

Corpus Oudnederlands application manual

Table of contents

Introduction	4
Information about the corpus	4
Lemmatization	4
Part of speech tagging	4
Metadata	5
Date	5
Witness Year	5
Permissive / Strict	5
Localization	5
Country	5
Area	5
Place	6
Kloeke location code	6
Text type	6
Fictionality	6
Genre	6
Subgenre	6
Title and author	6
Title	6
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	11
Wildcards	13
Upload a list of values	13
Part of speech dialog box	14
Cliticity	14
Complete word or word part	14
Starting a new search	15
Filter search by	15
Advanced search	16
The query builder	16
The tab search	16
Token attributes	17
Adding attributes to a token box	17
Function of the two +-buttons in a token box	18
The tab options	19
Managing sequences of token boxes	19
Uploading value lists in the query builder	20
Copy to CQL editor	20
Expert search	20
Copy to query builder	21
Import query	21
Gap filling	21
Viewing results	23
Per Hit view	23
Sorting results	23
Grouping results	24
Per Document view	26
Sorting results	26
Grouping results	27
Exporting results	27
Information about a document	27
Content	27
Metadata of a document	28
Statistics	28

Exploring the corpus	28
Documents	28
N-grams	29
Options	29
Example	29
Statistics (frequency lists)	30
Options	30
Example	30
Appendix: Corpus Query Language	32
CQL support	32
Supported features	32
Differences from CWB	33
(Currently) unsupported features	34
Using Corpus Query Language	34
Matching tokens	34
Sequences	35
Regular expression operators on tokens	35
Case- and diacritics-sensitivity	36
Matching XML elements	36
Labeling tokens, capturing groups	37
Global constraints	37

Introduction

This manual describes the corpus exploitation environment for the *Corpus Oudnederlands*. 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 Oudnederlands* in the current release is a collection of all remaining Dutch word material from the period 475-1200 that served as source material for the *Oudnederlands Woordenboek* (ONW; *Dictionary of Old Dutch*). This collection of material consists of various components: three large texts (Wachtendonkse Psalmen, Leidse Willeram, Mittelfränkische Reimbibel) and numerous smaller Old Dutch texts and text fragments (including rune inscriptions), glosses and single words, Frankish material and toponymic material. More detailed information on how the corpus was compiled can be found [here](#).

A first online accessible version of the corpus was launched on 29 February 2012, in the form of a quotation database, in which it was not possible to search for consecutive words. The reason was that each word in a quotation had a record in the database, with linguistic annotation (part of speech and lemma), the full quotation and the metadata of the quotation. This version is no longer available.

In 2018, for the Nederlab project, the data from the relational database were converted into token by token linguistically annotated text, including corrections of the linguistic annotation and additional text metadata.

In this new version, several corrections have been made to the added metadata in the corpus and the linguistic annotation has been mapped to the TDN-tagset (see below).

Lemmatization

The Old 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 het GiGaNt-lexicon](#).

Part of speech tagging

The original part of speech tagging of the Corpus Oudnederlands was done according to the guidelines developed for the Dictionary of Old Dutch (ONW). The Corpus was tagged manually by the editors of the ONW.

In the context of the CLARIAH+ project, a tagset and tagging principles for the annotation of diachronic corpora of historical Dutch have been developed: *Tagset voor Diachroon corpusmateriaal van het Nederlands (TDN)*. A detailed description can be found [here](#). The original part of speech tagging has been converted into the TDN, and is used in the current application.

Metadata

The *Corpus Oudnederlands* has also 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

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, *Le Compte Général de 1187, connu sous le nom de "Gros Brief"* can be dated exactly to 1187, while the *Liber Traditionum Sancti Petri Blandiniensis* originated between 639-1200.

Permissive / Strict

It is possible to do a permissive and strict search for Witness 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 800 and 825.

If you choose to do a Strict search by Witness Year, the search query will only result in manuscripts that were produced later than 800 but before 825. The *Corpus Oudnederlands* has two such documents: *Nederbergse doopbelofte*, which was handed down in a manuscript dating from the period 811-812 and *Runeninscriptie Bernsterburen*, dating from 800.

