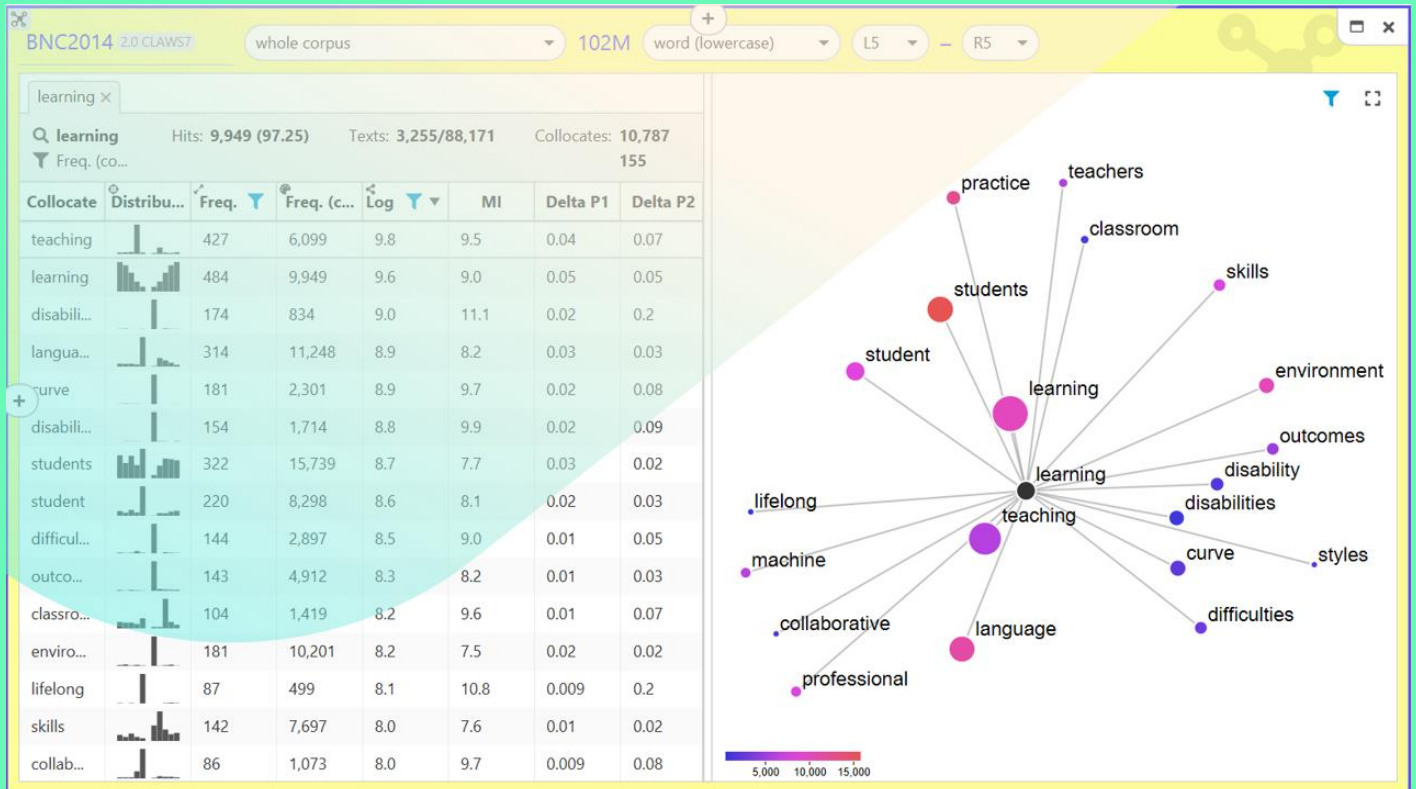# #LancsBox X



Data Driven Learning
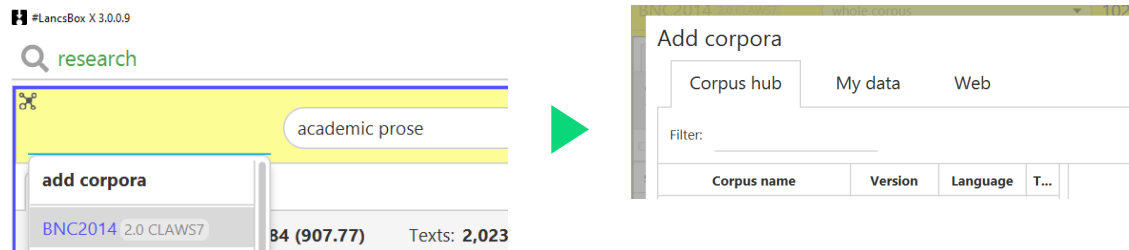
Professor Vaclav Brezina
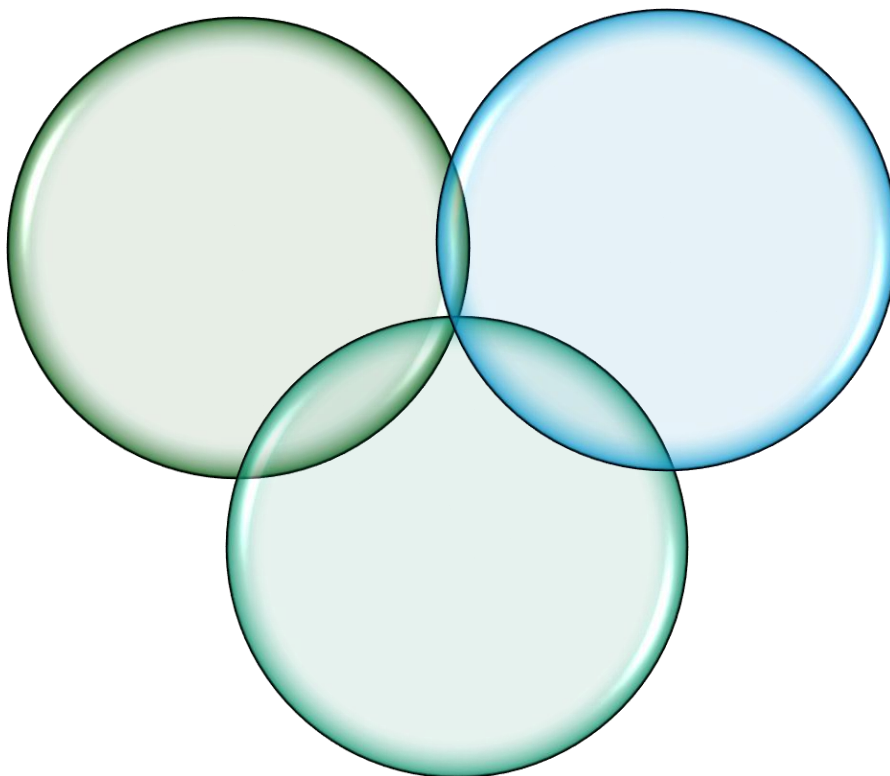
@ Lancaster University

# Starting with #LancsBox X

#LancsBox X is a powerful software tool for the analysis of large amounts of language. It can be used with your own data or the data provided.

The tool is very easy to use with an intuitive and flexible  UI.

1.  Download the most recent version of #LancsBox X from https://lancsbox.lancs.ac.uk
2.  Go to 'add corpora'> 'Corpus Hub' and select and download the British National Corpus 2014 version 2.

# Why is DDL important?

# Check your vocabulary knowledge

In this section, we will focus on the exploration of concordance lines to find meaning of words (their uses in context). We will also discuss frequency information. You will learn

- to use the KWIC tool
- to sort and filter concordance lines
- to use the summary table

**KWIC** — frequency — context — meaning

## Task 1 — Vocabulary test

**Test your vocabulary knowledge.** Check a box next to a word if you know the word. If you don't know the word, leave the box blank.

☐ back            ☐ lens            ☐ stuff
☐ moderate        ☐ neen            ☐ head
☐ dependent       ☐ footballer      ☐ desponation
☐ at firse        ☐ creature        ☐ memory
☐ stone           ☐ to register     ☐ to wrose
☐ efficiency      ☐ to sit          ☐ housing
☐ measurement     ☐ classical

Words checked: _____

**Find the meaning of the words from Task 1.** Using the KWIC tool, search for the words you were not sure about in the British National Corpus 2014. Note down their frequencies, genres where they typically occur and their meaning.

| Word | Frequency | Dominant genre/register? |
|------|-----------|--------------------------|
| moderate | | |
| | | |
| | | |
| | | |
| | | |
| | | |

_____

_____

_____

_____

_____

_____

Developed at
**Lancaster University**

# Collocations in context with GraphColl

In this section, you will explore collocation graphs and networks using the GraphColl tool in #LancsBox X. The GraphColl tool identifies collocations and displays them in a table and as a collocation graph or network. It can be used 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.

