots.

Notes:
- The _ phrases option allows you to combine the results of searches such as for reporting phrases: *according_to states stated claims claimed* will show results for "according to" as well as "states", "stated", "claims" and "claimed".
- When searching a corpus for multiword units, it is faster to find the items with the lowest frequency first and then filter out those not containing the required combinations, rather than looking through all the hits for a high frequency item.  If the word or words entered in the first box have higher frequencies than one of the other boxes, the Optimize Node button will light up and you can either click it to automatically choose the lower frequency box as the node or click the tool button again to proceed with your original node(s).

**N grams tab**
- Find strings of words (n-grams) which occur multiple times with no other words in between.
- View n grams from your DIY corpus with their frequencies in the online readymade corpus.  You can use all the n grams in that corpus or only n grams from a single category.



For 2 or 3 grams, a higher minimum frequency is recommended.

**Overview**
You can perform a range of operations on DIY corpora using tPM.  Other corpus tools generally have more flexibility in terms of specific settings, such as the statistical measures or choice of parameters for various functions.  However, tPM offers the ability to use readymade corpora as a baseline reference for many functions, and all these functions are available on all platforms.

| Focus | Function(s) |
|---|---|
| Measuring the difficulty of a text or a collection of texts | **Files & Tools tab: File list**<br>• STTR;<br>• Average sentence length;<br>• Average word length.<br>**Files & Tools tab: Wordlist stats**<br>• General Service lists;<br>• Academic Word List. |
| Determining whether a text or a collection of texts is similar to a readymade corpus (basic register analysis) | **Files & Tools tab: File list**<br>• STTR;<br>• Average sentence length;<br>• Average word length.<br>**Files & Tools tab: Wordlist stats**<br>• Modals;<br>• First and second person pronouns. |
| Finding useful vocabulary or the topics of a text or collection of texts | **Files & Tools tab: Keywords**<br>• Words which are repeated more often than expected (given the huge differences in normal word frequency);<br>• Words specific to a text or the collection, often related to topics;<br>• Names and places repeated in the text.<br>**Files & Tools tab: Key Keywords**<br>• Words which are repeated in several texts more often than expected;<br>• Topic words, names and places which are more widespread than just one text. |
| Locating specific sentences in a single text which may show unusual combinations of words. | **Files & Tools tab: Highlight MWUs**<br>• Colour-coding sentence by sentence, to show the strength of the combinations of words, using a readymade corpus as a baseline;<br>• Lists of collocations found in each sentence;<br>• A sense of what might be prominent or creative or a miscollocation because of unusual combinations. |

In order to use the following tools, you need to first download vocabulary profile information using the button on the Files & Tools tab under Main Corpus.  This collects four sets of information which can be stored in the tPMHDDIY file the next time you save it: vocabulary profile data for your DIY corpus using Nation's COCA+BNC lists; vocabulary ranking data for your DIY corpus using the currently selected online readymade corpus (default BNC:Academic); word family information using these two resources; and wordlist statistics for your DIY corpus on a word-by-word basis.  Once these data have been downloaded, Dictionary Style and Links Within Texts concordance line ranking scores will automatically be calculated.  This means you can use these additional concordance line ranking methods for DIY Corpora.  Grids for patterns of lexis in each DIY corpus text are also calculated, following a simplified procedure from the one described by Hoey (1991).  In tPM, only simple repetition of word form and word family is included (no substitution is carried out).

**Show Lex. Cohesion**
Select the DIY corpus text for analysis using the arrow buttons on the Files & Tools tab; you will see a preview of the beginning of the currently selected DIY corpus text.  When you click the Show Lex. Cohesion button, two tables of results will be generated. The first will show a matrix of links between sentences based on the repetition of words from the same families.  In this matrix, any sentence pairs with fewer than 3 links will be shown in light grey (as only those with 3 or more links are counted as "bonds".  This grid is limited to 1000x1000 (i.e. texts with no more than 1000 sentences).  The second table will show each sentence of the text along with its bond score.  Words contributing to bonds will be shown in green highlight.  The table is sorted so the top scoring sentences appear at the top.

