

Corpus Gysseling application manual

Table of contents

Introduction	4
Information about the corpus	4
Lemmatization	4
Part of speech tagging	4
Metadata categories	5
Date	5
Witness Year	5
Text Year	5
Permissive / strict	5
Localization	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
Application user manual	8
Getting started	8
Searching the corpus	9
Simple Search	9
Search	9
Wildcards	9
Reset	10

History	10
Global settings	10
Extended Search	12
Wildcards	13
Upload a list of values	13
Part of speech dialog box	14
Numerical part of speech code	14
Cliticity	14
Complete word or word part	15
Starting a new search	15
Filter search by	15
Advanced search	17
The query builder	17
The tab search	17
Adding attributes to a token box	18
Function of the two +-buttons in a token box	19
The tab options	20
Managing sequences of token boxes	20
Uploading value lists in the query builder	21
Copy to CQL editor	21
Expert search	21
Copy to query builder	22
Import query	22
Gap filling	22
Viewing results	24
Per Hit view	25
Sorting results	25
Grouping results	26
Per Document view	28
Sorting results	28
Grouping results	28
Exporting results	28
Information about a document	28
Content	28

Metadata of a document	29
Statistics	29
Exploring the corpus	29
Documents	30
N-grams	31
Options	31
Example	31
Statistics (frequency lists)	31
Options	32
Example	32
Appendix: Corpus Query Language	33
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
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 INT. The backend of the application is the BlackLab Lucene based search engine developed for corpora with token-based annotation (<http://inl.github.io/BlackLab/>). The web-based frontend is a further development of the corpus-frontend application developed by INT (<https://github.com/INL/corpus-frontend>) in CLARIN and CLARIAH projects. Its design is inspired by the first version of the OpenSoNaR user interface by Tilburg 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.

Lemmatization

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 lemmatization principles can be found in Marijke Mooijaart, [Het lemma in the GiGaNT lexicon](#).

Part of speech tagging

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](#).

In the context of the CLARIAH+ project, a tagset and tagging principles for the annotation of diachronic corpora of historical Dutch has been developed. 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](#).

Metadata categories

The *Corpus Gysseling* has been enriched with an elaborate set of metadata categories. These metadata will all be 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 came into being 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.)

Localization

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. 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 (see the paragraph [Filter search by](#)). 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 Heinric van Veldeke.

Application user manual

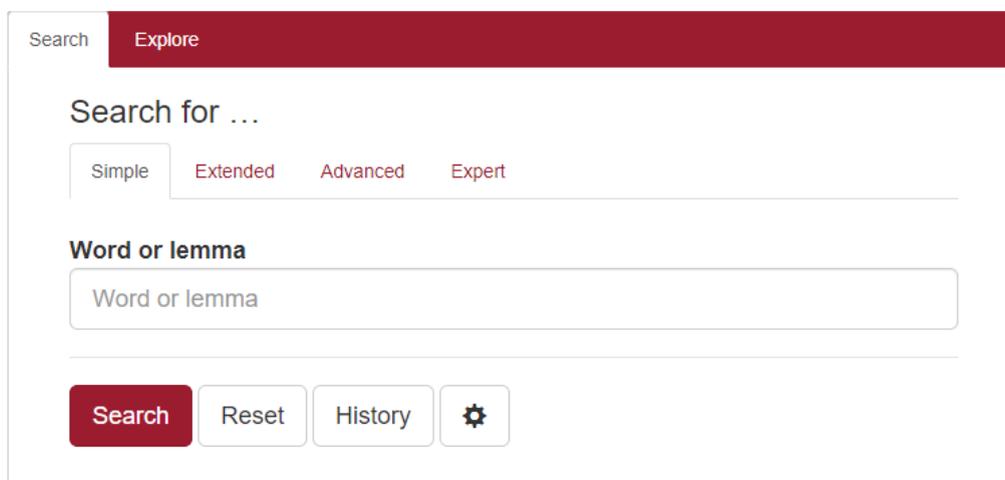
Getting started

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 the attribute Word in Extended Search.
 - Simple Search for [scepe](#)
 - Extended Search for Word [dochter](#)
- To search by lemma form (i.e. the canonical form or citation form of a set of forms (headform)), you can use Simple Search or else the attribute Lemma in Extended Search.
 - Simple Search for [leven](#)
 - Extended Search for Lemma [leven](#)
- To search for words 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 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 the Group hits by feature.
 - example Group by Word: [different adjectives before the lemma huis](#)
 - example Group by Lemma before: [lemmata preceding the word god](#)
- To explore the distribution of document properties in the corpus, use the Explore feature
 - example: [characteristics about the date](#)
 - example: [genres 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. *ship*). It is also possible to enter a phrase: *het ship* or *dit zijn de schepen*. To start the search simply press enter or press the Search button.

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

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 (see over there).

Note that in Simple Search the patterns will be matched case-insensitively: *ship* will deliver the same results as *ship* 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.

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* in Word or lemma 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*), *anden* (lemma *aan+de/het*)

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

Global settings ×

Results per page:

Sample size:

Seed:

Context size:

Wide View

Extended Search

The Extended Search allows you to find all occurrences of a *token* with its specific *attributes*. A *token* - usually just a single word - is the smallest unit within a corpus, whereas *attributes* are the different values that together make up a token.

In this corpus the six 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. All supported attributes are shown in the search form:

Search **Explore**

Search for ...

