

# *Corpus Gysseling* application manual

## **Table of contents**

<b>Introduction</b>	<b>4</b>
Information about the corpus	4
Linguistic Annotation	4
Metadata categories	5
Date	5
Witness Year	5
Text Year	5
Permissive / strict	5
Localisation	6
Region	6
Place	6
Kloeke location code	6
Text type	6
Fictionality	6
Genre	6
Subgenre	6
Authenticity	7
Title and author	7
Words in title	7
Title of containing publication	7
Author	7
<b>Application user manual</b>	<b>8</b>
Getting started	8
Searching the corpus	9
Simple Search	9
Search	9
Wildcards	9
Reset	10
History	10
Global settings	10
Extended Search	11
Upload a list of values	13
Part of Speech dialogue box	13
Numerical part of speech code	13
Cliticity	13
Complete word or word part	14

Filter search by	15
Filter by date	15
Filter by other metadata	15
Advanced search	16
The query builder	16
The tab Search	17
Token attributes	18
Adding attributes to a token box	18
Function of the two +-buttons in a token box	19
The tab Context	19
Managing sequences of token boxes	20
Uploading value lists in the query builder	20
Copy to CQL editor	20
Expert search	21
Copy to query builder	22
Import query	22
Gap filling	22
Viewing results	24
Per Hit view	24
Sorting results	24
Grouping results	25
Per Document view	26
Sorting results	26
Grouping results	27
Exporting results	28
Information about a document	28
Content	28
Metadata	29
Statistics	29
Exploring the corpus	29
Documents	29
N-grams	30
Options	30
Example	31
Statistics (frequency lists)	32
Options	32
Example	32
<b>Appendix: Corpus Query Language</b>	<b>33</b>
CQL support	33
Supported features	33
Differences from CWB	34
(Currently) unsupported features	35

Using Corpus Query Language	35
Matching tokens	35
Sequences	36
Regular expression operators on tokens	36
Punctuation	36
Case- and diacritics-sensitivity	37
Matching XML elements	37
Labeling tokens, capturing groups	38
Global constraints	38

# Introduction

This manual describes the corpus exploitation environment for the *Corpus Gysseling*. The corpus application is developed by the Dutch Language Institute (Instituut voor de Nederlandse Taal or INT). The backend of the application is the BlackLab Lucene based search engine developed for corpora with token-based annotation (<https://blacklab.ivdnt.org/>). The web-based frontend is a further development of the corpus-frontend application developed by INT (<https://github.com/instituutnederlandsetaal/blacklab-frontend>). Its design is inspired by the first version of the OpenSoNaR user interface by Tilburg University and Radboud University (<https://github.com/Taalmonsters/WhiteLab2.0>).

## Information about the corpus

The *Corpus Gysseling* is the collection of all thirteenth-century texts that served as source material for the *Vroegmiddelnederlands Woordenboek* (VMNW). It is the digital edition, enriched with word type and lemma, of the thirteenth-century material from the *Corpus of Middle Dutch texts (up to and including the year 1300)*, published in the period 1977-1987 by the Ghent linguist Maurits Gysseling. The linguistic annotation was manually verified.

A first online accessible version of the corpus was launched on 25 April 2012.

In 2018, the corpus was integrated in the Nederlab portal, in a new version, containing corrections of the linguistic annotation and additional metadata to the texts in the corpus.

In this new version several corrections have been made to the added metadata in the corpus.

## Linguistic Annotation

The Part of Speech tagging of the *Corpus Gysseling* was originally done using the same principles for annotation and tagset as for the Corpus Van Reenen Mulder, using a system of three digits. This numerical encoding can be used to search the online corpus. A detailed description can be found [here](#).

The original Part of Speech tagging has been converted into the Part of Speech tagging that has been done using the tagset and tagging principles for the annotation of diachronic corpora of historical Dutch, developed in the context of the CLARIAH+ project. This annotation layer has been added to the corpus, and can also be used to search the online corpus. A detailed description can be found [here](#).

The Early Middle Dutch word forms all have a modern Dutch lemma. For words no longer used in modern Dutch, a modern lemma has been constructed using the same linguistic principles applicable to still existing words. More information about the used lemmatisation principles can be found in [Lemmatiseerprincipes voor GiGaNT, het centrale lexicon van het INT](#).

## Metadata categories

The *Corpus Gysseling* has been enriched with an elaborate set of metadata categories. These metadata are described below. In the corpus application it is possible to limit a search by filtering on metadata categories.

### Date

With respect to dating the source, a distinction is made between the date of the manuscript in which a text was handed down (Witness Year) and the period in which a text was produced (Text Year).

### Witness Year

For each document in this corpus, we indicate the period in which the manuscript, providing us the text, was written. Witness Year does not necessarily refer to the period in which the text itself was written. It only concerns the carrier of the text.

Witness Year cannot be stated with the same accuracy for every document. For example, the manuscript named *Corp. I, 0002B, Gent* can be dated exactly to 1250, while the source in which the text of *Wrake van Ragisel* is included, originated between 1260-1280.

### Text Year

Text Year represents the year or the period in which a text was conceived. For autographs - manuscripts written by a writer himself - Witness Year and Text Year generally coincide; this is often the case with non-fictional administrative and legal texts.

In fictional texts, Witness Year and Text Year often differ from each other. For example, the text of *Reinaert G* was conceived between 1185-1191 (Text Year), while the copy in the Municipal Library of Rotterdam dates from the period 1260-1280 (Witness Year).

### *Permissive / strict*

It is possible to do a permissive and strict search, both for Witness Year and Text year. What exactly is the difference between the two options? An example can clarify this. Suppose you want to investigate sources that were written between 1200 and 1225.

If you choose to do a Strict search by Witness Year, the search query will only result in manuscripts that were produced later than 1200 but before 1225. The *Corpus Gysseling* has only one such text: *Floyris ende Blantseflur*. (Note that a Strict search in Text Year gives a different result: *Aiol*.)

If, on the other hand, you choose the option Permissive in Witness Year, four documents are found, one of which is Heinric van Veldekes *Sente Servas*, which was handed down in a manuscript dating from the period 1190-1210. This dating does indeed partly fall within the indicated period 1200-1225. (Note that a Permissive search in Text Year results in eighteen documents.)

## Localisation

### Region

Within the four different countries, the following regions are distinguished: Brabant-Noord, Brabant-Oost, Brabant-West, Limburg, Nederrijn, Oost-Holland, Oost-Nederland, Oost-Vlaanderen, regio Venlo, Utrecht, West-Holland, West-Vlaanderen and Zeeland. Sometimes a document mentions two regions, for example a charter with an agreement between the Count of Holland and the Count of Flanders.

### Place

If it was possible to determine exactly where a particular text originated, this Place (in modern spelling) will be mentioned.

### Kloeke location code

In the 1920s, the Dutch dialectologist G.G. Kloeke designed a system with unique designations for thousands of places and hamlets in the Netherlands, Flanders, French Flanders and north-western Germany. Each Kloekecode consists of an uppercase letter followed by three digits and a lowercase letter, for instance K244p (Antwerp). It is possible to filter documents based on this so-called *Kloekecode*. More information about (searching with) the Kloekecode can be found [here](#).

### Text type

All texts in this corpus are provided with metadata to help determine fictionality, genre, subgenre and authenticity of the text. These metadata can be filtered during the search.

### Fictionality

The documents have been divided into fictional texts (fiction) and non-fictional texts (non-fiction). Generally speaking, this contemporary classification will also apply to Middle Dutch texts. However, there are cases in which we will classify a text as fiction and medieval people as non-fiction. Jacob van Maerlant's *Der naturen bloeme* for instance is an impressive overview of what was available in his time in terms of knowledge of nature, but it also contains descriptions of miraculous creatures. For this corpus, *Der naturen bloeme* has been classified as non-fiction.

### Genre

The texts are divided into two main genres: prose and verse.

### Subgenre

The texts can be sorted out using one or more of the eighteen different subgenre labels. These subgenre labels may indicate a general text category as well as a more specific one; they may touch either the content of a text or its form (e.g. *stanzaic*).

- *Artes*: non-religious and non-fictional texts written for utilitarian and instructive purposes
- *Biblical texts*: Bible, lectionaria, diatessara
- *Biology*: texts with biological information
- *Chivalric romance*: mostly versified stories about the idealized world of knights
- *Drama*: stage plays

- *Epic*: narrative literary texts
- *Legendary hagiography*: descriptions of the life of saints
- *Lexicon*: texts offering lexicographical information (a vocabulary or wordlist)
- *Linguistics*: texts with linguistic information
- *Medicine*: texts with medical information or instructions
- *Metallurgy*: texts with information on (the use) of metals
- *Natural history*: texts providing in a popular form the scientific study of nature in general or animals, plants, even stones, in particular
- *Philosophy / ethics*: texts with life lessons, instructions of how to behave in life and social traffic, based on specific philosophical or secular grounds
- *Religious*: texts dealing with religion or religious matters
- *Secular*: texts dealing with non-religious matters
- *Science*: texts offering scientific information
- *Stanzaic*: texts composed in the form of stanzas
- *Theology / ethics*: texts with life lessons, instructions of how to behave in life and social traffic, based on theological or religious grounds

#### Authenticity

The degree of authenticity is indicated for each text. Seven categories are distinguished: clean copy (i.e. a copy of a document without editing notations), copy, dorsal (i.e. additional or restrictive notes on the back of a deed), draft, original, reworking and translation.

#### Title and author

##### Words in title

Every document has a title. Usually this is the title as we know it from the editions used for this corpus.

This search field is provided with a list, which contains suggestions for possible search terms in alphabetical order, based on the characters typed in.

##### Title of containing publication

All 2231 documents of this corpus come from nine different publications. A vast majority can be found in the *Corpus van Middelnederlandse teksten I* (2192 documents).

##### Author

It is possible to search by author name. However, for most of the documents in this corpus the author is unknown or uncertain. Only five names of authors of Early Middle Dutch texts in this corpus have been handed down: Willem van Affligem, Broeder Geraert, Jacob van Maerlant, Willem and Henric van Veldeke.