**Keywords / Wordlist Statistics / N-grams on a text-by-text basis**
With the Vocabulary Profile Information for your DIY corpus stored in memory, Keywords and Wordlist Statistics can be calculated on a text-by-text basis.  N-grams also have an option for text-by-text results.  These will show all the results for all texts in your DIY corpus, with an additional column allowing you to filter results by text name or by keyword/feature/n-gram.  This is a convenient way to generate a master spreadsheet of keywords, wordlist statistics and n-grams using each individual text as the level of study.

**Show Vocab Profile**
After you have collected the Vocabulary Profile Information for your DIY corpus, you can use this button to show two sets of results: the vocabulary profile of each DIY text compared against Nation's COCA+BNC wordlists and the matches for vocabulary in each DIY text compared to the word rankings in the readymade online corpus (the corpus selected when vocabulary profile information was downloaded).  These two lists handle word families differently: Nation's COCA+BNC list is based on word families; the word rankings for the readymade online corpus are based on specific word forms (i.e. *disaster* may be ranked differently from *disastrous*).  From either of these tables of results, you can highlight vocabulary by level or rank using the highlight button.
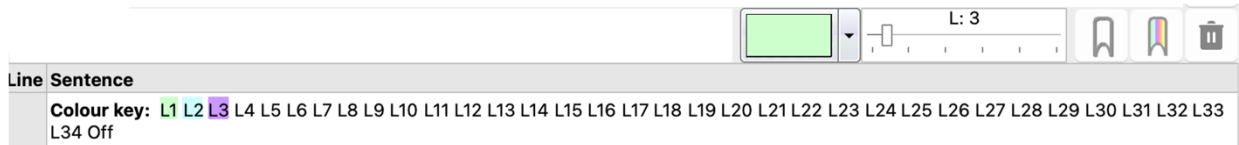
The drop-down colour picker allows you to select a colour for highlighting.
The track bar allows you to select which range of levels or rankings will be highlighted (default L1-L3 or K1-K3).
The white bookmark button will highlight all words which have levels or rankings up to and including the track bar selection using the currently selected colour.  Only those words which are not currently highlighted will be changed.

The multi-coloured bookmark button will highlight all words up to the levels or rankings shown in the track bar, moving through the list of available colours.  When the last colour has been reached, the process will stop.
The rubbish bin button simply clears all highlighting.

| Line | Sentence |

Colour key:  L1 L2 L3 L4 L5 L6 L7 L8 L9 L10 L11 L12 L13 L14 L15 L16 L17 L18 L19 L20 L21 L22 L23 L24 L25 L26 L27 L28 L29 L30 L31 L32 L33 L34 Off

| What you can do | How to do it |
|---|---|
| Highlight vocabulary in the top three levels using green for level 1, light blue for L2 and purple for L3. | Click the rainbow bookmark (the default settings for the track bar and colour picker don't need to be changed). |
| Highlight vocabulary in the top three levels using green.  Highlight all other words on the vocabulary profile in blue.  Highlight any words not on the profile lists in red. | Step 1: Click the white bookmark button (the default settings for the track bar and colour picker don't need to be changed).<br><br>Step 2: Change the colour to light blue and move the track bar to L34. Click the white bookmark button.<br><br>Step 3: Change the colour to red and move the track bar to the end (all).  Click the white bookmark button. |

Finally, if you use tPM for academic research, please do cite the app in your list of academic references.

Jeaco, S. (2017). Concordancing Lexical Primings: The rationale and design of a user-friendly corpus tool for English language teaching and self-tutoring based on the Lexical Priming theory of language. In M. Pace-Sigge & K. J. Patterson (Eds.), *Lexical Priming: Applications and Advances* (pp. 273-296). John Benjamins.

For some of the background to these methods and approaches, please see the tPM Help Selected Bibliography available from https://www.theprimemachine.net/help.html.



Dr. Stephen Jeaco - 杰大海
www.theprimemachine.net

Hoey, M. (1991). *Patterns of Lexis in Text*. Oxford University Press.

Nation, I.S.P. (2017). The BNC/COCA word family lists. Available from http://www.victoria.ac.nz/lals/staff/paul-nation.aspx

First published: Thursday, 17 March 2022

Last updated: Sunday, 10 September 2023