**GraphColl**    collocation graphs   collocation networks
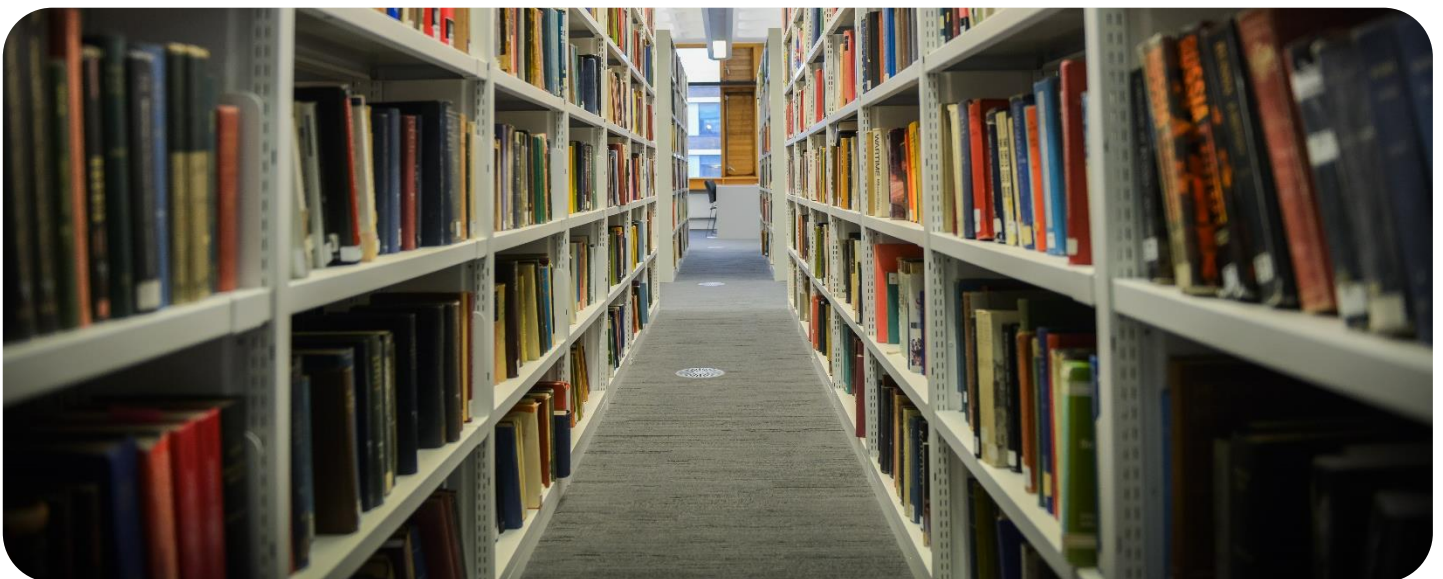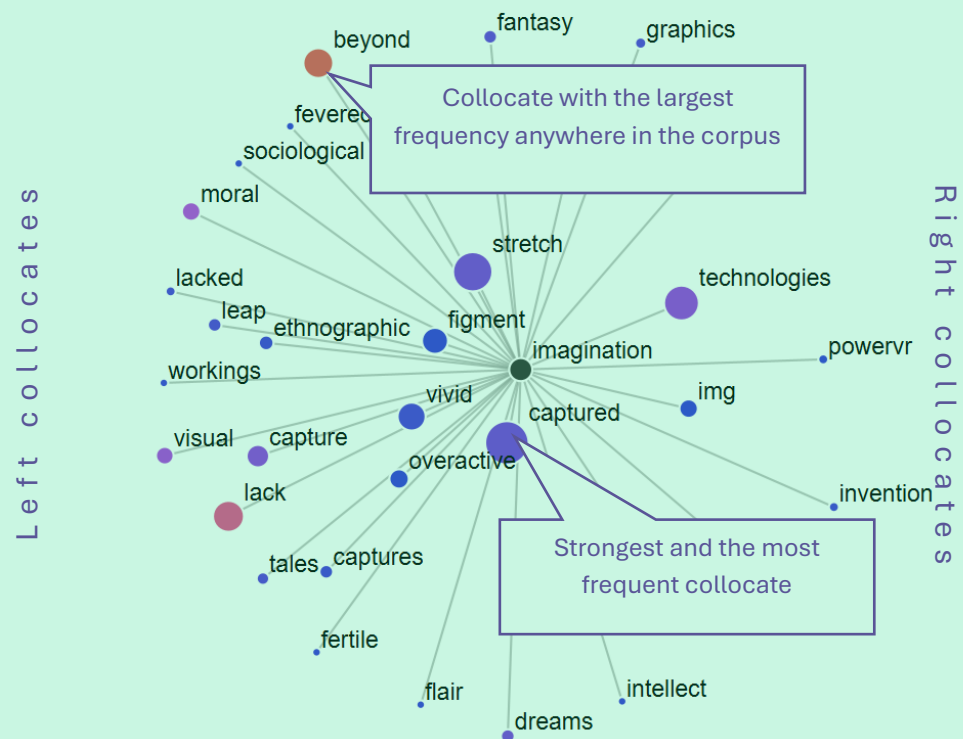
**Task 3**    Finding collocates