If, on the other hand, you choose the option Permissive in Witness Year, more documents are found, one of which is *Traditiones et antiquitates Fuldenses*, which was handed down in a manuscript dating from the period 822-845. This dating does indeed partly fall within the indicated period 800-825.

Localization

Country

In this corpus, three countries are distinguished: België, Duitsland and Nederland.

Area

Within the above mentioned countries, the following areas are distinguished: Friesland, Groningen, Limburg, Nederberg, Noord-Holland, Noord-Nederrijn, Oostnederrijn-Westfalen, Utrecht, and Zuid-Nederrijn.

Place

If it was possible to determine exactly where a particular text originated or was found, this Place (in modern spelling) will be mentioned: Arum, Britsum, Egmond, Essen, Garnwerd, Loppersum, Munsterbilzen, Nederberg, Oostum, Raskwerd, Toornwerd, Werden and Westeremden.

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 *Kloekencode*. More information about (searching with) the Kloekencode can be found [here](#).

Text type

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

Fictionality

All the documents in the *Corpus Oudnederlands* are considered to be non-fictional texts (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 different subgenre labels (see the paragraph Filter search by). These 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.

- *Biblical text*: Bible, lectionaria, diatessara, bible books (Psalm, Song of Songs)
- *Glossary*: collection of glosses
- *History*: text with historical information
- *History, theology*: text with both historical and theological information
- *Legislative/administrative document*: official records, charters
- *Legendary biography*: description of the life of worldly leaders
- *Legendary hagiography*: description of the life of saints
- *Miscellaneous*
- *Probatio pennae*
- *Religious*: texts dealing with religion or religious matters
- *Religious, secular*: texts dealing with both religious and non-religious matters
- *Runic Inscription*: inscription, dedication

Title and author

Title

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

All 90 documents of this corpus come from different sources (see also the About).

Author

It is possible to search by author name. However, for almost all documents in this corpus the author is unknown or uncertain. Only seven names of authors of Old Dutch texts in this corpus have been handed down: Alfrid, Einhard, Galbert de Bruges, Hariulf, Jean de Klerk, Rudolfus Trudonensis and Willeram van Ebersberg.