Simple Extended **Advanced** Expert

Word Word Case- and diacritics-sensitive

Lemma Lemma Case- and diacritics-sensitive

Part of speech Part of speech

Numerical part of speech code Numerical part of speech code

Cliticity Cliticity

Complete word or word part Complete word or word part

Filter search by ...

Date Localization **Text type** Title and author

Witness Year From To Permissive **Strict**

Text Year From To Permissive **Strict**

Selected subcorpus:
Total documents: 2,231 (100%)
Total tokens: 1,547,893 (100%)

Search Reset History

In the search fields Word, Lemma and Numerical part of speech code enter the value of the attributes (or Upload a list of values; see below) you are looking for. In the search fields Part of speech, Cliticity and Complete word or word part select the desired values. Then press enter or click the Search button below to execute the search and view the results. Note that the default setting for Word and Lemma in Extended search is case- and diacritics-insensitive. For example, searching for the Word *Maria* will result in 127 occurrences of this name. By ticking the box Case-sensitive in - for instance - Group by Word in Results you will not only find the Word *Maria* (68x), but also the variants *maria* (58x) and *MARIA* (1x). In order to directly find only occurrences of the Word (form) *Maria*, tick the box Case- and diacritics-sensitive under the search field Word (as shown below).

Search for ...

Simple

Extended

Advanced

Expert

Word

Maria

 Case- and diacritics-sensitive

Lemma

Lemma

 Case- and diacritics-sensitive

Please note that there is an important difference between the search fields Word and Lemma. As an example: entering the value *ship* in Word will only provide you with occurrences of that exact string of characters. When you enter *ship* in the search field Lemma you will - besides the lemma *ship* - also find all word forms that are linked to that lemma, such as the spelling variants *schep*, *scep*, *scip*, and inflected forms as *scepen*, *scepe* and *sceps*.

Wildcards

In Extended 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* in Word matches all word forms that start with an *a* and end with a *n*, e.g. *an*, *Allen*, *alsmen*, *anden* and *andren*. Note that the same query in Lemma will give other results.
- ? The question mark matches a single character once. Therefore, searching for *a?n* in Lemma matches *only* three-letter values starting with an *a* and ending with a *n*, e.g. *aan*.

This wildcard can be used more than once. Thus *a???n* (in Word) matches *allen*, *anden*, *adden*, *armen* en *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.]

In the search fields Word and Lemma it is possible to 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|aan*|hond*.

For the search field Word it is also possible to search for a series of tokens by entering multiple values - including wildcards - separated by a space, e.g. *goede man*, *goede ** or ** man*. Note that searching for *goed man*, *goed ** and ** man* in the search field Lemma will give different results!

Values at the same position in different fields are grouped together as a single token, meaning that all values in the first position of each field are grouped to match a single token.

- A single-token example: searching for the Word(form) *man* together with Part of speech NOU-C (*number=sg*, *inflection=0*) will result in a list of all occurrences of the non-inflected,

singular noun *man*. The syntax of your query is shown in the results: [word="man"&pos="nou-c"&pos_number="sg"&pos_inflection="0"].

- A multi-token example: searching for *ieghen redene* in the Word(form) field and *jegen reden* in the Lemma field should find those occurrences of the bigram in which the first word is the spelling variant of the proposition *jegen* and the second word the declined form of the common noun *reden*.

Upload a list of values

At the right sight 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 dialog box

Clicking on the pencil next to the search field Part of speech provides you with the Part of speech dialog box.

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, queden, 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. A detailed description of the encoding can be found [here](#).

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. *aue*, *of*).

Starting a new search

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

The search fields Word, Lemma and Numerical part of speech code are provided with a list, which contains suggestions for possible search terms in alphabetical order, based on the characters typed in.

If you use the fields Word, Lemma or Numerical part of speech code, there are two possibilities to start a search: fill in the desired value and then press enter, or click the Search button. The only way to start a new search after a change in Part of speech, Cliticity or Complete word or word part is to click the Search button.

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, Localization, Text type and Title and author, which will be explained hereafter. (To view the results for all documents simply leave the attributes in the filtering form empty.)

By means of a number at the top of Filter search by, the number of values used to filter on, is displayed:

Filter search by ...

Date Localization **1** Text type Title and author Metadata

Region

Brabant-Noord, Brabant-Oost, Brabant-West

Place

Place

Kloeke location code

Kloeke location code

Region (Localization): "Brabant-Noord", "Brabant-Oost", "Brabant-West"

Selected subcorpus:
 Total documents: 225 (10.1%)
 Total tokens: 308.729 (19.9%)

There are two different ways to specify a filter, depending on the field type. You can either fill in a value yourself - for instance Date Witness Year - 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. 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.)

When on the other hand the set of possible values is relatively large (e.g. Title and author: 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.