# Application user manual

The language of the corpus application is set to Dutch by default. Press the globe icon in the top right corner to select English.

## Getting started

Here are a few examples of what you can do with the corpus application (the links will take you to the application):

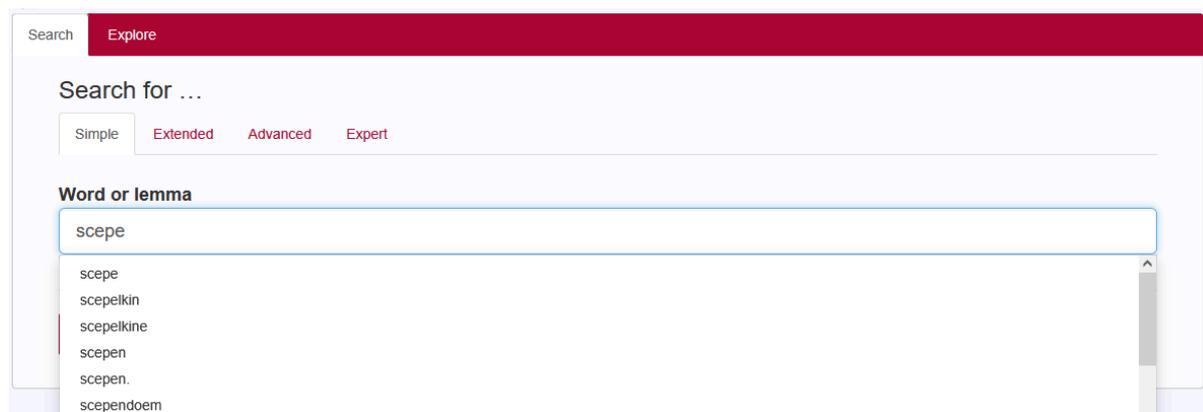
- To search for a word literally in the form you specify, use Simple Search or Word in Extended Search.
  - Simple Search for [scepe](#)
  - Extended Search for Word [dochter](#)
- To search for different spellings and inflections/conversions of a given lemma, use Simple Search or Lemma in Extended Search.
  - Simple Search for [leven](#)
  - Extended Search for Lemma [leven](#)
- To search for words or lemmata satisfying a certain pattern, use *wildcards* in Simple Search or Extended Search, or *regular expressions* in Advanced Search or Expert Search.
  - words and lemmata starting with *ver* and ending with *len* in [Simple Search](#)
  - only word forms starting with *ver* and ending with *len* in [Extended Search](#)
  - only lemmata starting with *ver* and ending with *len* in [Extended Search](#)
  - lemmata starting with *ver*, ending in *eren* with (mostly) one syllable in between in [Expert Search](#)
  - lemmata starting with *ver*, ending in *eren* with (mostly) one syllable in between in [Expert Search](#)
- To search for a multi-word pattern, e.g. all adjectives appearing before a given lemma as a noun, use the query builder in Advanced Search or use Expert Search.
  - adjectives before the lemma *huis* in [query builder](#) in Advanced Search
  - adjectives before the lemma *huis* in [Expert Search](#)
- To see which unique forms occur as a result of your search, use Group Results.
  - example Group by annotation: [different adjectives before the lemma huis](#)
  - example Group by annotation: [different words preceding the word god](#)
- To explore the distribution of document properties in the corpus, use the Explore feature.
  - example: [characteristics of the date](#)
  - example: [subgenres in the corpus](#)

# Searching the corpus

## Simple Search

### Search

The Simple Search allows you to quickly search for specific words (e.g. *scepe*) or lemmata (e.g. *schip*). After entering a search term, a running list appears, which contains suggestions for possible search terms in alphabetical order, based on the characters typed in.



Note that this only works with a single word, like *schip*. If you want to use the autocomplete option for a lemma of which the word form contains one or more spaces (e.g. *schep laken*), the search query must be put within quotation marks (e.g. "schep laken"). Otherwise the search will yield no results (note that this differs for phrases that occur in the corpus, see below).

It is also possible to enter a phrase: *het schip* or *dit zijn de schepen*. You will then find all occurrences of that exact phrase. Furthermore, you can search for different values simultaneously by separating them without spaces by a vertical line, e.g. *god|man|lief* or - with the use of wildcards - *god|\*schip|vrouw*.

Keep in mind that when a historical word form corresponds with a modern Dutch lemma, you will not only find the desired historical word form, but also all word forms that can be traced back to that homonymous lemma. For instance, the search term *man* does not only result in all occurrences of *man*, but also in word forms such as *manne*, *man*, *mans*, which after all also belong to the lemma *man*. In order to only find the word form *man*, use the attribute *Word* in Extended Search.

Note that in Simple Search the patterns will be matched case-insensitively: *schip* will deliver the same results as *schip* or *Schip*. See the paragraph [Grouping results](#) in Per Hit view to see how it is nevertheless possible to distinguish between uppercase and lowercase letters, or go to Extended Search.

### Wildcards

In Simple Search, the use of wildcards can prove good service to search for specific word forms or lemmata. A wildcard is a symbol used to replace or represent one or more characters. The following two wildcards are supported:

\* The asterisk matches any character zero or more times. Therefore, searching for *a\*n* matches all word forms and lemmata that start with an *a* and end with a *n*, e.g. *an*, *allen*, *Alsmen*, *alreheleghen*, *alsmen*, *arresteren* but also *ane* (lemma *aan*), *allene* (lemma *alleen*) and *antwerdde* (lemma *antwoorden*).

? The question mark matches a single character once. Therefore, searching for *a?n* in Word or lemma matches *only* three-letter word forms or lemmata starting with an *a* and ending with a *n*, e.g. *an* (lemma *aan*) and *anden* (lemma *aan+de/het*).

This wildcard can be used more than once. Thus *a???n* matches *allen*, *anden*, *adden*, *armen* and *aluen*.

Note that searching with wildcards is limited to Simple Search and Extended Search. [In Advanced Search and Expert Search you can use so-called regular expressions instead of wildcards.]

### Reset

You can start a new search by pressing the Reset button. By doing so, both the search query and the hits found will be cleared. Your search history, however, will remain unchanged.

Note that it is also possible to start a new search by entering a new word or phrase in the search field Word or lemma.

### History

The History button will display your query history. Per search query there are several possibilities (as shown in the screenshot below): you can perform the search query again (Search), you can copy the search query as a link (Copy as link), you can download the search query as a file (Download as file), you can delete a single search query (Delete) or delete all search queries (Delete all).



Every search query has its own url. If you copy this url via History (Copy as link) or directly from the address bar of your browser, you can send it to someone else who can import this link via Import from a link. It offers that person the possibility to run the search on his own computer.

### Global settings

The Global settings dialogue, activated by pressing the wheel button, allows you to configure five settings: Results per page, Sample size, Seed, Context size and Wide View.

- *Results per page*: you can choose whether you want 20, 50, 100 or 200 results to be shown;

- *Sample size*: selecting a value here will instruct the search engine to return a random sample drawn from the complete result set. The sample size can be limited by
  - a percentage of the total number of search results (percentage);
  - the number of results displayed (count).
- *Seed*: a ‘random seed’ is a number used to initialize a so-called pseudo-random number generator. Keeping the same seed will ensure that two samples drawn from the same result set are identical. A new seed will most likely result in a different sample;
- *Context size*: by entering a number you can determine the number of words Before hit and After hit;
- *Wide View*: the default setting is ‘small view’; you can change to Wide View by ticking the checkbox.

The image shows a 'Global settings' dialog box with the following elements:

- Results per page:** A dropdown menu showing '20 results per page'.
- Sample size:** A dropdown menu showing 'percentage' and a text input field containing 'Sample size'.
- Seed:** A text input field containing 'Seed'.
- Context size:** A text input field containing 'Context size'.
- Wide View:** An unchecked checkbox.
- Close:** A button at the bottom right.

## Extended Search

Like in Simple search, Extended Search allows you to quickly search for specific word forms or lemmata. The search is performed in the same way as described for Simple Search. In this corpus the six main attributes you can search for are Word (more precise: word form), Lemma, Part of Speech, Numerical part of speech code, Cliticity and Complete word or word part.

In the search fields Word, Lemma and Numerical part of speech code enter the value of the attributes (or Upload a list of values) you are looking for. After entering a search term, a running list appears, which contains suggestions for possible search terms in alphabetical order, based on the characters typed in. Like in Simple Search, you can also enter a phrase for Word and/or Lemma. In the search fields Part of Speech, Cliticity and Complete word or word part you can select the desired values. All supported attributes are shown in the search form:

Simple **Extended** Advanced Expert

**Word**

Case- and diacritics-sensitive

**Lemma**

Case- and diacritics-sensitive

**Part of Speech**

**Numerical Part of Speech code**

**Cliticity**

**Complete word or word part**

In Extended Search it is also possible to search case- and diacritics-sensitive. Note that the default setting for search is case- and diacritics-insensitive. For example, searching for the Word *maria* will result in 127 occurrences of this name (*Maria*, *MARIA*, *maria*). By ticking the box Case- and diacritics-sensitive you will only find 56 hits of the Word *maria*, but none of *Maria* and *MARIA*. In order to directly find only occurrences of the Word (form) *Maria* (68x), use the search term *Maria* and tick the box Case- and diacritics-sensitive under the search field Word (as shown below).

Search for ...

Simple **Extended** Advanced Expert

**Word**

Case- and diacritics-sensitive

**Lemma**

Case- and diacritics-sensitive

Like in Simple Search, wildcards are supported in Extended Search. (See for a short explanation of wildcards [Simple Search](#).)

You can combine multiple fields (such as Word and Part of Speech) to form a single query. For example, searching for the Word (form) *was* together with the Part of Speech VRB will result in a list of all verbs containing the word form *was*. The syntax of your query is shown in the results: [\[word="was"&pos="vrb"\]](#).

## Upload a list of values

At the right side of the search fields Word, Lemma and Numerical part of speech code there is an option to Upload a list of values; those values must all be separated by a white space. Note that this function only works for \*.txt-files. (If you are using a text editor like Word, you have to save your file as a \*.txt-file first.)