In this task, you will practice finding collocates and interpreting collocation statistics. Go to the GraphColl tool in #LancsBox X, select the BNC2014 and search for the expressions in the table below.

Note down top collocates according to log Dice and the collocation frequency.

| Search term | Top 3 log Dice collocates | Most frequent collocate |
|---|---|---|
| `vision` | | |
| `avid` | | |
| `unclear` | | |
| `[hw="disagree" pos="V.*"]` | | |
| `[hw="ring" sem="B.*"]` | | |
| `[hw="ring" sem="Q.*"]` | | |

# Collocation graph

A collocation graph shows the relationship between a node, which is in the middle of the graph, and its collocates, which are displayed around the node. The closer the collocate is to the node, the stronger the association. The position of the collocates indicates the position in the text, before or after the node, while the size of the collocate reflects the frequency of co-occurrence. Finally, the colour indicates the frequency of the word anywhere in the corpus on the scale from blue (small) to red (large).
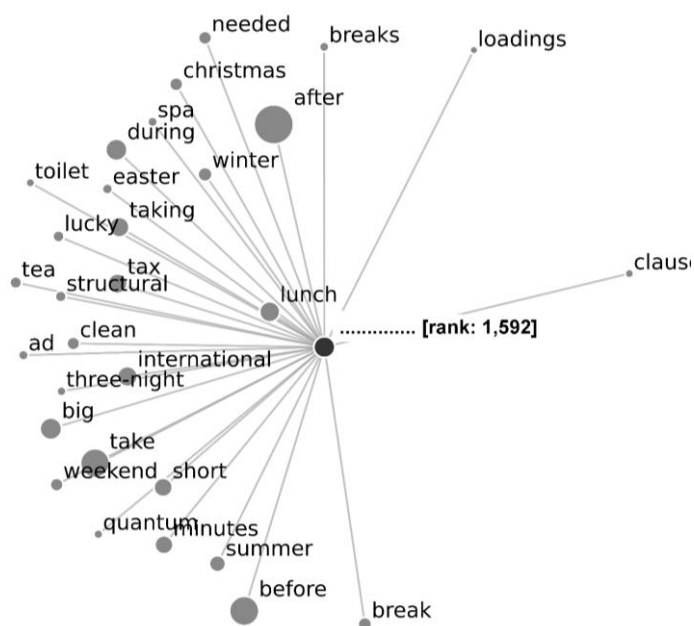
fantasy
graphics
beyond

fevered
Collocate with the largest
frequency anywhere in the corpus
sociological

moral

stretch
lacked
technologies
leap
figment
ethnographic
imagination
workings
powervr
vivid
img
captured
visual    capture
overactive
lack
invention
Strongest and the most
frequent collocate
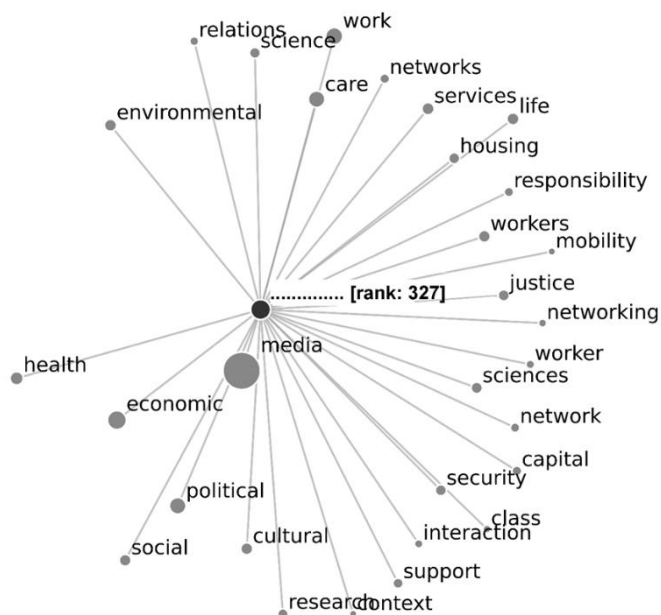tales  captures