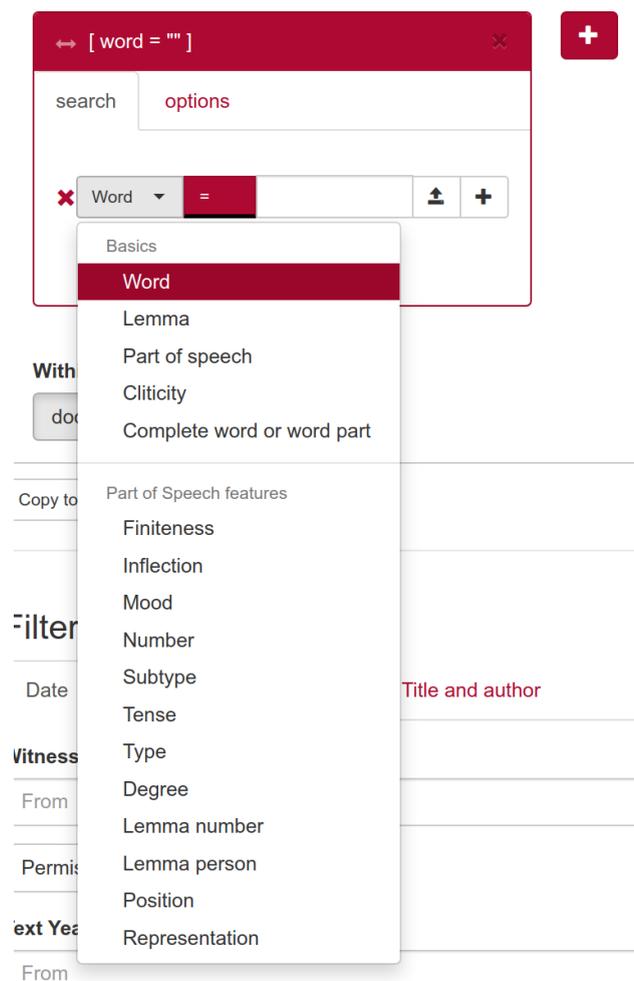
For a detailed description of the metadata, see the section Metadata categories at the beginning of this manual.

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 options.

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* and 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 query builder window with a red header bar containing the query: `[lemma = "lief" & pos = "nou\c"]`. Below the header, there are two tabs: 'search' and 'options'. The 'options' tab is active. The main area contains two conditions stacked vertically, connected by an 'AND' operator. The first condition is 'Lemma = lief', and the second is 'Part of s = Common noun'. Each condition has a red 'X' icon on the left and a '+' icon on the right. A '+' button is also located below the second condition.

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.

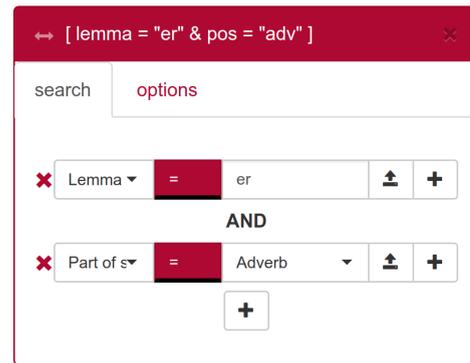
Adding attributes to a token box

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

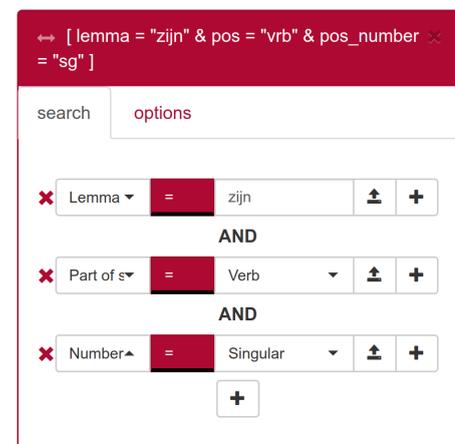
The *AND* option creates a new attribute restriction that a token must match in addition to the ones which were already there. As an example: suppose we want to match *zijn* ('to be') as a verb, not as a pronoun. First, fill in the attribute Lemma with value *zijn*, then click +, choose *AND*, and choose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

The screenshot shows a query builder window with a red header bar containing the query: `[lemma = "zijn" & pos = "vrb"]`. Below the header, there are two tabs: 'search' and 'options'. The 'options' tab is active. The main area contains two conditions stacked vertically, connected by an 'AND' operator. The first condition is 'Lemma = zijn', and the second is 'Part of s = Verb'. Each condition has a red 'X' icon on the left and a '+' icon on the right. A '+' button is also located below the second condition.

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*.



With the use of the AND-button it is possible to combine different attributes to make specific searches possible, as is shown in the token box on the right. By choosing Lemma (*zijn*), Part of speech (*Verb*) and Number (*Singular*), this will result in hits like *eist, esse, est, hets, waert en wart*.



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=NOU 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=*lief* OR Lemma=*goed*), thus resulting in the correct query [pos = "NOU-C" & (lemma = "lief" | lemma = "goed")], as shown in the right screenshot below.

← [(pos = "aa" & pos_number = "pl") | pos_number = "sg"]

search options

Part of s = Adjective-Adve

AND

Number = Plural

OR

Number = Singular

← [pos = "aa" & (pos_number = "sg" | pos_number = "pl")]

