# *SketchEngine Workshop on Wordlists and Concordances*
# *Prof. Fiona Farr, UL*

**Note 1:** If you haven't had a previous introduction to corpus linguistics or SketchEngine, it is advisable that you watch this recorded session before attempting the tasks below: https://media.heanet.ie/page/e83bca1503294580a716dab4c310fc7e
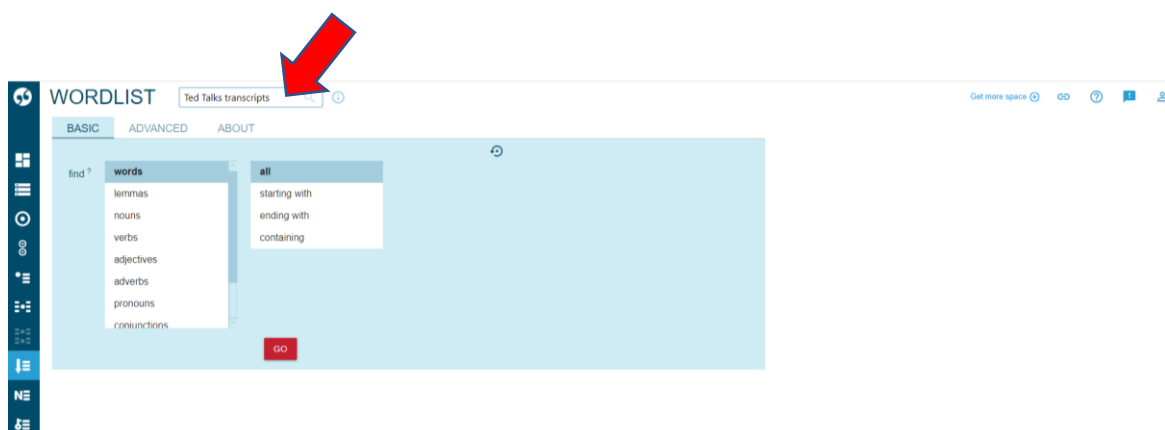
**Note 2:** In order to complete the tasks below you need to have access to SketchEngine.eu through the institution to which you are affiliated. You can check if your institution has a licence by clicking on the log-in button at SketchEngine.eu, going to institutional log-in, searching for your institution, clicking on it and signing in with your regular institution credentials.

**Note 3:** The SketchEngine Platform is well supported through integrated video links (click the 'about' tab in the various functions) and through the SketchEngine Channel on YouTube: https://www.youtube.com/channel/UCo2fn2_SNxCikCSAFCBcWBw

## **Wordlists**

Watch this Video: https://www.youtube.com/watch?v=nqpCIICCEdw&feature=emb_logo

1. Using the TedTalks corpus, run a wordlist for all the **words** (not lemmas – what is the difference?)



What are the top 5 most frequent words?

Is there a limit on the number of words returned in a frequency list in SketchEngine? The answer is on the screen.
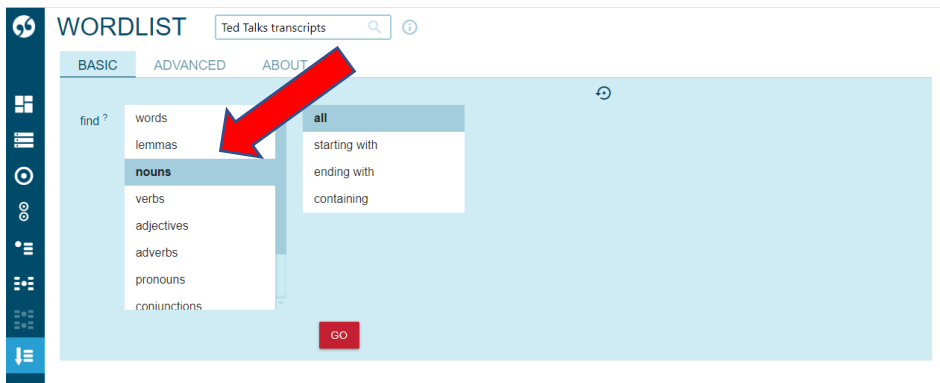
2. From here, run a concordance for the word 'this' by clicking on the three dots next to it on the frequency list.