Every word in the uploaded file will be added to the list of values to search for. To remove the word list simply delete all text in the search field or press the Reset button.

## Part of Speech dialogue box

Clicking on the pencil next to the search field Part of Speech provides you with the Part of Speech dialogue box.

Part of Speech	Type	Finiteness	Tense	Mood	Inflection
Adjective-Adverb	<input type="checkbox"/> Main	<input type="checkbox"/> Finite	<input type="checkbox"/> Present	<input type="checkbox"/> Indicative Conjunctive	<input type="checkbox"/> No inflection
Adposition	<input type="checkbox"/> Main or copula or auxiliary	<input type="checkbox"/> Infinitive	<input type="checkbox"/> Past	<input type="checkbox"/> Imperative	<input type="checkbox"/> -a
Adverb	<input type="checkbox"/> Copula or auxiliary	<input type="checkbox"/> Present participle	<input type="checkbox"/> Unclear	<input type="checkbox"/> Unclear	<input type="checkbox"/> -e
Conjunction	<input type="checkbox"/> Unclear	<input type="checkbox"/> Past participle			<input type="checkbox"/> -n
Interjection		<input type="checkbox"/> Unclear			<input type="checkbox"/> -r/-re
Common Noun					<input type="checkbox"/> -s/-th
Proper Noun					<input type="checkbox"/> -t
Numeral					<input type="checkbox"/> Other
Pronoun					
Residual					
<b>Verb</b>					

pos="VRB"

Reset OK

For most of the categories on the left you can tick certain features to further specify your search query. By doing so you can for instance delimit your search, as shown in the above screenshot. The query 'Verb - finite - present' will result in *es*, *kennen*, *quedden*, *selen* and numerous other hits.

## Numerical part of speech code

Each word has been given a morphological code consisting of three digits, the so-called 'numerical part of speech code'. The first digit indicates the Part of Speech; the second and third digits differ per Part of Speech. This is the original encoding of the corpus (see [Linguistic Annotation](#)).

## Cliticity

This attribute enables you to distinguish between clitical and non-clitical forms in your search. For instance if you are interested in all clitical wordforms containing the modern lemma *ik* ('I') you

should fill in *ik* at Lemma and choose clitic at Cliticity. Both search queries will be combined, as can be seen in the search query:

Results for: "[lemma="ik"&isclitic="clitic"]" within all documents

This search results in hits such as *hebbic* and *icse*.

### Complete word or word part

This option makes it possible to search for words that are split into two or more parts. Think of separable compound verbs as the infinitive *afsnijden* (‘to cut off’; conjugated form *sneet ... af*, ‘cut off’) and pronominal adverbs as *daeromme* (*daer* + *omme*). Keep in mind that you can only find both parts at the same time using Lemma (*afsnijden*) and the option part. If you are specifically looking for just one of the composing parts (e.g. *af*), you can enter that separate part in Word and click on the option part. In order to find all occurrences with that word part, it is necessary to take into account the different spelling variants of that word part (e.g. *wech*, *wegh*). A possible [search query](#) to find all these forms is:

The screenshot shows a search interface with a dark red header containing 'Search' and 'Explore' tabs. Below the header, the search criteria are displayed under the heading 'Search for ...'. There are four tabs: 'Simple' (selected), 'Extended', 'Advanced', and 'Expert'. The search criteria are as follows:

- Word:** Input field contains 'we\*' with a download icon.
- Lemma:** Input field contains 'weg\*' with a download icon. Below it is a checkbox labeled 'Case- and diacritics-sensitive' which is unchecked.
- Part of Speech:** A dropdown menu showing 'Part of Speech' with an edit icon.
- Numerical Part of Speech code:** Input field contains 'Numerical Part of Speech code' with a download icon.
- Cliticity:** A dropdown menu showing 'Cliticity' with a downward arrow.
- Complete word or word part:** A dropdown menu showing 'part' with a downward arrow.

At the bottom of the search criteria section, there are four buttons: 'Search' (dark red), 'Reset', 'History', and a gear icon for settings.

Below the search criteria, the results are shown as: Results for: "[word="we.\*"&lemma="weg.\*"&iswordpart="part"]" within all documents

## Filter search by

At the right side you will find the option to limit your query to a subset of documents with specific metadata values. You can apply different filters for Date (*Witness Year*, *Text Year*), Localisation (*Region*, *Place*, *Kloeke location Code*), Text type (*Fictionality*, *Genre*, *Subgenre*, *Authenticity*) and Title and author (*Words in title*, *Title of containing publication*, *Author*). To view the results for all documents, simply leave the attributes in the filtering form empty.

## Filter by date

You can select specific dates by entering the same data in the ‘from’ row as in the ‘to’ row (see screenshot below). This can be done for Witness Year, Text Year or even both. If you do not enter a date, the entire corpus is searched. If you want to filter by another date or another period, please press the reset button.

---

### Filter search by ...

Date **1** Localisation Text type Title and author

**Witness Year**

1300 1300

Permissive **Strict**

**Text Year**

From To

Permissive **Strict**

Witness Year (Date): 1300-1300

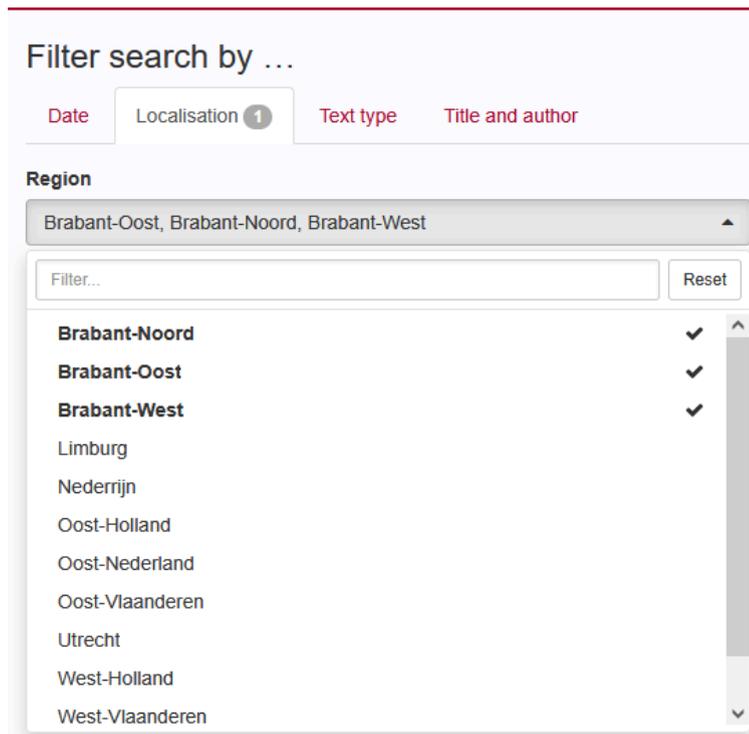
Selected subcorpus:

Total documents: 85 (3.81%)

Total tokens: 29.266 (1.89%)

## Filter by other metadata

There are two different ways to specify a filter, depending on the field type. You can either fill in a value yourself – beside Year above only for Words in Title – or choose one or more values from a drop-down list – for instance Region. The drop-down list has been applied especially when the number of values to choose from is relatively small. Region for instance has only twelve possibilities. You can pick one of these values by clicking on it; your choice will be marked with a tick. It is possible to choose several values. If you want to delete a selection, you can click on the corresponding line again. To close the drop-down list, you can either press the upward pointing arrow in the upper right corner or simply press escape.



By means of a number at the top of Filter search by, the number of values used to filter on, is displayed as can be seen in the screenshot above.

When on the other hand the set of possible values is relatively large (Words in title), you have to type a specific value in a search field. After entering a single character, a list of possible values is suggested. Clicking on an auto-completed value will paste that value in the field. Note that this only works with a single word, like *der*. In order to search for an exact phrase, i.e. a multiple word value, it must be surrounded by double quotes. For instance, in the field Words in title “*Der naturen bloeme*” will result in two manuscripts of this text.

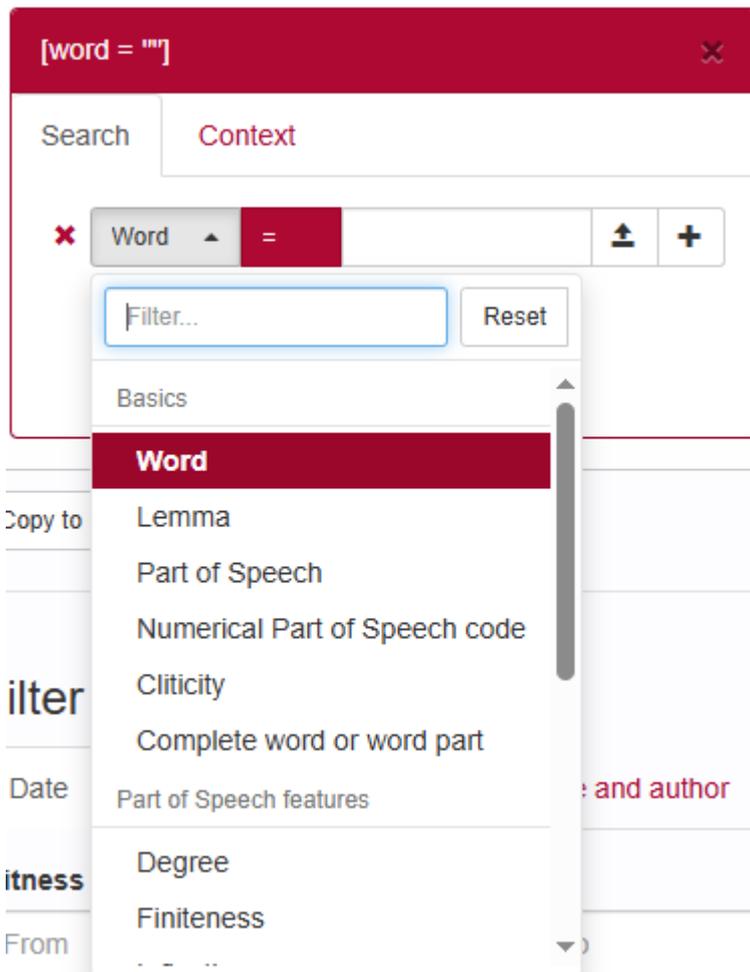
For a detailed description of the metadata, see the section [Metadata categories](#).