Application user manual

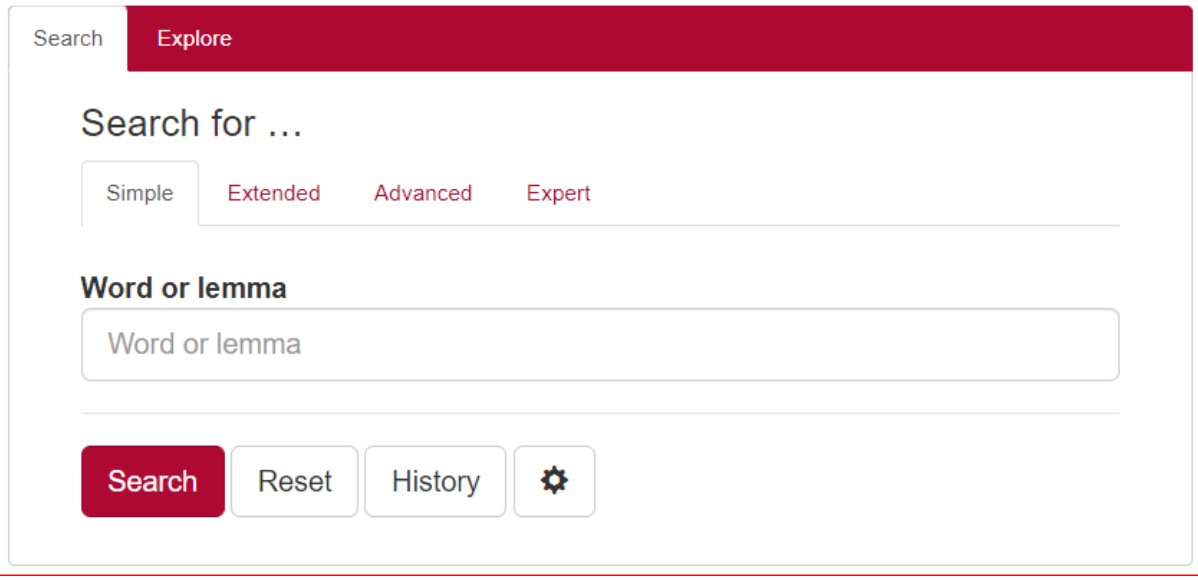
Getting started

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 [gisund](#)
 - Extended Search for Word [gisund](#)
- 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 [onrecht](#)
 - Extended Search for Lemma [onrecht](#)
- 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 *an* in [Simple Search](#)
 - only words starting with *ver* and ending with *an* in [Extended Search](#)
 - only lemmata starting with *ver* and ending with *en* in [Extended Search](#)
 - lemmata starting with *be*, ending in *en* with one syllable in between in [Expert Search](#)
- To search for a multi-word pattern, e.g. all determiners appearing before a given lemma as a noun, use the query builder in Advanced Search or use Expert Search:
 - adpositions before the lemma *huis* in [query builder](#) in Advanced Search
 - adposition 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 word man](#)
 - example Group by Lemma before: [words preceding the lemma god](#)
- To explore the distribution of document properties in the corpus, use the Explore feature
 - example: [characteristics about subgenres](#)
 - example: [localization](#)

Searching the corpus

Simple search



The screenshot shows a search interface with a dark red header containing 'Search' and 'Explore' tabs. Below the header, the text 'Search for ...' is displayed. There are four tabs: 'Simple' (selected), 'Extended', 'Advanced', and 'Expert'. A text input field labeled 'Word or lemma' contains the placeholder text 'Word or lemma'. At the bottom, there are four buttons: 'Search' (dark red), 'Reset', 'History', and a gear icon for settings.

Search

The Simple Search allows you to quickly search for specific words (e.g. *banuerc*) or lemmata (e.g. *banwerk*). It is also possible to enter a phrase: *ad banuerc* or *ad banuerc constituti sunt*. 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 e.g. *manne*, *manno*, *mannon*, *men*, *min*, 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: *banuerc* will deliver the same results as *BANUUERC* or *Banuerc*. 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 an *n*, e.g. *amman*, *annen* and *Athalbern* but also *ana* (lemma *aan*), *antsceine* (lemma *antschijn*) and

anther (lemma *aan+die*)

? The question mark matches a single character once. Therefore, *b?n* matches *only* three-letter values starting with a *b* and ending with an *n*, e.g. *bin*, *ban*, *ben* but also *bannum* (lemma *ban*), *buna* (lemma *bun*)

This wildcard can be used more than once. Thus *b???n* matches *began*, *baden*, *bedon*, *bacon*, but also *brun* (lemma *bruin*), *bethis* (lemma *beren*) and *geboren* (lemma *beren*).

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.]

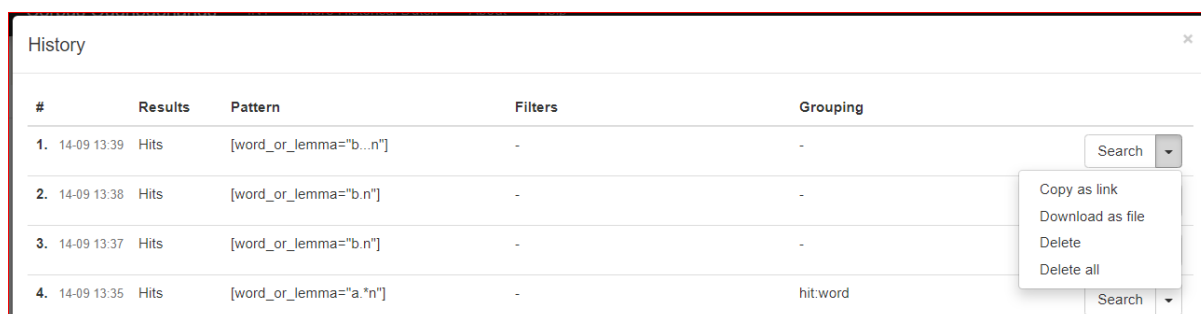
Reset

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 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.