search options

Part of s = Adjective-Adverb

AND

Number = Singular

OR

Number = Plural

Before hit	Hit	After hit	Lemma	Part of speech + features
Corp. I, 0065A, Holland, grafelijke kanselarij, 24 oktober 1268	graue	van holland...	GRAAF	NOU-C(number=sg,inflectic
...zelant maken	cont	alle den...	KOND	AA(degree=und,position=p
...die dese	letten	sullen sien...	LETTER	NOU-C(number=sg,inflectic
...Florense onsen	neue	onsen moeyen...	NEEF	NOU-C(number=sg,inflectic
...neue onser	moyen	[sone] veren...	MOEJ	NOU-C(number=sg,inflectic
...onsen moeyen	[sone]	veren aleyden...	ZOON	NOU-C(number=sg,inflectic
...moyen [sone]	veren	aleyden van...	VROUW	NOU-C(number=sg,inflectic
...van henengowe	ser	... Johans wif...	DE+HEER	PD(type=dem,subtype=art,)
...ser ... Johans	wif	van Auennis...	WIJF	NOU-C(number=sg,inflectic
...was ... onse	woninghe	[dij] hare...	WONING	NOU-C(number=sg,inflectic
...al dat	lant	dat wil...	LAND	NOU-C(number=sg,inflectic
...tussen der	ouder	scie, en[de]...	OUD	AA(degree=und,position=p
...nuwen sciedamme	hofstaden	... ende al...	HOFSTEDE	NOU-C(number=sg,inflectic
...al dat	gherechte	ende visserie...	GERECHT	NOU-C(number=sg,inflectic
...gherechte ende	visserie	ende sluse...	VISSERIJ	NOU-C(number=sg,inflectic
...visserie ende	sluse	... die ons...	SLUIS	NOU-C(number=sg,inflectic
...hout te	lene	ende in...	LEEN	NOU-C(number=sg,inflectic
...in derre	maniren	dat gheulle...	MANIER	NOU-C(number=sg,inflectic
...stone sonder	kint	dat hel...	KIND	NOU-C(number=sg,inflectic
...xp sinen	ousten	broeder naest...	OUD	AA(degree=und,position=p

Before hit	Hit	After hit	Lemma	Part of speech + features
Corp. I, 0065A, Holland, grafelijke kanselarij, 24 oktober 1268	cont	alle den...	KOND	AA(degree=und,position=ps
...zelant maken	ouder	scie, en[de]...	OUD	AA(degree=und,position=ps
...tussen der	ousten	broeder naest...	OUD	AA(degree=und,position=ps
...xp sinen	ousten	of hi...	OUD	AA(degree=und,position=ps
...johanne den	vorghesegt	vrove aleyd...	VOORGEZEGD	AA(degree=und,position=ps
...sal dese	vorgheseghede	goed met...	VOORGEZEGD	AA(degree=und,position=ps
...al di	gheshouden	ende ghestade...	GEHOUDEN	AA(degree=und,position=ps
...dinc si	gheshouden	so hebbe...	GESTADE	AA(degree=und,position=ps
...ghehouden ende	ghestade	lukes dach...	GESTADE	AA(degree=und,position=ps
...sonsdags na	sente		SINTSINT-	AA(degree=und,position=ps
Corp. I, 0065AA, Holland, grafelijke kanselarij, 24 oktober 1268	cont	alle den...	KOND	AA(degree=und,position=ps
...zelant maken	ouder	scie, en[de]...	OUD	AA(degree=und,position=ps
...tussen der	ousten	broeder naest...	OUD	AA(degree=und,position=ps
...xp sinen	ousten	of hi...	OUD	AA(degree=und,position=ps
...johanne den	vorghesegt	vrove aleyd...	VOORGEZEGD	AA(degree=und,position=ps
...sal dese	vorgheseghede	goed met...	VOORGEZEGD	AA(degree=und,position=ps
...al di	gheshouden	ende ghestade...	GEHOUDEN	AA(degree=und,position=ps
...dinc si	gheshouden	so hebbe...	GESTADE	AA(degree=und,position=ps
...ghehouden ende	ghestade	lukes dach...	GESTADE	AA(degree=und,position=ps
...sonsdags na	sente		SINTSINT-	AA(degree=und,position=ps
Corp. I, 0102, Brugge, 26 februari 1271?	grote	Clais f...	GROOT	AA(degree=und,position=ps

The tab options

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

← [lemma = "optatief"]

search options

Optional

Begin of sentence

End of sentence

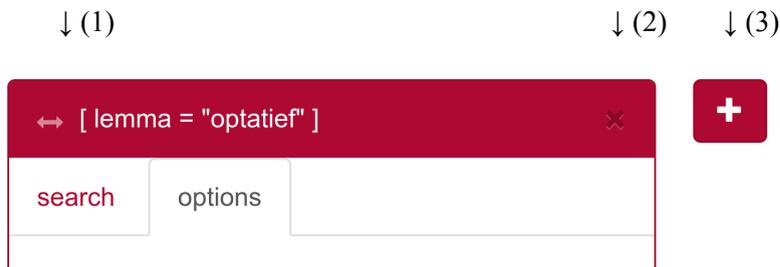
repeats 1 to 1 times

Managing sequences of token boxes

There are three ways to manage the sequence and the number of token boxes:

- *Rearrange* a token by clicking and dragging the little arrow handle in the top-left corner simultaneously (1).
- *Delete* a token by clicking the x in the top-right corner (2).

- *Create a new token box* by clicking the + -button next to the upper right corner of the utmost right token box (3).



Uploading value lists in the query builder

It is also possible to upload a list of values, separated by a white space. To do so, click the upload button (with the arrow pointing upwards) and select a text file. Tokens will then be matched for any of the values from the file.

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 or you can copy and paste the values into a *.txt file first.)

After uploading a file, the text can be edited by clicking the yellow marked file name in the text field. Editing the text is temporary and will not modify your original file.

To remove an uploaded file and go back to typing a value, click on the cross (x) next to the yellow text box. Another possibility to clear the uploaded values is by clicking the yellow marked text field and then press the Clear button on the bottom left corner of the Edit box. Using the Reset button will start a complete new search.

Copy to CQL editor

It is possible to copy a query - like [\[pos="aa"\]\[lemma="goed"\]](#) - to the CQL editor using the *Copy to CQL editor* button. This will take you automatically to the Expert Search screen, after which you can start the search or adjust the query if desired.

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"]` 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.
- The 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 (in fact, some of the links on the examples will direct you to there), 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 CQL manual in the appendix, 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.