## Advanced search

### The query builder

The basic building block in the query builder is the *token box* (see below). Each box represents a token - usually just a single word - or a simple repetition of tokens; when multiple tokens are used, they are matched in order from left to right.

You can use the query builder to create complex queries without writing CQL (here: Corpus Query Language). Therefore, it is easy to use.



A token box in the querybuilder has two tabs: Search and Context.

The tab Search

The tab Search contains a set of attributes a token in the corpus must have to be matched by the query. By clicking the +-button on the right hand side of this token, you can add new attributes (see below). Then enter a value that the attribute must have for the token to be found. The search command Lemma=*lief* & Part of Speech=Common Noun for example excludes all forms of *lief* as an adjective.

The CQL query generated to match this token (the *token query*) in the corpus is displayed in the top bar of the box, to help you understand what is happening internally. The following applies to our example:

The screenshot shows a search query builder window with a red header containing the query "[lemma = 'lief' & pos = 'NOU-C']" and a close button. Below the header are two tabs: "Search" and "Context". Under the "Search" tab, there are two rows of criteria. The first row has a red 'x' icon, a dropdown menu set to "Lemma", an equals sign, the text "lief", and two buttons: a plus-minus icon and a plus icon. Below this row is a checkbox labeled "Case- and diacritics-sensitive". The second row is separated by the word "AND" and has a red 'x' icon, a dropdown menu set to "Part ...", an equals sign, a dropdown menu set to "Common Noun", and the same two buttons. It also has a "Case- and diacritics-sensitive" checkbox. At the bottom center is a plus icon button.

### Token attributes

Specifying token attributes is similar to the Extended Search form. Select which attribute a token should have, and enter the value that the attribute must have for the token to be matched. Attributes in the query builder are interpreted as *regular expressions*. Note that this is different from the Extended Search, where token patterns use wildcards.

Going beyond single-attribute token queries, a token box also allows you to combine several attributes and to specify repetition options.

### Adding attributes to a token box

Using the +-button, new attributes can be added. Two options exist: *AND* and *OR*.

The *AND* option creates a new attribute restriction that a token must match in addition to the ones which were already there. As an example: suppose we want to match *zijn* ('to be') as a verb, not as a pronoun. First, fill in the attribute Lemma with value *zijn*, then click +, choose *AND*, and choose the value Verb for Part of Speech.

The screenshot shows a search query builder window with a red header containing the query "[lemma = 'zijn' & pos = 'VRB']" and a close button. Below the header are two tabs: "Search" and "Context". Under the "Search" tab, there are two rows of criteria. The first row has a red 'x' icon, a dropdown menu set to "Lemma", an equals sign, the text "zijn", and two buttons: a plus-minus icon and a plus icon. Below this row is a checkbox labeled "Case- and diacritics-sensitive". The second row is separated by the word "AND" and has a red 'x' icon, a dropdown menu set to "Part ...", an equals sign, a dropdown menu set to "Verb", and the same two buttons. It also has a "Case- and diacritics-sensitive" checkbox. At the bottom center is a plus icon button.

Similarly, creating a new attribute using *OR* will create a token query matching tokens that have the original attribute or the new attribute. For instance, enter *Word=er* first, add a new attribute with the *OR* option and enter Adverb as Part of Speech to match tokens with Part of Speech tag adverb or with word form equal to *er*.

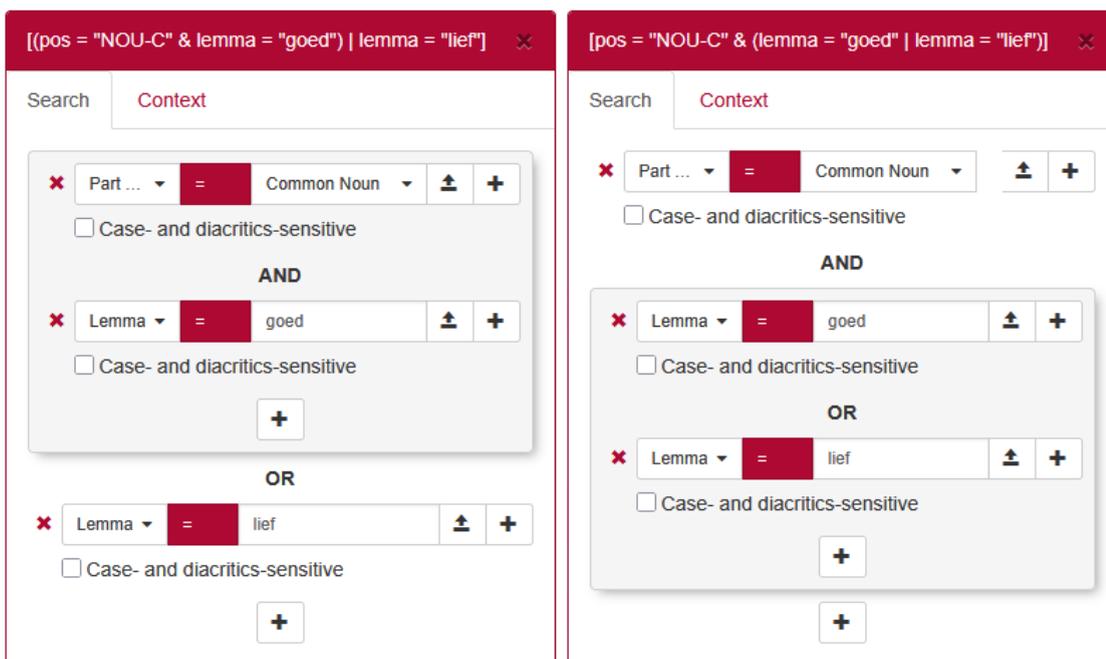
The screenshot shows a search query builder window with a red header containing the query "[lemma = 'er' | pos = 'ADV']" and a close button. Below the header are two tabs: "Search" and "Context". Under the "Search" tab, there are two rows of criteria. The first row has a red 'x' icon, a dropdown menu set to "Lemma", an equals sign, the text "er", and two buttons: a plus-minus icon and a plus icon. Below this row is a checkbox labeled "Case- and diacritics-sensitive". The second row is separated by the word "OR" and has a red 'x' icon, a dropdown menu set to "Part ...", an equals sign, a dropdown menu set to "Adverb", and the same two buttons. It also has a "Case- and diacritics-sensitive" checkbox. At the bottom center is a plus icon button.

### Function of the two +-buttons in a token box

The difference between the +-sign on the right of an attribute and the one below it, is that the +-sign on the right keeps the newly added attribute 'within a subclause'. This is most easily explained by means of an example.

Suppose we want to search for either *goed* or *lief*, used as a noun. If we add the attributes in the order Part of Speech=Common Noun AND Lemma=*goed*, OR Lemma=*lief* using the +-signs **below** the attributes, as in the left screenshot below, we get the token query [(pos = "NOU-C" & lemma = "goed") | lemma = "lief"]. This will also match adjective forms of *lief*, as in *Florens graue van hollant die groet sine lieue scepene ende die ghemeene port van dordrecht*, where *lieue* is an adjective, so this is not what we were after.

If, on the other hand, we add OR lemma=*lief* with the +-sign to the **right** of the attribute Lemma=*goed*, it will be inserted in a subclause (Lemma=*goed* OR Lemma=*lief*), thus resulting in the correct query [pos = "NOU-C" & (lemma = "goed" | lemma = "lief")], as shown in the right screenshot below.



The tab Context

The tab Context specifies the contextual properties, such as whether the token occurs at the end of a sentence, and the repetition pattern:



by hand. This will take you automatically to the Expert Search screen, after which you can start the search or adjust the query if desired.

Copy to CQL editor

## Expert search

The Corpus Query Language (CQL) editor allows you to type your own CQL query, to copy your query into the query builder (in Advanced Search), to import a previously downloaded query and to upload a tab separated list of values to substitute for gap values (see below for further explanation).

CQL queries are expressions built up with the help of a few sequence operators and brackets from basic blocks enclosed by square brackets, in each of which one or more token attributes are specified (these correspond to the token boxes in the query builder).

In CQL, spaces only affect a search if they are included in quotes. Whether the search command is `[word="man"]` or `[ word = "man" ]` (or just "man") does not make any difference to the result. However, there is a difference between the queries `[word="man"]` and `[word=" man"]`. The first search results in almost 3000 hits, but the second one in zero!

Some examples:

- Simple: `[word="schip"]`, e.g. the attribute word matches the regular expression *schip*; `[word!="schip"]`, e.g. the attribute word does **not** match the regular expression *schip*; `[lemma="*.schip"]` matches all lemmata ending with *schip*, including *schip* itself. (Note that `[lemma="*schip"]` will not give any results, because in Expert Search an asterisk is not a wildcard but a repetition operator.)
- Simple sequence: `[pos="PD"][lemma="willen"]` matches all occurrences of the lemma *willen* preceded by a pronoun.
- Combination of attributes (combining operators are &, |, !), e.g. `[word_or_lemma="hope" & pos!="NOU-C"]` or - equivalently - `[word_or_lemma="hope" & !pos="NOU-C"]` matches all occurrences of *hope*, not being a noun.
- Repetition operators: `[pos="AA"]{3}` matches a sequence of 3 adjectives, `[pos="AA"]{2,4}` matches a sequence of 2 to 4 adjectives, `[pos="AA"]{3,}` matches a sequence of 3 or more adjectives.
- The empty [] matches any token, e.g. `[pos="AA"] []{2} [pos="AA"]` matches two adjectives with 2 arbitrary tokens in between.
- Operators |, & and parentheses () and the repetition operators (+, \*, ? and {} ) can be used to build complex sequence queries. Example: `"lieue" "god" | "ons" "here"`, or even `("lieue" "god" | "ons" "here")+`, matching any sequence of *lieue god* or *ons here*. Note that, while most queries up to this point could also have been constructed with the query builder, we really need the power of CQL from here on.