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 five attributes you can search for are Word (more precise: word form), Lemma, Part of speech, Cliticity and Complete word or word part. All supported attributes are shown in the search form:

In the search fields Word and Lemma 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 seven 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* (5x), but also the variant *maria* (2x).

Per Hit | Per Document

Hits / Grouped by hit:word

Total hits: 7 (0.0114%)
Total groups: 2
Search time: 0.2s

Group by Word Case-sensitive

« 1 »

table hits

Group	#hits in group	Relative frequency (hits)
Maria	5	0.00811%
maria	2	0.00325%

Sort by... Export Export for Excel

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 | Explore

Search for ...

Simple Extended Advanced Expert

Word Case- and diacritics-sensitive

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 *berg* in Word will only provide you with occurrences of that exact string of characters. When you enter *berg* in the search field Lemma you will - besides the lemma *berg* - also find all word forms that are linked to that lemma, such as the spelling variant *bergh* and inflected forms as *bergon*, *bergo*, *berga*, *berge*, *bergan* and *bergas*.

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 an *n*, e.g. *amman*, *Albin*, *annen*, *annen*, *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 *b?n* in Lemma matches *only* three-letter lemmata starting with a *b* and ending with an *n*, i.e. *ban* (with word forms *bannum*, *ban*, **ban*) and *bun* (word form *buna*).

This wildcard can be used more than once. Thus *b???n* (in Word) matches *began*, *baden*, *bedon*, *bacon* and *boben*.

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. *then bergan*, *then ** or ** bergan*. Note that searching for *then bergan*, *then ** and ** bergan* in the search field Lemma will give no 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 Noun Common and Number Singular will result in a list of all occurrences of the singular noun *man*. The syntax of your query is shown in the results:
[word="man"&pos="nou\c"&pos_number="sg"]
- A multi-token example: searching for *kinde hiez* in the Word(form) field and *kind heten* in the Lemma field should find those occurrences of the bigram in which the first word is the spelling variant of the noun *kind* and the second the declined form of the verb *heten*.

Upload a list of values

At the right side of the search fields Word and Lemma, 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.

Part of speech	Type	Finiteness	Tense	Mood	Number	Person
Adjective-Adverb	<input type="checkbox"/> Copula	<input type="checkbox"/> Finite	<input checked="" type="checkbox"/> Present	<input type="checkbox"/> Indicative	<input type="checkbox"/> Singular	<input type="checkbox"/> 1
Adposition	<input type="checkbox"/> Auxiliary	<input type="checkbox"/> Infinitive	<input type="checkbox"/> Past	<input type="checkbox"/> Conjunctive	<input checked="" type="checkbox"/> Plural	<input type="checkbox"/> 1 or 3
Adverb	<input type="checkbox"/> Unclear	<input type="checkbox"/> Present participle	<input type="checkbox"/> Unclear	<input type="checkbox"/> Imperative	<input type="checkbox"/> Unclear	<input type="checkbox"/> 2
Conjunction		<input type="checkbox"/> Past participle		<input type="checkbox"/> Unclear		<input type="checkbox"/> 3
Interjection		<input type="checkbox"/> Unclear				<input type="checkbox"/> Unclear
Noun Common						
Noun Proper						
Numeral						
Pronoun-Determiner						
Residual						
Verb						

pos="vrb"&pos_tense="pres"&pos_number="pl"

Ok Reset

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 - present - plural will result in *hebban* and *[u]mbida[n]* and numerous other hits.

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 *námir* [NA+IK] and *woldik* [WILLEN+IK].

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 *afnemen* ('to take off'; conjugated form *nam ... áue*, 'took off') and pronominal adverbs as *thar umbe* (daarom). Keep in mind that you can only find both parts at the same time using Lemma (*afnemen*) and the option part. If you are specifically looking for just

one of the composing parts (e.g. *áue*), 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. *aua*, *abe*).

