# #LancsBox X
# Text

**Corpus overview**

In this task, you will explore the properties of individual texts in a corpus. Go to the Text tool in #LancsBox X and select the BNC2014 corpus (whole corpus). Provide the following information.

1. Number of files in the BNC2014 is _____88,171_____.

2. The **largest** file has _____123,259_____ tokens.

3. The **smallest** file has _____30_____ tokens.

4. The number of files that are equal or larger than 10,000 words is _____1,820_____.

5. The most lexically diverse file in **Academic prose** is _AcaMedRv88.xml_ with MATTR$_{50}$ _____0.88_____.

6. The least lexically diverse file in **Informal speech** is _Sp0m2f99.xml_ with MATTR$_{50}$ _____0.64_____.

| Overview | ▼ 2,879 | | | </> ▼ |
|---|---|---|---|---|
| **Name** | **Tokens** ▼ | **MATTR$_{50}$** | **MTLD** | **genre** ▼ + |
| AcaNatBk13.xml | 51,157 | 0.73 | 43.38 | academic pr |
| AcaMedBk9.xml | 49,716 | 0.78 | 71.88 | acad |
| AcaPleBk15.xml | 49,629 | 0.82 | 101.7 | |
| AcaSocBk13.xml | 49,298 | 0.80 | 86.02 | |

**Tip**: To find files with given properties click on the filter icon ▼ and apply an appropriate filter . Then click on a relevant column to sort files. Columns can be added by clicking on the + sign.

## Lexical diversity

There are a number of lexical diversity measures showing the range of different words in a text. For the comparison of files of varying sizes, we need to go beyond a simple Type/token ratio (TTR) and compute more sophisticated measures such as Moving average type/token ration (MATTR) or a Measure of textual lexical diversity (MTLD).

**Type/token ratio (TTR)** expresses the proportion of types relative to the proportion of tokens. It is calculated by dividing the number of types in a text or corpus by the number of tokens. It decreases with text size so it cannot be used to compare texts of different sizes in a corpus.

**Moving average type/token ration (MATTR)** is calculated by dividing a text into standard sized overlapping segments (e.g. 50 words in MATTR$_{50}$) as a window moves through the file one token at a time. TTR is calculated for each overlapping segment and then the mean value of the TTRs is taken. MTTR is suitable for comparing texts of different sizes.

**Measure of textual lexical diversity (MTLD)** is the mean number of words in a text that maintain a given TTR value of .72.

**Distribution of linguistic features in texts**

In the BNC2014, search for occurences of the past tense using the smart search PAST_TENSE ( don't' forget to include the underscore). Answer the following questions:

1. In how many texts does the **past tense** occur? _____81,397_____

2. In how many texts does the **past tense** occur with a relative frequency that is higher than the average

relative frequency for the whole corpus?__Number of texts: _____29,024_____

3. In how many <u>newspaper</u> texts does the past tense occur **at least once**? _____46,112_____

4. In how many <u>newspaper</u> texts does the past tense **not occur at all**? _____4,098_____

---

Task 3

**Analysing individual texts**

In the BNC2014, find the text with the largest relative frequency of the search term fuck*. Provide information about this text.

1. Name of the text file: _____MagCla1338.xml_____

2. Genre: _____magazines_____

3. Source: _____Classic Rock_____

4. The swearword fuck* occurs ____13_____times in the text, which has ____148____tokens.

   This means on average, an f-word occurs every ___11.38_____ words.

5. The function of the swearwords in this context is

   To describe songs on an album which contains a high frequency of the swearword_____

   fuck*'. The article reports the release of the album and also highlights the frequent use ____

   of the term in an individual song._____

Developed at
Lancaster
University