3. What are the top 5 nouns in the corpus?



4. Download the noun frequency list to an Excel file and rearrange alphabetically in Excel. What is the first and the last noun on the list?

5. Try to get the top 5 regular past tense verbs from the wordlist function. Try this:



What interference do you find in the results?
Try this (full list of tags is here: https://www.sketchengine.eu/penn-treebank-tagset/):

What did you get?

After you press 'go' in the previous, run a concordance from the next screen (three dots next to the frequency number):



Do you get what you want? How do you find the most frequent verb from here? Try this:



Did you find exactly what you were looking for? If not, what can you now do manually to get exactly what you are looking for?

## Concordances
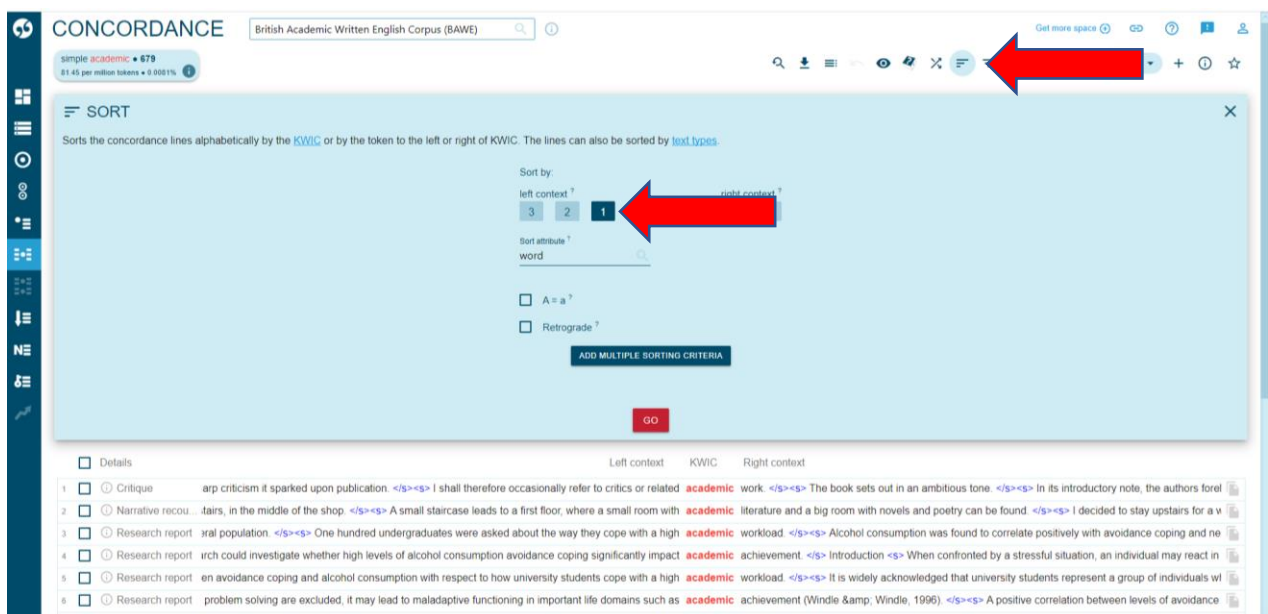
Watch this Video: https://www.youtube.com/watch?v=FzI6tbO5EvQ&feature=emb_logo

1. Using the BAWE corpus, search for the work 'academic'. How many hits did you get?



2. Which word occurs most frequently one place to the right of 'academic'? Try this:



What results did you get? Can you find which is most frequent (manually) from here?

3. Here's another way to do it more automatically using the collocation function:

What do you get? The results appear alphabetically by default. To rearrange the words by frequency, click on the column heading that says 'cooccurences':



4. Try to find the top 5 regular past tense verbs through the concordance function, using a basic search for *ed. What are the issues with the results?

5. To resolve some of these issues, let's rerun the concordance with some more specific criteria (filters) from the advanced tab. Try this to tell SE to include only examples of *ed which has pronouns one place to the left of it:

6.  The results you get from 5 above, will appear randomly. You can rearrange them alphabetically by clicking on the arrow next to the KWIC column heading, and to the left and right by using the sort function as above. To get a list of which of the past tense verbs is most frequent, click on the 'frequency' tab and go:



Does this give you what you need? Is there any interference in the results that you get?