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 and Lemma 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 and Lemma, there are two possibilities to start a new search: fill in the desired value and 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. (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

Country
Country

Area
Limburg

Place
Place

Kloeke location code
Kloeke location code

Area (Localization): *Limburg*

Selected subcorpus:
Total documents: 1 (1.11%)
Total tokens: 6 (0.00974%)

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 Area. 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.)

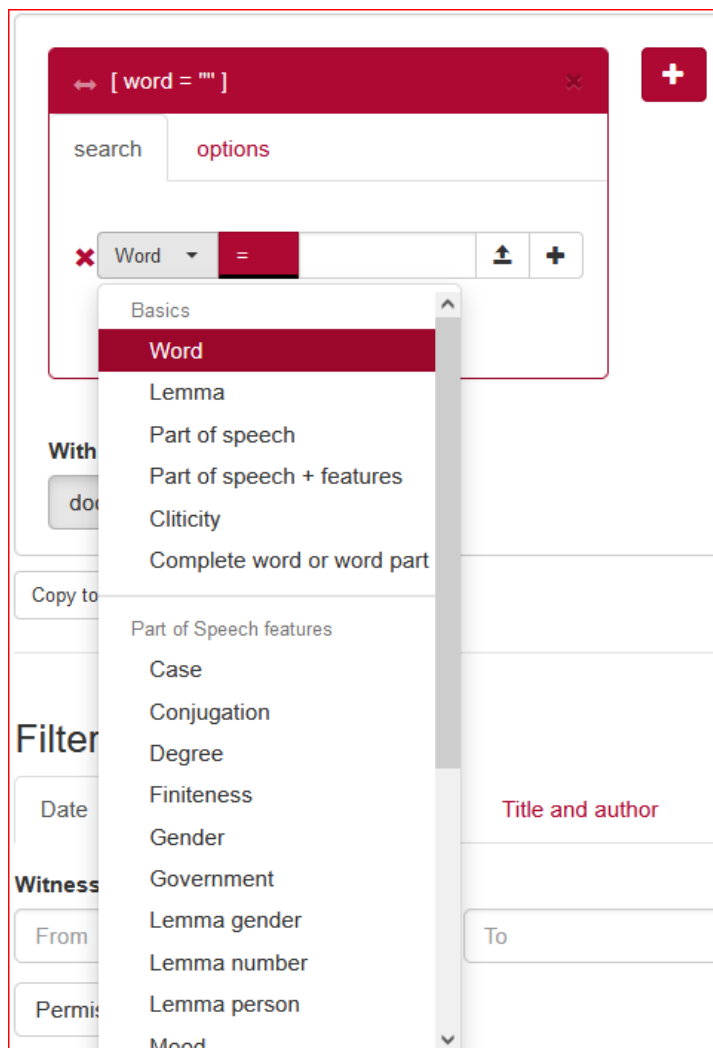
For a detailed description of the metadata, see the section Metadata categories at the beginning of this manual.

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.



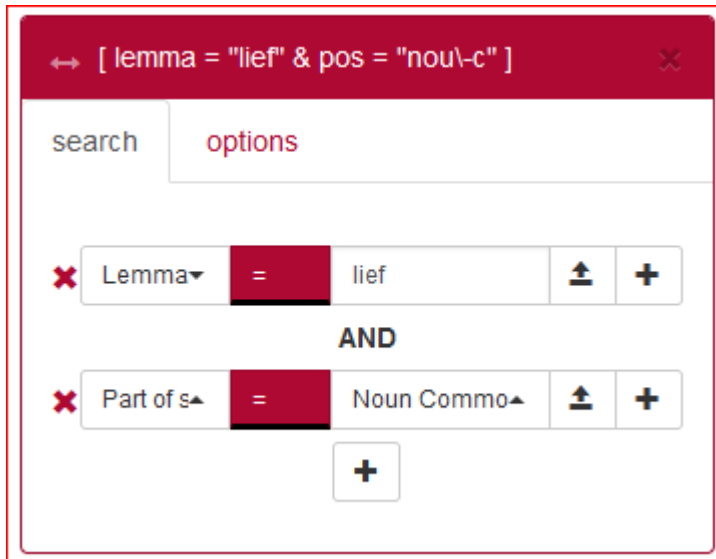
A token box in the querybuilder has two tabs: search and options.