This short list does not cover all CQL features. For more detailed information on how to write CQL, please consult the short [Appendix: Corpus Query Language](#), which contains further pointers.

## Copy to query builder

When the query is relatively simple - like [\[pos="AA"\]](#)[\[lemma="god"\]](#) - it can also be imported into the querybuilder using the *Copy to query builder* button. This will take you automatically to the Advanced Search screen, after which you can start the search or adjust the query if desired.

A message will be displayed next to the button if the query couldn't be parsed.

## Import query

If you have entered a search query, you can find it back by clicking the History button. On the right hand side you can select Download as file in the drop-down menu (default value is Search) and save the file. (For a more elaborate description of the History button see [Simple Search](#).)

Previously saved queries can be used again by uploading them through the Import query button.

## Gap filling

Use this button to upload a Tab Separated Values (TSV) file, which is a simple text format for storing data in a tabular structure. Each record in the table is one line of the text file. Each field value of a record is separated from the next by a tab character. It is also possible to upload a plain text file (.txt) that has the same properties.

A .tsv file or a comparable .txt file enables you to complete a query with marked gaps.

If, for instance, you are interested in the distribution of adjectives you can create this query in the Corpus Query Language field:

```
[lemma="@@"][pos="AA"]
```

By clicking Gap-filling you can upload a file with a tab-separated list of values from your computer to substitute them for the gap values, i.e. the at signs (@@) in your query. After the upload your values will appear in a separate box:

## Search for ...

Simple

Extended

Advanced

Expert

## Corpus Query Language:

```
[lemma="@@"][pos="AA"][lemma="@@"]
```

Copy to query builder

Import query

Gap-filling



```
de god
een vrouw
het dier
```

The values in the first column - *de, een, het* - will be entered at the position of the first gap (@@) and the values in the second column - *god, vrouw, dier* - at the position of the second gap. With these values, gap-filling yields the following results (titles are hidden):

Before Hit	Hit	After Hit	Lemma	Part of Speech with features	Numerical Part of Speech code
...steit ghescreuen van der versteruenessen	<b>eynre</b> <b>eersamer</b> <b>vrouwen</b>	ver zyben wilen was horre...	EEN EERZAAM VROUW	PD(type=indef,subtype=art,position=prenom,inflection=r/re) AA(degree=uncl,position=prenom postnom pred,number=sg,inflection=r/re) NOU-C(number=sg,inflection=n)	485 105 004
...onsen seghelen van der versteruenessen	<b>eynre</b> <b>eersamer</b> <b>vrouwen</b>	wilen was ver zyben horre...	EEN EERZAAM VROUW	PD(type=indef,subtype=art,position=prenom,inflection=r/re) AA(degree=uncl,position=prenom postnom pred,number=sg,inflection=r/re) NOU-C(number=sg,inflection=n)	485 105 004
...pasdach . Dar vor ons quamen .	<b>eyne eersame</b> <b>vrouwe</b>	. ver zybe van tricht . Godeuert...	EEN EERZAAM VROUW	PD(type=indef,subtype=art,position=prenom,inflection=e) AA(degree=uncl,position=prenom postnom pred,number=sg,inflection=e) NOU-C(number=sg,inflection=e)	481 101 001
...wi , In alse selch goet , alse	<b>ene Edele</b> <b>vrowe</b>	, vrowe ver Marie , wilen vrowe...	EEN EDEL VROUW	PD(type=indef,subtype=art,position=prenom,inflection=e) AA(degree=uncl,position=prenom postnom pred,number=sg,inflection=e) NOU-C(number=sg,inflection=e)	481 101 001
...onsen seghelen van der versteruenessen	<b>eynre</b> <b>eersamer</b> <b>vrouwen</b>	wilen was . ver zyben horre...	EEN EERZAAM VROUW	PD(type=indef,subtype=art,position=prenom,inflection=r/re) AA(degree=uncl,position=prenom postnom pred,number=sg,inflection=r/re) NOU-C(number=sg,inflection=n)	485 105 004

This mimics the functionality to upload a list of values in the Extended Search and Advanced Search interfaces.

Please note that for this to work, you do need to enter @@ in the field where you want the substitution to take place. An empty field ([]) will match any term.

## Viewing results

Results can be viewed in two ways: Per hit (hit is defined as one token or a group of tokens that matched the query), or Per document (each document listed contains at least one hit).

### Per Hit view

Click a hit - i.e. a line with the bold word(s) in the column Hit - to display the properties and values of the hit (in the following example **edelre vrouwen**). Click the hit again to close.

Document id	Before Hit	Hit	After Hit	Lemma	Part of Speech with features	Numerical Part of Speech code												
Document id: INT_350b3c93-1717-48c4-a5ba-f7e287b93035																		
Corp.I, 1149, Holland, grafelijke kanselarij, (21 april 1291-4 april 1292) by Unspecified																		
...enen tuist hadden , met eenre	<b>edelre</b>	<b>vrouwen</b>	ende groeter , ver katerinen vrouwe...	EDEL	AA(degree=uncl,position=prenom postnom pred,number=sg,inflection=r/re) C(number=sg,inflection=n)	NOU- 105 004												
<p>...Wi Jan van renjisse , ende didderic here van brederode , maken cont alle den ghenen die desen brief zullen zien iof horen lesen , dat wi enen tuist hadden , met eenre <b>edelre vrouwen</b> ende groeter , ver katerinen vrouwe van vorne , ende burgrauinne van zelant , dien wi , bi onsen consente , ende onsen vrijen wille , ende bi rade onser vrienden , ende onser maghen , vereffent hebben met haer in deser manire , alsoe , dat wi haer helpen zullen ghetrouwelike , ende raden , ende wisen , bi gheloeffder trouwen , tieghens...</p> <table border="1"> <thead> <tr> <th>Property</th> <th>value</th> </tr> </thead> <tbody> <tr> <td>Word</td> <td>edelre</td> </tr> <tr> <td>Word id</td> <td>w.616937</td> </tr> <tr> <td>Lemma</td> <td>EDEL</td> </tr> <tr> <td>Part of Speech with features</td> <td>AA(degree=uncl,position=prenom postnom pred,number=sg,inflection=r/re)</td> </tr> <tr> <td>Numerical Part of Speech code</td> <td>105</td> </tr> </tbody> </table>							Property	value	Word	edelre	Word id	w.616937	Lemma	EDEL	Part of Speech with features	AA(degree=uncl,position=prenom postnom pred,number=sg,inflection=r/re)	Numerical Part of Speech code	105
Property	value																	
Word	edelre																	
Word id	w.616937																	
Lemma	EDEL																	
Part of Speech with features	AA(degree=uncl,position=prenom postnom pred,number=sg,inflection=r/re)																	
Numerical Part of Speech code	105																	

Hit rows are always preceded by a row containing the document title in which those hits occurred, in this case *Corp.I, 1149, Holland, grafelijke kanselarij, (21 april 1291-4 april 1292)*. The document titles can be toggled on or off by using the Hide Titles (or Show Titles when titles are hidden) button at the bottom of the page. If you hover the mouse over the title, the identification number of the document appears, in this case: INT\_350b3c93-1717-48c4-a5ba-f7e287b93035.

### Sorting results

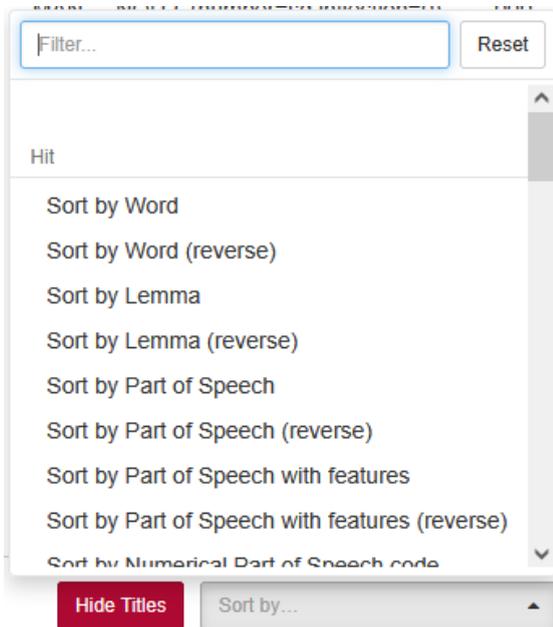
Click on any of the column headings to sort the hits on values within that column, clicking again inverts the sorting. Extra sorting options are given when clicking on Before hit, Hit and After hit: you can sort by Word, Lemma, Part of Speech, Part of Speech with features and Numerical Part of Speech Code, as shown below.

Before Hit ▾	Hit ▲	After Hit ▾	Lemma
Corp.I, 1323, Maastricht, 4 ap			
...steit ghescreuen			ERZAAM
van der	ee		W
versteruenessen	vr		
Corp.I, 1324, Maastricht, 6 ap			
...onsen seghelen .			ERZAAM
van der	ee		W
versteruenessen	vr		
Corp.I, 1333, Maastricht, 9 m			
...paspdach . Dar vor	eyne	eerzame	ERZAAM
ons quamen .	vrouwe	tricht .	VROUW
		Godeuert...	

Sort by...

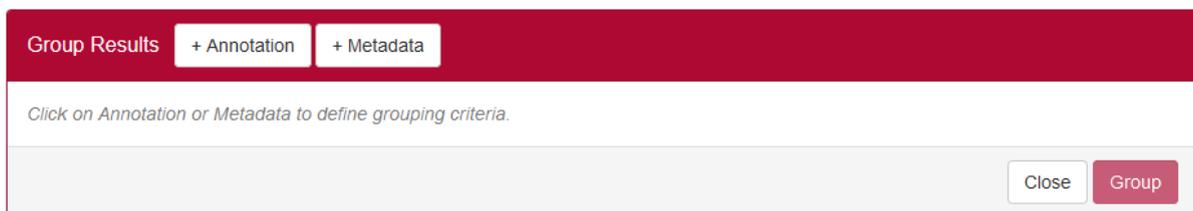
- Basics
- Word
- Lemma
- Part of Speech
- Part of Speech with features
- Numerical Part of Speech code

You can also sort the results by means of the drop-down menu at the bottom of the page (Sort by...), which offers you the possibility to sort by various attributes for Hit, Before hit, After hit, Date, Localisation, Text type and Title and Author.



