# 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



**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.

▶ 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.