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=Noun common 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 interface. At the top, a red header bar contains the query: [lemma = "lief" & pos = "nou\c"]. Below the header, there are two tabs: "search" and "options". The "options" tab is active. The interface consists of two rows of attribute-value pairs, each with a red 'X' on the left and a '+' button on the right. The first row is: Lemma = lief. The second row is: Part of s = Noun Commo. Between the two rows, the word "AND" is centered. Below the second row, there is a '+' button.

Token attributes

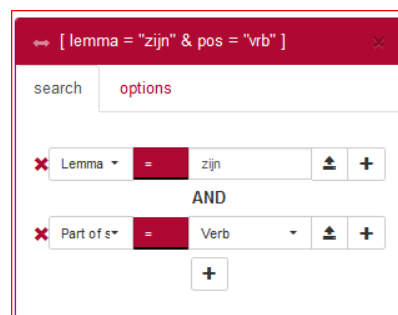
Specifying token attributes is similar to the Extended Search form. Select which attribute a token should have, and enter the value that the attribute must have for the token to be matched. Attributes in the query builder are interpreted as *regular expressions*. Note that this is different from the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

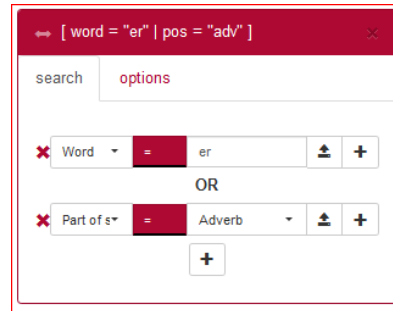
Using the +-button, new attributes can be added. Two options exist: *AND* and *OR*.

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 the verb *zijn* ('to be'). 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 interface. At the top, a red header bar contains the query: [lemma = "zijn" & pos = "vrb"]. Below the header, there are two tabs: "search" and "options". The "options" tab is active. The interface consists of two rows of attribute-value pairs, each with a red 'X' on the left and a '+' button on the right. The first row is: Lemma = zijn. The second row is: Part of s = Verb. Between the two rows, the word "AND" is centered. Below the second row, there is a '+' 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*.