#### Grouping results

It is possible to group the results by clicking on the button Group Results, after which the following menu appears:



Results can be grouped by Annotation and by Metadata.

By clicking +Annotation you can group by the first word, by all words or by specific words, whether before the hit, within the hit or after the hit, and based on the annotation Word, Lemma, Part of Speech, Part of Speech with features or Numerical Part of Speech Code. When grouping by the first word or specific words, you can also group from the end of the hit. The default grouping is grouping all words within the hit using annotation Word. Clicking +Metadata allows you to group by metadata assigned to the document (Date, Localisation, Text type, Title and Author).

By clicking the Case sensitive box it is possible to distinguish between case sensitive and case insensitive.

The example below is grouped by the first word before the hit. The example dynamically updates when the grouping options are changed.

Group Results + Annotation + Metadata

Word within hit ✕

I want to group on all words within the hit using annotation Word

Case sensitive:

also Scoutete & hebben ghelijc andren **mannen** . Van manne te ontlouene jof  
also Scoutete hebben ghelijc andren Van manne te ontlouene jof

Clear Group

Click a group to show or hide hits within that group, as shown below. Click once more on the group to close it again. If more than twenty hits are found in a document, you can make them appear by clicking on Load more concordances.

Group	#hits in group	Relative frequency (hits)
desen	624	0.0403%
dese	96	0.0062%

« View detailed concordances Load more concordances

Before Hit	Hit	After Hit
...recht ende onser graueleker . dese	<b>brief</b>	word ghegheuen ende ghezegheld tot...
...besegelen met minen segele . dese	<b>brief</b>	waes ghegheuen tote herlem jnde...
...naesten , achte daghen , na dat dese	<b>brief</b>	was ghegheuen , sullen si hem...
...bider vorscreuene , scepenen , seghele . dese	<b>brief</b>	was ghegheuen jnt jaer ons...
...m'ho.cc'ho.bxxx'ho . in hoymaent was dese	<b>brief</b>	ghemaket gheselm vanden houtekine ghaf...
...besegheld met onsen seghel . dese	<b>brief</b>	was ghegheuen vp onser vrouwen...
...dese	<b>brief</b>	es vanden wine vandesen iaere...
...besegheld met onsen seghele . dese	<b>brief</b>	wart ghegheuen in den jaren...
...metter stede heimeleken seghele . dese	<b>brief</b>	was ghegheuen jnden seluen jare...
...daniels ende jans seghele . dese	<b>brief</b>	was ghegheuen , jnt jaer ons...
...min wif dar vp heuet . dese	<b>brief</b>	wart ghescreuen in den jar...
...mijn wijf daer up heuet . dese	<b>brief</b>	wart ghescreuen in den iaere...
...te borghe hebben geset ¶ dese	<b>brief</b>	wart gegheuen in sinte bartholomeus...
...houden . aldusgedane vorrewarde also dese	<b>brief</b>	houdet . dar dese dur gesteken...
...met onser stede seghele . dese	<b>brief</b>	was ghegheuen jnt jaer vanden...
...van leyden onse clerc . dese	<b>brief</b>	was ghegheuen jn die haghe...
...jn orkonde minen seghele . dese	<b>brief</b>	was ghegheuen jnt jaer ons...
...bezegheilen met onsen seghele . dese	<b>brief</b>	es ghegheuen en den jare...
...besegheilen met onsen seghele . dese	<b>brief</b>	es ghegheuen in den iaere...
...besegheld met onsen seghele . dese	<b>brief</b>	wart ghegheuen tote broecele jn...

« View detailed concordances Load more concordances

den	19	0.00123%
openen	17	0.0011%
iegghenwordighen	14	0.000904%
tieghenwordighen	14	0.000904%

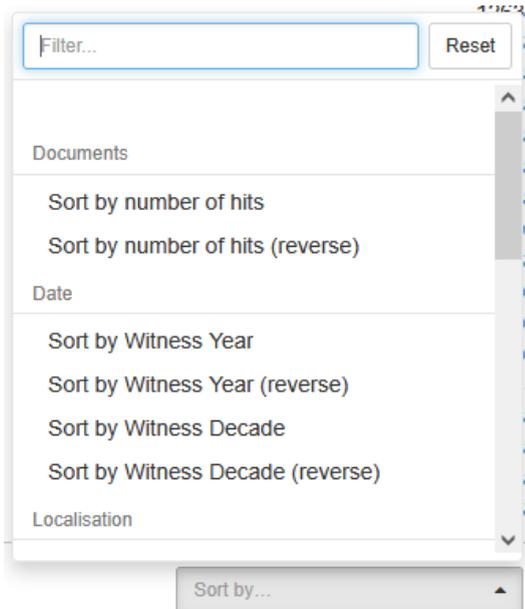
Click on View detailed concordances to go back to the normal hits view to see more detailed information for the hits in this group. The button Go back to grouped results brings you back to the list of groups.

## Per Document view

### Sorting results

Click on any of the column headings to sort the documents by Document (name), Witness year or Hits within that column, clicking again inverts the sorting.

You can also sort the results by means of the drop-down menu at the bottom of the page (Sort by...), which offers you the possibility to sort by various attributes such as Hit (Documents), Witness Date (Date) and Genre (Text type).



## Grouping results

Results Per Document can be grouped by metadata assigned to the document (Date, Localisation, Text type, and Title and author). The example below shows all documents in which the Lemma *penning* occurs grouped by Witness Decade.

Results for: [lemma="penning"] within all documents

---

Per Hit | Per Document

Documents / Grouped by Document Witness Decade Total documents: 829 (37.2%)  
Total groups: 11  
Search time: 0.002s

Group Results + Metadata

Document Witness Decade ✕

Select the metadata to group on.

Group by Witness Decade ▾

Case sensitive:

Clear Group

« 1 » table docs

Group	#docs in group	Relative frequency (docs)
1290	406	49%
1280	283	34.1%
1270	69	8.32%
1300	37	4.46%
1260	23	2.77%
1250	6	0.603%
1230	2	0.241%
1310	1	0.121%
1320	1	0.121%
1240	1	0.121%
1380	1	0.121%

## Exporting results

The search results - both Per hit as Per document - can be exported by using the Export or the Export for Excel button at the bottom right of the page. The first button transfers the search results - including all metadata - to a Comma-Separated Values-file. These CSV-files consist only of text data, which makes it easy to implement (read and/or write) them into a spreadsheet or database program. The second button offers the possibility to export the results - including all metadata - to a CSV-file for use with Excel.

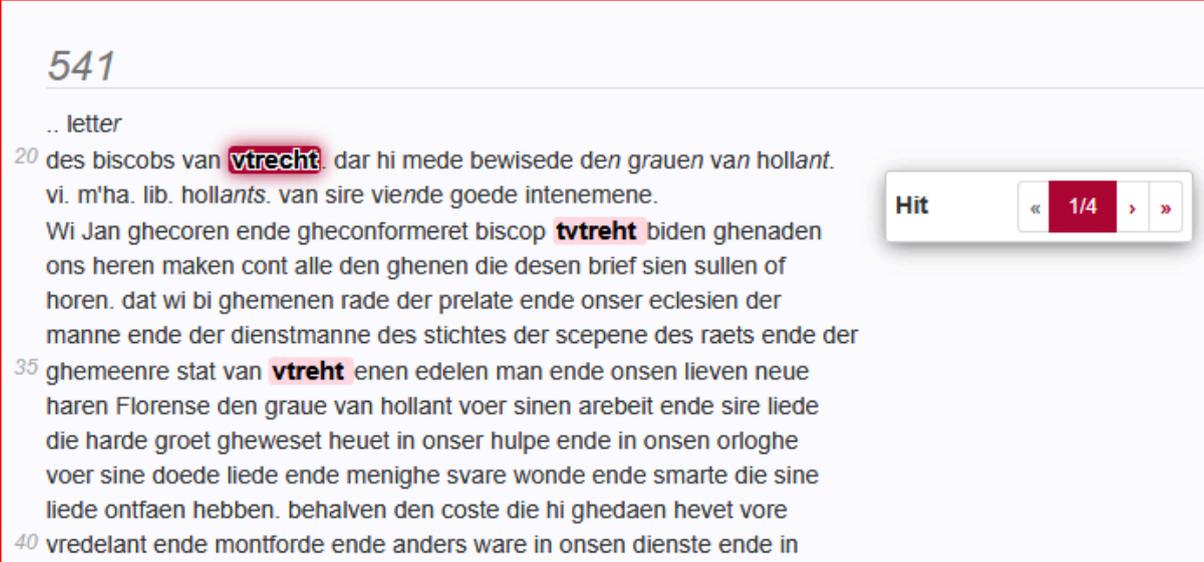
Grouped results can be exported in the same way. However, if you would like to have the metadata with each concordance of a group, you must first click on the red bar of a specific group and then on View detailed concordances. The results you then see can be exported by the use of the Export buttons. This operation must be carried out for each individual group you wish to export.

## Information about a document

Click on a document title or the chain icon in the per hit view to open this document in a new window: the Content window.

### Content

Hits from the current query will be highlighted in bold in the opened document. In the case of several hits only the current hit will also appear in shadow (such as *vtrecht* in the example below). You can navigate from one hit to another by using the arrows at the Hits button (this button can be dragged around), and you can also browse through the pages if the document consists of more than one page.



541

.. letter

20 des biscobs van **vtrecht**, dar hi mede bewisede den grauen van hollant.  
vi. m'ha. lib. hollants. van sire viende goede intenemene.