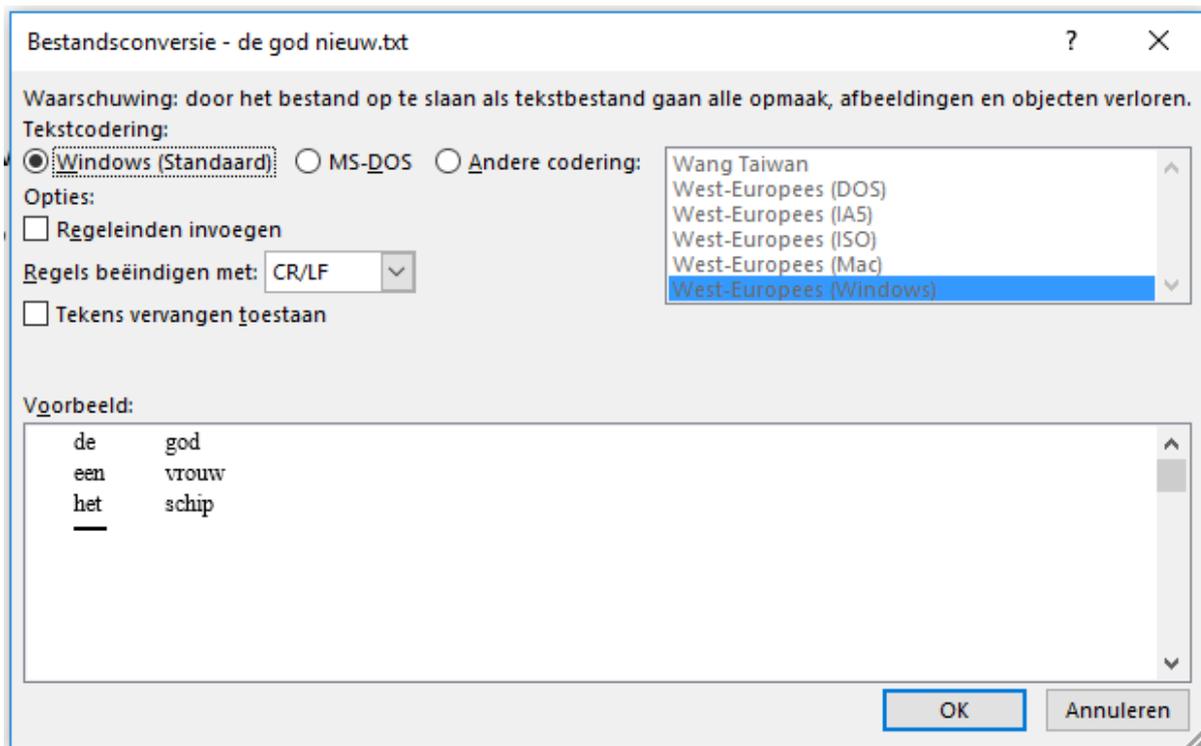
Import query

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 afterwards 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, as is shown in the following screenshot:



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"][lemma="@@"]
```

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 schip
```

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, schip* - at the position of the second gap. With these values, gap-filling yields the following results:

Before hit ▾	Hit ▾	After hit ▾	Lemma	Part of speech + features
Corp.I, 0171, Gent, 3e kwart 13e eeuw ...Tote	ere goeder vrouwen	ende ere...	EEN GOED VROUW	PD(type=undef,subtype=art,position=prenom,inflection=r/re) AA(degree=uncl,position=prenom px C(number=sg,inflection=n)
Corp.I, 1497, Mechelen, 10 mei 1296 ...goet , alse	ene Edele vrowe	, vrowe ver...	EEN EDEL VROUW	PD(type=undef,subtype=art,position=prenom,inflection=e) AA(degree=uncl,position=prenom pos C(number=sg,inflection=e)
Corp.I, 1149, Holland, grafelijke kanselarij, (21 april 1291-4 april 1292) ...hadden , met	eenre edelre vrouwen	ende groeter...	EEN EDEL VROUW	PD(type=undef,subtype=art,position=prenom,inflection=r/re) AA(degree=uncl,position=prenom px C(number=sg,inflection=n)
...worden , tieghens	eenre edelre vrouwen	ende groter...	EEN EDEL VROUW	PD(type=undef,subtype=art,position=prenom,inflection=r/re) AA(degree=uncl,position=prenom px C(number=sg,inflection=n)
Corp.I, 1315, Maastricht, 7 maart 1294 ...der versteruenessen	eynre eersamer vrouwen	wilen was...	EEN EERZAAM VROUW	PD(type=undef,subtype=art,position=prenom,inflection=r/re) AA(degree=uncl,position=prenom px C(number=sg,inflection=n)
Corp.I, 1323, Maastricht, 4 april 1294 ...der versteruenessen	eynre eersamer vrouwen	ver zyben...	EEN EERZAAM VROUW	PD(type=undef,subtype=art,position=prenom,inflection=r/re) AA(degree=uncl,position=prenom px C(number=sg,inflection=n)
Corp.I, 1324, Maastricht, 6 april 1294 ...der versteruenessen	eynre eersamer vrouwen	wilen was...	EEN EERZAAM VROUW	PD(type=undef,subtype=art,position=prenom,inflection=r/re) AA(degree=uncl,position=prenom px C(number=sg,inflection=n)
Corp.I, 1333, Maastricht, 9 mei 1294 ...ons quamen .	eyne eersame vrowe	. ver zybe...	EEN EERZAAM VROUW	PD(type=undef,subtype=art,position=prenom,inflection=e) AA(degree=uncl,position=prenom pos C(number=sg,inflection=e)
Corp.I, 1562, Mechelen, 19 januari 1297 ...goed alse	ene edele vrowe	, vrowe ver...	EEN EDEL VROUW	PD(type=undef,subtype=art,position=prenom,inflection=e) AA(degree=uncl,position=prenom pos C(number=sg,inflection=e)

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 - the bold words in the column Hit - to display the properties and values of the hit. Click the hit again to close.