fertile

flair
intellect
dreams

**Left collocates**

**Right collocates**

Look at the collocation graphs below. Each graph represents collocates around a key word, which has been hidden. Choose an appropriate key word from the box below.

> break (n), eat (v), good (adj), play (v), rude (adj), social (adj)



relations work
science
networks
care services life
environmental housing
responsibility
workers mobility
............. [rank: 327]
justice
networking
health media worker
sciences
economic network
capital
security
political class
interaction
social cultural support
research context



needed breaks loadings
christmas
spa after
during
toilet winter
easter
lucky taking
tea tax clause
structural lunch
ad clean ............. [rank: 1,592]
international
three-night
big
take
weekend short
quantum minutes
summer
before
break

# Keywords and Key phrases

In this section, you will explore words and phrases important in a particular corpus when compared to a reference corpus. These are called *keywords* and *key phrases*. You will learn how to:

- Create and visualize a frequency list.
- Create a keyword list.
- Understand keyword statistics.

**Wordlists**      ( keywords )  ( key phrases )

**Task 5**      Understanding wordlists and n-gram lists

Select the Academic prose subcorpus of the BNC2014. Create and visualize words and n-grams.

1. First create a wordlist based on lemmas. How many lemmas does the list contain? _____

2. Click on the filter icon ▼ and select nouns (type in _N). How many nouns are there? _____

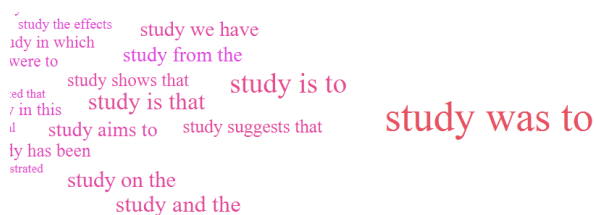3. What are the top 5 academic nouns?  _____

   _____

4. Change the settings as follows [ word (lowercase) ▼ ] [ 3-grams ▼ ] and identify all trigrams.

   What are the top 3 academic trigrams?  _____

5. Click on the filter icon ▼ and select trigrams starting with 'study', the most frequent academic noun.

   How many trigrams have you identified? _____

6. Visualize and interpret these results.

study the effects    study we have
dy in which    study from the
were to
study shows that    study is to
ed that
r in this    study is that
study aims to    study suggests that
ly has been
strated    study on the
study and the

study was to

Keywords and key phrases are words that occur with a considerably higher frequency in a given (sub)corpus compared to a reference (sub)corpus. In this task, you will explore key 2-grams in the Academic prose subcorpus compared to the whole BNC2014.

1.  First create a 2-gram list based on the Academic prose subcorpus and note down the first three items:

    _____

2.  Click on the keyword icon 🔑 and select the whole BNC2014 as a reference corpus.

3.  Note down the top 10 key bigrams, i.e. combinations typical of the academic prose subcorpus.

    _____

    _____

4.  Write a short description of key features of academic writing based on the bigrams above.

    _____

    _____

    _____

    _____

    _____

    _____

Developed at
Lancaster
University

# #LancsBox X in your research

In this section, you will design a small study in the area of your interest. The focus is on:

- Conceptual grounding.
- Research question(s).
- Operationalization and study design.
- Data collection.
- Data analysis.

**Research** | design | methodology

## Task 9    Designing a corpus study

In this task, you will design a mini-study. First, think about a topic in the area of your interest.

**Topic:** _____

Then, think of a specific question you would like to ask:

Specific question: _____

Is it a yes/know question? If no, think of an aspect of your question that can be formulated as a yes/no question.

**Yes/No RQ**: _____

**Data:** Is there an available corpus that can be used to answer the RQ? Yes/No _____

If no, can the data be obtained online? URL_____

**Operationalization:**

#LancsBox X tools: _____

Search terms: _____

Comparisons: _____

Possible challenges: _____