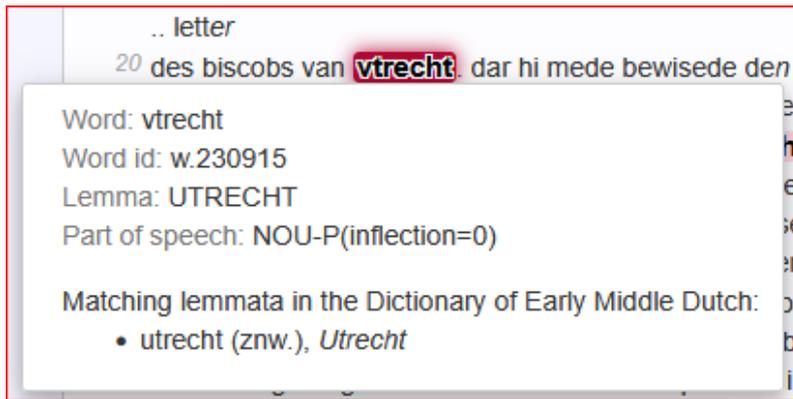
Wi Jan ghecoren ende gheconformeret biscop **vtrecht** biden ghenaden  
ons heren maken cont alle den ghenen die desen brief sien sullen of  
horen. dat wi bi ghemenen rade der prelate ende onser eclesien der  
manne ende der dienstmanne des stichtes der scepene des raets ende der

35 ghemeenre stat van **vtrecht** enen edelen man ende onsen lieuen neue  
haren Florense den graue van hollant voer sinen arebeit ende sire liede  
die harde groet gheweset heuet in onser hulpe ende in onsen orloghe  
voer sine doede liede ende menighe svare wonde ende smarte die sine  
liede ontfaen hebben. behalven den coste die hi ghedaen hevet vore

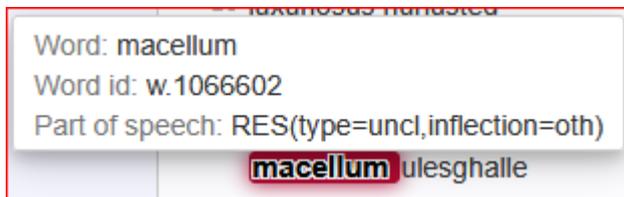
40 vredelant ende montforde ende anders ware in onsen dienste ende in

Hit « 1/4 » »

When you hover with your mouse over a specific word in the document a pop-up will appear with the modern lemma in capitals and the option Show details. By clicking this link you will see extra information on word level:



Note that for the foreign (Latin) words there is no option for details, as these words do not have a modern lemma in the Dictionary of Early Middle Dutch. Here, all the information is directly shown by hovering over or clicking on the specific word with your mouse:



## Metadata

In the Metadata tab all metadata properties of the document are displayed.

## Statistics

The Statistics tab shows several document statistics: the number of Tokens, Types (unique word forms), Lemmas, Part of Speech and the Type/Token ratio. It is possible to print or to download these statistics via the menu symbol right of the title Token/Part of Speech Distribution respectively the title Vocabulary Growth.

## Exploring the corpus

The Explore tab has three subdivisions: Documents, N-grams and Statistics.

### Documents

This subtab allows you to investigate the documents. It consists of two drop-down menus to specify the grouping of the metadata and to specify the way the groups are to be shown.

A simple example: suppose we want to obtain information about medieval books on medicine within the *Corpus Gysseling*.

- In the Group documents by metadata drop-down menu, choose Group by Subgenre (Text type)
- In Show groups as, select *Docs*
- In the metadata search form (Filter search by), select in Subgenre (Text type) *medicine*
- Press Search

Search Explore

## Explore ...

Documents **N-grams** Statistics

**Group documents by metadata** Group by Subgenre

**Show groups as** Docs

---

## Filter search by ...

Date Localisation **Text type 1** Title and author

**Fictionality** Fictionality

**Genre** Genre

**Subgenre** medicine

You will get this result:

Per Hit Per Document

Documents / Grouped by Document Subgenre

Total documents: 2 (100%)  
Total groups: 1  
Search time: 0.001s

Group Results + Metadata

Document Subgenre ✕

Select the metadata to group on.  
Group by Subgenre

Case sensitive:

Clear Group

« 1 » table docs tokens

Group	#docs in group	Relative frequency (docs)
artes · medicine · secular	2	100%

## N-grams

An *N-gram* is a sequence of *N* items. This option will list the frequency of different N-grams in a (sub-)corpus.

## Options

- *N-gram size*: the length of the sequence (a number from 1 to 5; default setting is 5).

- *N-gram type*: the attribute to search for. You can choose: Word (i.e. word form), Lemma, Part of Speech, Part of Speech with features or Numerical Part of Speech code. If you do not specify the search term further, a series of arbitrary tokens equal to the n-gram size will be searched for.
- It is also possible to restrict to, for instance, n-grams with some slots already specified, as is shown in the following example.
- By using the Filter search by... you can create a subcorpus within the *Corpus Gysseling* for specific metadata.

## Example

The screenshot shows the 'Explore' interface with the following settings:

- Documents** (selected), N-grams, Statistics
- N-gram size**: 5
- N-gram type**: Word
- Filters: Lemma (ik), Part of Speech (Verb), Part of Speech (Verb), Word (Word), Word (Word)

Within all the documents of the *Corpus Gysseling*, you will find more than 900 occurrences of this so-called 5-gram:

The screenshot shows the search results page with the following details:

- Per Document** (selected), Per Hit
- Hits** / Grouped by Word within hit
- Total hits: 984 (0.0584%)
- Total groups: 871
- Search time: 0.02s
- Group Results** (selected), + Annotation, + Metadata
- Word within hit:
- I want to group on: all words, within the hit, using annotation, Word
- Case sensitive:
- Text snippet: `sullen sien ende horen dat ic hebbe ghedoget bi minen goeden wille dat een leec man`
- Navigation: 1, 2, 3, 4, 6, 11, table, hits

Group	#hits in group	Relative frequency (hits)
ics vermaent worde van hem	6	0.000388%
ic belouet hebbe ende beloue	4	0.000258%
ic vermach dade ic dat	4	0.000258%
ic gheset hebbe ende ghegheuen	3	0.000194%
ic hebbe ontfanghen jn lene	3	0.000194%
ic ghelouet hebbe vor mi	3	0.000194%
ic doen mach Ende dese	2	0.000129%
ic ghehad hebbe leghens minen	2	0.000129%
ic verboert hebbe mijn lijf	2	0.000129%
mi comen sal So ebbe	2	0.000129%
Jc ben gesendet hir uan	2	0.000129%
ic ghelouet hebbe minen lieuen	2	0.000129%
mi beuolen hebben Ende om dat	2	0.000129%
hic hebbe vercocht alse van	2	0.000129%
mi gheproeuen mach met den	2	0.000129%
ic salne don op herstaen in	2	0.000129%
mi helpen moghen ofte mochten	2	0.000129%
ic hebbe mesdaen vor Gode	2	0.000129%
ic ghelouet hebbe ende gheloue	2	0.000129%
ic vermach dadic dat soe	2	0.000129%

## Statistics (frequency lists)

Here, you can produce frequency lists for the corpus. It is rather similar to the previous option, but restricted to 1-grams.

### Options

- *Frequency list type*: choose for lists of Word (i.e. word form), Lemma, Part of Speech, Part of Speech with features or Numerical Part of Speech code.
- By using the Filter search by... you can create a subcorpus within the *Corpus Gysseling* for specific metadata.

### Example

It is possible to determine the use of the most frequently used words in *Het Luikse Diatesseron* in the *Corpus Gysseling* by searching for Frequency list type Word and by filtering search by Title of containing publication: *Het Luikse Diatesseron*. This results in:

Results for: Word frequency within documents where Title of containing publication: Het Luikse Diatesseron

Per Hit **Per Document**

Hits / Grouped by Word within hit Total hits: 69.269 (100%)  
Total groups: 7.506  
Search time: 0.04s

Group Results + Annotation + Metadata

Word within hit

I want to group on all words within the hit using annotation Word

Case sensitive:

 SCHONE HISTORIE VAN DEN WESENE  
SCHONE HISTORIE VAN DEN WESENE

Clear Group

« 1 2 3 4 6 11 » table hits

Group	#hits in group	Relative frequency (hits)
ende	3.660	5.28%
die	1.885	2.72%
dat	1.868	2.7%
hi	1.276	1.84%
van	1.249	1.8%
den	1.095	1.58%
in	1.058	1.53%
so	990	1.43%
Ende	856	1.24%
de	836	1.21%
en	779	1.12%
hem	775	1.12%
te	722	1.04%
si	684	0.987%
Ihesus	643	0.928%

# Appendix: Corpus Query Language

BlackLab supports Corpus Query Language, a full-featured query language introduced by the IMS Corpus WorkBench (CWB) and also supported by the Lexicom Sketch Engine. It is a standard and powerful way of searching corpus.

The basics of Corpus Query Language is the same in all three projects, but there are a few minor differences in some of the more advanced features, as well as some features that are exclusive to some projects. For most queries however, this will not be an issue.

This page will introduce the query language and show all features that BlackLab supports. If you want to learn even more about CQL, see [CWB CQP Query Language Tutorial](#) and [Sketch Engine Corpus Query Language](#).

## CQL support

For those who already know CQL, here's a quick overview of the extent of BlackLab's support for this query language. If there is a feature we don't support, yet is important to you, please let us know. If it's quick to add, we may be able to help you out.