Corp.I, 1149, Holland, grafelijke kanselarij, (21 april 1291-4 april 1292)			
...hadden , met	eenre edelre vrouwen	ende groeter...	EEN EDEL VROUW PD(type=undef,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
<p>Wi Jan van renisse ende didderic here van brederode , maken cont alle den ghenen die desen brief zullen sien isf horen lesen , dat wi enen tuist hadden , met eenre edelre vrouwen ende groeter , ver katerinen vrouwe van vome , ende burgrauinne van zelant , dien wi , bi onsen consentie , ende onsen vrien wille , ende bi rade onser vrien , ende onser maghen , verreffent hebben met haer in deser manire , alse , dat wi haer helpen zullen ghetrouwelike , ende raden , ende wisen , bi gheboefder trouwen , tieghens</p>			
Property	Value		
Word	eenre	edelre	vrouwen
Lemma	EEN	EDEL	VROUW
Part of speech + features	PD(type=undef,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)
Numerical part of speech code	485	105	004
Word id	w.616936	w.616937	w.616938

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.

Sorting results

Click on any of the column headings to sort the hits on values within the column, clicking again inverts the sorting. Extra sorting options are given when clicking on Before hit, Hit and After hit: you can sort by various attributes, as shown below.

Hits

Group hits by...

« 1 2 3 5 »

Before hit	Hit	After hit	Lemma	Part of speech + features
Corp.I, 0171, Gent, 3e kwart 13e eeuw ...Tote	ere goeder vrouwen		VROUW	PD(type=indef,subtype=art,positior C(number=sg,inflection=n)
Corp.I, 1497, Mechelen, 10 mei 1296 ...goet , also	ene Edele vrouwe		VROUW	PD(type=indef,subtype=art,positior C(number=sg,inflection=e)
Corp.I, 1149, Holland, grafelijke kanselarij, (21 april 1291-4 april 1292) ...hadden , met	eenre edelre vrouwen	ende groeter...	EEN EDEL VROUW	PD(type=indef,subtype=art,positior C(number=sg,inflection=n)
...worden , tieghens	eenre edelre vrouwen	ende groter...	EEN EDEL VROUW	PD(type=indef,subtype=art,positior C(number=sg,inflection=n)
Corp.I, 1315, Maastricht, 7 maart 1294 ...der versteruenessen	eynre eersamer vrouwen	wilen was...	EEN EERZAAM VROUW	PD(type=indef,subtype=art,positior C(number=sg,inflection=n)
Corp.I, 1323, Maastricht, 4 april 1294 ...der versteruenessen	eynre eersamer vrouwen	ver zyben...	EEN EERZAAM VROUW	PD(type=indef,subtype=art,positior C(number=sg,inflection=n)
Corp.I, 1324, Maastricht, 6 april 1294 ...der versteruenessen	eynre eersamer vrouwen	wilen was...	EEN EERZAAM VROUW	PD(type=indef,subtype=art,positior C(number=sg,inflection=n)

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, both by Hit, as Before hit, as After hit.

Sort by... Hide Titles Export CSV

Filter...

Hit

- Sort by Word
- Sort by Word (descending)
- Sort by Lemma
- Sort by Lemma (descending)
- Sort by Part of speech
- Sort by Part of speech (descending)
- Sort by Part of speech + features
- Sort by Part of speech + features (descending)

Before hit

- Sort by Word before

Grouping results

Results Per Hit can be grouped by properties of Hit, Before hit, After hit and by the metadata of the documents in which those hits occur (Date, Localization, Text type, Title and author). Grouping is

facilitated by the drop-down menu Group hits by.... By selecting one of the properties a tick box appears that makes it possible to distinguish between case-sensitive and case-insensitive.

Per Hit **Per Document**

Hits / Grouped by hit:word Total hits: 102 (0.00659%)
Total groups: 2

Group by Word **Case-sensitive**

« 1 »

table hits

Group	#hits in group	Relative frequency (hits)
ghemaket	101	0.00652%
Ghemaket	1	0.0000646%

Sort by... Export CSV

In the Per hit view, advanced grouping options are available by selecting the option Context (advanced). This option allows you to group the results by up to 5 tokens before or after the hit. It also allows you to group the results based on (parts of) the hits. By pressing the New context group you can group the results by another property or another range.

We will work that out using an example. A search for groups of three verbs - in Expert Search: [pos="VRB"]{3} - produces hits like the following (Titles are hidden):

...den ghoenen . die dese lettren	sullen zien & horen	lesen . dat quamen voer hons...	ZULLEN ZIEN & HOREN	VRB(type=coplaux, finiteness=fin, tense=pres, mood=ind conj, inflection=n) VRB(type=mai, finiteness=inf, inflection=0) & VRB(type=mai, finiteness=inf, inflection=0)
...ghoenen . die dese lettren sullen	zien & horen lesen	... dat quamen voer hons als...	ZIEN & HOREN LEZEN	VRB(type=mai, finiteness=inf, inflection=0) & VRB(type=mai, finiteness=inf, inflection=0) VRB(type=mai, finiteness=inf, inflection=0)
...up sinen wettachtighen oer . die hosen	es tewetene si	dat sake dat die uorseide...	ZUN TE+WETEN ZUN	VRB(type=mai, finiteness=fin, tense=pres, mood=ind conj, inflection=0) ADP(type=uncl, inflection=0)+VRB(type=mai, finiteness=inf, inflection=e) VRB(type=coplaux, finiteness=fin, tense=pres, mood=ind conj, inflection=0)

It is now possible to group the hits by the second and third tokens of those hits. See below.

Per Hit **Per Document**

Hits / Grouped by context:word:i:H2-3 Total hits: 3,725 (0.241%)
Total groups: 2,659

Context (advanced)

Word Case-sensitive

From end of hit

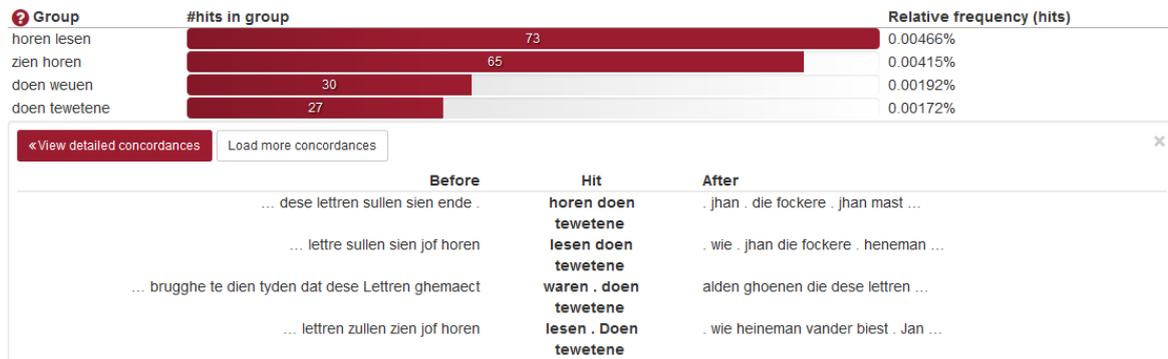
New context group

« 1 2 3 4 6 11 »

table hits

Group	#hits in group	Relative frequency (hits)
horen lesen	74	0.00478%
zien horen	63	0.00407%
doen weuen	29	0.00187%
Doen tewetene	27	0.00174%
sal cesseren	24	0.00155%
soude cesseren	19	0.00123%

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.

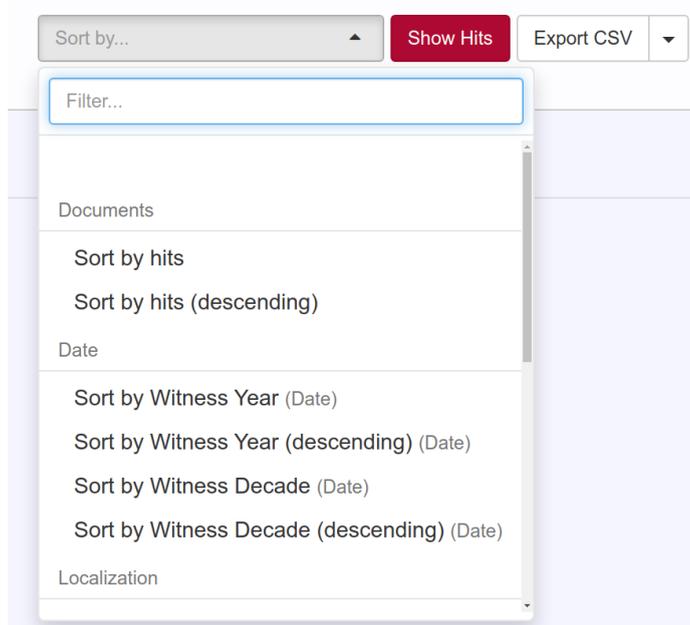


Click on View detailed concordances to go back to the normal hits view to see more detailed information for the hits in this group. The button Go back to grouped view brings you back to the list of groups.

Per Document view

Sorting results

Results can be sorted by means of the drop-down menu at the bottom of the page, which enables you to sort on Documents and on Date, Localization, Text type, and Title and author.



Grouping results

Results Per Document can be grouped by the metadata of the documents in which those hits occur (Date, Localization, Text type, and Title and author). Here, grouping is facilitated by the drop-down menu Group docs by....

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 (see screenshot below). 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.

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

Per Hit **Per Document**

Hits / Grouped by hit:word Total hits: 3.457 (0.223%)
Total groups: 37

Group by Word Case-sensitive

« 1 2 »

table hits

Group	#hits in group	Relative frequency (hits)
gode	846	0.0547%
god	828	0.0535%
got	621	0.0401%
gods	571	0.0369%
godes	332	0.0214%

«View detailed concordances» Load more concordances

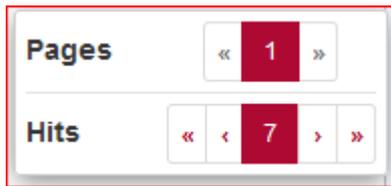
Before	Hit	After
ghemeender aelmoesene achter tof . Jn	godes	name si hem allen cont
soe gheven wi hem . in	godes	Ere . Ende omme der port
heren flormans sone visirde . JN	Godes	name amen , Want die mensce
sente Johan bi der ghenade	godes	, hebben onsen seghel ghedaen , an
Jn	Godes	name . Amen . Went die mensce
unsen here Jan uan der	godes	ge naden Hertoge uan Lothrigge ende