## Supported features

BlackLab currently supports (arguably) most of the important features of Corpus Query Language:

- Matching on token annotations (also called properties or attributes), using regular expressions and =, !=, !. Example: [word="bank"] (or just "bank")
- Case/accent-sensitive matching. Note that, unlike in CWB, case-INsensitive matching is currently the default. To explicitly match case/accent-insensitivity, use "(?i)...". Example: "(?i)Mr\." "(?i)Banks"
- Combining criteria using &, | and !. Parentheses can also be used for grouping. Example: [lemma="bank" & pos="V"]
- Match-all pattern [] matches any token. Example: "a" [] "day"
- Regular expression operators +, \*, ?, {n}, {n,m} at the token level. Example: [pos="AA"]+
- Sequences of token constraints. Example: [pos="AA"] "cow"
- Operators |, & and parentheses can be used to build complex sequence queries. Example: "happy" "dog" | "sad" cat"
- Querying with tag positions using e.g. <s> (start of sentence), </s> (end of sentence), <s/> (whole sentence) or <s> ... </s> (equivalent to <s/> containing ...). Example: <s> "The" . XML attribute values may be used as well, e.g. <ne type="PERS"/> ("named entities that are persons").
- Using within and containing operators to find hits inside another set of hits. Example: "you" "are" within <s/>
- Using an anchor to capture a token position. Example: "big" A:[]. Captured matches can be used in global constraints (see next item) or processed separately later (using the Java interface; capture information is not yet returned by BlackLab Server). Note that BlackLab can actually capture entire groups of tokens as well, similarly to regular expression engines.

- Global constraints on captured tokens, such as requiring them to contain the same word.  
Example: "big" A:[] "or" "small" B:[] :: A.word = B.word

See below for features not in this list that may be added soon, and let us know if you want a particular feature to be added.

## Differences from CWB

BlackLab's CQL syntax and behaviour differs in a few small ways from CWBs. In future, we'll aim towards greater compliance with CWB's de-facto standard (with some extra features and conveniences).

For now, here's what you should know:

- Case-insensitive search is currently the default in BlackLab, although you can change this if you wish. CWB and Sketch Engine use case-sensitive search as the default. We may change our default in a future major version.  
If you want to switch case-/diacritics-sensitivity, use "(?-i).." (case-sensitive) or "(?i).." (case-insensitive, usually the default). CWBs %cd flags for setting case/diacritics-sensitivity are not (yet) supported, but will be added.
- If you want to match a string literally, not as a regular expression, use backslash escaping: "e.g.". %l for literal matching is not yet supported, but will be added.
- BlackLab supports result set manipulation such as: sorting (including on specific context words), grouping/frequency distribution, subsets, sampling, setting context size, etc. However, these are supported through the REST and Java APIs, not through a command interface like in CWB. See [BlackLab Server overview](#)).
- Querying XML elements and attributes looks natural in BlackLab: <s/> means "sentences", <s> means "starts of sentences", <s type='A'> means "sentence tags with a type attribute with value A". This natural syntax differs from CWBs in some places, however, particularly when matching XML attributes. While we believe our syntax is the superior one, we may add support for the CWB syntax as an alternative.  
We only support literal matching of XML attributes at the moment, but this will be expanded to full regex matching.
- In global constraints (expressions occurring after ::), only literal matching (no regex matching) is currently supported. Regex matching will be added soon. For now, instead of A:[] "dog" :: A.word = "happy|sad", use "happy|sad" "dog".
- To expand your query to return whole sentences, use <s/> containing (...). We don't yet support CWBs expand to, expand left to, etc., but may add this in the future.
- The implication operator -> is currently only supported in global constraints (expressions after the :: operator), not in regular token constraints. We may add this if there's demand for it.
- We don't support the @ anchor and corresponding target label; use a named anchor instead. If someone makes a good case for it, we will consider adding this feature.
- backreferences to anchors only work in global constraints, so this doesn't work: A:[] [] [word = A.word]. Instead, use something like: A:[] [] B:[] :: A.word = B.word. We hope to add support for these in the near future, but our matching approach may not allow full support for this in all cases.

## (Currently) unsupported features

The following features are not (yet) supported:

- intersection, union and difference operators. These three operators will be added in the future. For now, the first two can be achieved using & and | at the sequence level, e.g. "double" [] & [] "trouble" to match the intersection of these queries, i.e. "double trouble" and "happy" "dog" | "sad "cat" to match the union of "happy dog" and "sad cat".
- \_ meaning "the current token" in token constraints. We will add this soon.
- lbound, rbound functions to get the edge of a region. We will probably add these.
- distance, distabs functions and match, matchend anchor points (sometimes used in global constraints). We will see about adding these.
- using an XML element name to mean 'token is contained within', like [(pos = "N") & !np] meaning "noun NOT inside in an tag". We will see about adding these.
- a number of less well-known features. If people ask, we will consider adding them.

## Using Corpus Query Language

### Matching tokens

Corpus Query Language is a way to specify a "pattern" of tokens (i.e. words) you're looking for. A simple pattern is this one:

```
[word="man"]
```

This simply searches for all occurrences of the word "man". If your corpus includes the per-word properties lemma (i.e. headword) and pos (part-of-speech, i.e. noun, verb, etc.), you can query those as well. For example, to find a form of word "search" used as a noun, use this query:

```
[lemma="search" & pos="NOU-C"]
```

This query would match "search" and "searches" where used as a noun. (Of course, your data may contain slightly different part-of-speech tags.)

The first query could be written even simpler without brackets, because "word" is the default property:

```
"man"
```

You can use the "does not equal" operator (!=) to search for all words except nouns:

```
[pos != "NOU-C"]
```

The strings between quotes can also contain wildcards, of sorts. To be precise, they are [regular expressions](#), which provide a flexible way of matching strings of text. For example, to find "man" or "woman", use:

```
"(wo)?man"
```

And to find lemmata starting with "under", use:

```
[lemma="under.*"]
```

Explaining regular expression syntax is beyond the scope of this document, but for a complete overview, see [here](#).

## Sequences

Corpus Query Language allows you to search for sequences of words as well (i.e. phrase searches, but with many more possibilities). To search for the phrase "the tall man", use this query:

```
"the" "tall" "man"
```

It might seem a bit clunky to separately quote each word, but this allows us the flexibility to specify exactly what kinds of words we're looking for. For example, if you want to know all single adjectives used with man (not just "tall"), use this:

```
"an?|the" [pos="AA"] "man"
```

This would also match "a wise man", "an important man", "the foolish man", etc.

## Regular expression operators on tokens

Corpus Query Language really starts to shine when you use the regular expression operators on whole tokens as well. If we want to see not just single adjectives applied to "man", but multiple as well:

```
"an?|the" [pos="AA"]+ "man"
```

This query matches "a little green man", for example. The plus sign after [pos="AA"] says that the preceding part should occur one or more times (similarly, \* means "zero or more times", and ? means "zero or one time").

If you only want matches with two or three adjectives, you can specify that too:

```
"an?|the" [pos="AA"] {2,3} "man"
```

Or, for two or more adjectives:

```
"an?|the" [pos="AA"] {2,} "man"
```

You can group sequences of tokens with parentheses and apply operators to the whole group as well. To search for a sequence of nouns, each optionally preceded by an article:

```
("an?|the"? [pos="NOU-C"])+
```

This would, for example, match the well-known palindrome "a man, a plan, a canal: Panama!"

## Punctuation

In BlackLab, punctuation tends to not be indexed as a separate token, but as a property of a word token - CWB and Sketch Engine on the other hand tend to index punctuation as a separate token instead. You certainly could choose to index punctuation as a separate token in BlackLab, by the way -- it's just not commonly done. Both approaches have their advantages and disadvantages, and of course the choice affects how you write your queries.

It is possible to search for punctuation marks. E.g. to find occurrences of the word "want" preceded by a comma use the following query:

```
[punctBefore=", " & word="want"]
```

To find occurrences of the lemma "krant" that are followed by an exclamation mark, use:

```
[lemma="krant" & punctAfter="!"]
```

Some punctuation marks have a special function in regular expressions and therefore must be preceded by a backslash (\) when used in queries. For instance, to search for a period (.) after the word "geweest", use:

```
[word="sentence" & punctAfter="\."]
```

## Case- and diacritics-sensitivity

CWB and Sketch Engine both default to (case- and diacritics-)sensitive search. That is, they exactly match upper- and lowercase letters in your query, plus any accented letters in the query as well. BlackLab, on the contrary, defaults to \*IN\*sensitive search (although this default can be changed if you like). To match a pattern sensitively, prefix it with "(?-i)":

```
"(?-i) Panama"
```

If you've changed the default search to sensitive, but you wish to match a pattern in your query insensitively, prefix it with "(?i)":

```
[pos="( ?i) NOU-C"]
```

Although BlackLab is capable of setting case- and diacritics-sensitivity separately, it is not yet possible from Corpus Query Language. We may add this capability if requested.

## Matching XML elements

Corpus Query Language allows you to find text in relation to XML elements that occur in it. For example, if your data contains sentence tags, you could look for sentences starting with "the":

```
<s>"the"
```

Similarly, to find sentences ending in "that", you would use:

```
"that" </s>
```

You can also search for words occurring inside a specific element. Say you've run named entity recognition on your data and all person names are surrounded with <person>...</person> tags. To find the word "baker" as part of a person's name, use:

```
"baker" within <person/>
```

Note the forward slash at the end of the tag. This way of referring to the element means "the whole element". Compare this to <person>, which means "the element's open tag", and </person>, which means "the element's close tag".

The above query will just match the word "baker" as part of a person's name. But you're likely more

interested in the entire name that contains the word "baker". So, to find those full names, use:

```
<person/> containing "baker"
```

Or, if you simply want to find all persons, use:

```
<person/>
```

As you can see, the XML element reference is just another query that yields a number of matches. So as you might have guessed, you can use "within" and "containing" with any other query as well. For example:

```
([pos="AA"]+ containing "tall") "man"
```

will find adjectives applied to man, where one of those adjectives is "tall".

## Labeling tokens, capturing groups

Just like in regular expressions, it is possible to "capture" part of the match for your query in a "group".

CWB and Sketch Engine offer similar functionality, but instead of capturing part of the query, they label a single token. BlackLab's functionality is very similar but can capture a number of tokens as well. For example:

```
"an?|the" Adjectives:[pos="AA"]+ "man"
```

This will capture the adjectives found for each match in a captured group named "Adjectives".

BlackLab also supports numbered groups:

```
"an?|the" 1:[pos="AA"]+ "man"
```

## Global constraints

If you tag certain tokens with labels, you can also apply "global constraints" on these tokens. This is a way of relating different tokens to one another, for example requiring that they correspond to the same word:

```
A:[ ] "by" B:[ ] :: A.word = B.word
```

This would match "day by day", "step by step", etc.