Information about a document

Click on a document title to open the document in a new 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. You can navigate from one hit to another by using the arrows at the Pages and the Hits button:



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:

40 Tote ere goeder vrouwen ende ere vroeder mire vrouwen mijns ser

Word: spetale

Word id: w.53565

Lemma: SPITAAL

Part of speech: NOU-C(number=sg,inflection=e)

Matching lemmata in the Dictionary of Early Middle Dutch:

- spitael (znw.o.), *gasthuis*

ente pieters onbieden
e dat wi een spetael
elc spetael
n coep af hieuet ende al
wilken dat wi v bidden

1 dat ghi onsen **spetale** sine hervachtechede behout bedi ghi sijt sculdech
te doene. ende dat ghire so vele toe doet dat ons ghene noet ne es dat
wire voerder toe doen bedi wine ghedoghen niet. ons danckes dat
tspetael verliest. ende weet wel dat. al ware dat sake dat v scepenen den

5 **spetale** ontwijst hadden nochtan ware te verstante hoe uerre het taerliede
vonnese boert bedie dat kerkelijc goet es ende ter almoesenen boert.
ende ghi ende v schepenen wilen eene ebisie van doernike van den
seluen sticken gheadt hebt daer bi bidden wi v dat ghire so vele toe doet
dat wijs v danc weten god si met v

Metadata of a document

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

Statistics

The Statistics tab shows several document statistics: the number of Tokens, the number of Types (unique word forms), the number of Lemmas 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 or via the menu symbol right of 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 Text type: Subgenre *medicine*
- Press Search

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

- Navigation: Search, Explore
- Subtabs: Corpora, N-grams, Statistics
- Group documents by metadata: Group by Subgenre (Text type)
- Show groups as: docs
- Filter search by: Date, Localization, Text type (1), Title and author
- Fictionality: Fictionality
- Genre: Genre
- Subgenre: medicine

After pressing the bar with the number of docs in group, you will get this result:

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

- Navigation: Per Hit, Per Document
- Path: Documents / Grouped by field:subgenre
- Summary: Total documents: 2 (100%), Total groups: 1
- Grouping: Group by Subgenre (Text type), Case sensitive
- Page navigation: < 1 >
- View options: table, docs, tokens
- Table header: Group, #docs in group, Relative frequency (docs)
- Table content: medicine, 2, 100%
- Action: < View detailed concordances
- Document text: Nederbergse geneeskundige geneeskundige recepten, Noordlimburgse gezondheidsregels

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

N-grams

An *N-gram* is a sequence of *N* items: Word, Lemma and Part of speech (+features). 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*: choose for sequences of Word (i.e. word form), Lemma, Part of speech or Part of speech + features. If you do not specify the search term further, a series of five consecutive Words, Lemmas, Parts of speech or Parts of speech + features will be searched for.
- It is also possible to restrict to, for instance, 5-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

Search Explore

Explore ...

Documents N-grams Statistics

N-gram size: 5

N-gram type: Word

Word Part of speech Part of speech Word Word

ik Verb Verb Word Word

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

Per Hit Per Document

Hits / Grouped by hit:word

Total hits: 904 (0.0584%)
Total groups: 869

Context (advanced) Apply

New context group

« 1 2 3 4 6 11 »

table hits

Group	#hits in group	Relative frequency (hits)
ics vermaent worde van hem	6	0.000388%
ic vermach dade ic dat	4	0.000258%
ic belouet hebbe ende beloue	4	0.000258%
ic hebbe ontfanghen jn lene	3	0.000194%
ic ghelouet hebbe vor mi	3	0.000194%
ic gheset hebbe ende ghegheuen	3	0.000194%
ic salne don op herstaen in	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 or Part of speech + features
- 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 the manuscripts of *Der Naturen Bloeme* in the *Corpus Gysseling* by searching for Frequency list type Word and by filtering search by Title and author, Words in title *naturen*:

The screenshot shows a search interface with a dark red header containing 'Search' and 'Explore' tabs. Below the header, there are three tabs: 'Corpora', 'N-grams', and 'Statistics'. The 'Frequency list type' dropdown is set to 'Word'. Under 'Filter search by ...', there are four tabs: 'Date', 'Localization', 'Text type', and 'Title and author' (which is selected and has a '1' icon). Below these tabs, there are three input fields: 'Words in title' containing 'naturen', 'Title of containing publication' (a dropdown menu), and 'Author' (a dropdown menu).

This results in:

The screenshot shows the search results interface. At the top, there are two tabs: 'Per Hit' and 'Per Document'. Below this, there is a header 'Hits / Grouped by hit:word' and 'Total hits: 91,869 (100%)' and 'Total groups: 12,421'. There is a 'Context (advanced)' dropdown and an 'Apply' button. Below that, there is a 'New context group' input field. A pagination bar shows '« 1 2 3 4 6 11 »' with a pencil icon next to '1'. There are two tabs: 'table' and 'hits'. The 'hits' tab is selected, showing a table of results.

Group	#hits in group	Relative frequency (hits)
ende	3,867	4.21%
die	2,689	2.93%
dat	2,166	2.36%
es	1,693	1.84%
in	1,409	1.53%
van	1,338	1.46%

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 in 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 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-insensitively, 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 a 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!" (A note about 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.)

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 that forward slash at the end of the tag. This way of referring to the element means "the whole element". Compare 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.

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.