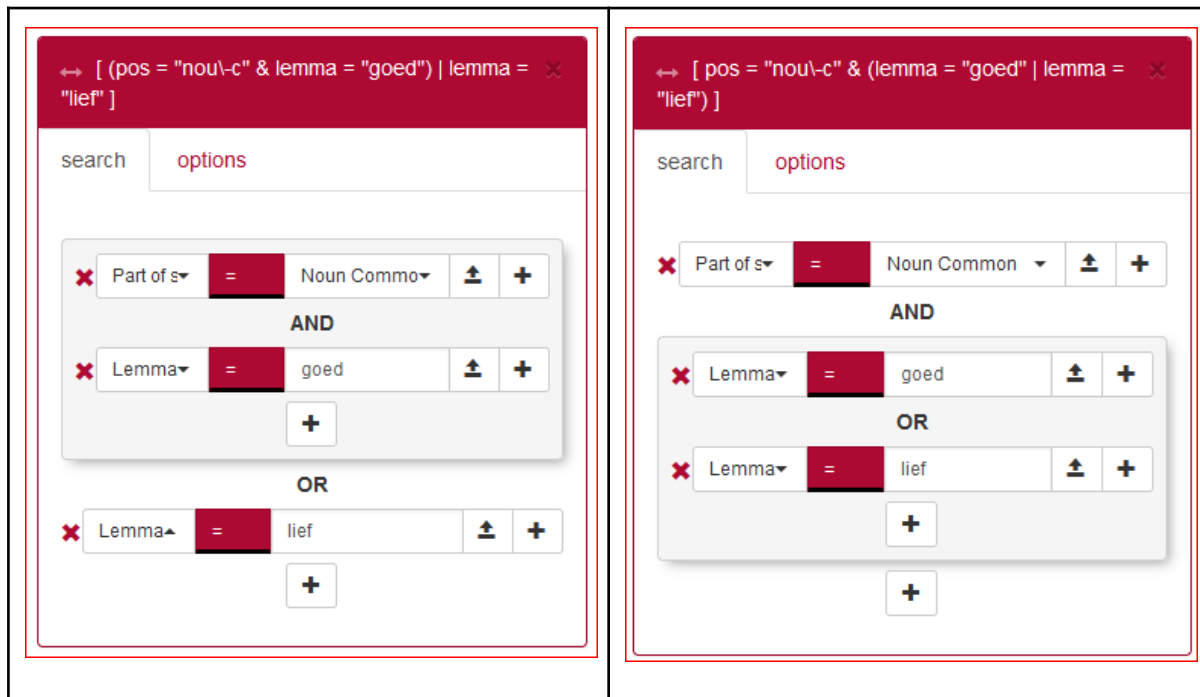


Function of the two +-buttons in a token box

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 “Thes scalt thu nu liebe uater ezzen turch thinen willen haben ich iz gewonnen”, where *liebe* 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 = "goed" | lemma = "goed", as shown in the right screenshot below.

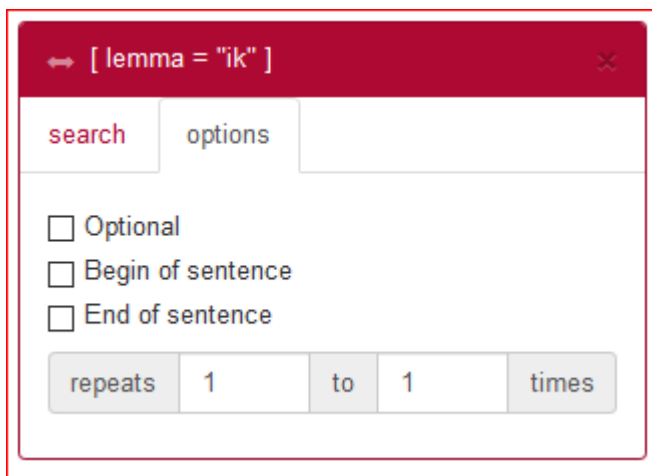


Before hit -	Hit -	After hit -	Lemma	Part of speech	features
Middelfrankische Reimbibel A ... si c'were ére else ... worden Getheerne nu quad her	liet lieue	so ther if The heide kint that thu gender unt ferges...	LIEF	AA(degree-pos, position-pred, case-unc1, number-unc1, gender-unc1)	LIEF
... gesot Thu haues veroren thaz	got	ande sal emer thion thise...	GOED	NOU-C(number-sg, case-acc, gender-n)	GOED
Middelfrankische Reimbibel B ... s[. h[] Thes scait thu nu	liebe	ualler ezzen lurch thinen euten...	LIEF	AA(degree-pos, position-precomp, position-pred, case-mon, number-sg, gender-m)	LIEF
... we gewinnest thu so siere ... indes viner ualer is so	liet liet	sun du Thei sun sprach... habet also wir haben gesaget	LIEF	AA(degree-pos, position-precomp, position-pred, case-mon, number-sg, gender-m)	LIEF
... froun stuzet unde thise werelelich	güt	stuzet Thu scait uon thinen...	GOED	NOU-C(number-pl, case-acc, gender-n)	GOED
(Expositio) Wilrammi Eberspergenis Abbas in Cantu's Cantuarum ... ih wole weg scait wole	liet	is ton haue thaz wole...	LIEF	AA(degree-pos, position-pred, case-unc1, number-unc1, gender-unc1)	LIEF
... Of thei menisco al sin	guod	hine geguet hz is himo angen...	GOED	NOU-C(number-sg, case-acc, gender-n)	GOED
... machot hin contemprom alles erthiscan	guodes	ande machot hin gregan thes...	GOED	NOU-C(number-sg, case-gen, gender-n)	GOED
... then thu mugest wanda mir	liet	ande lustich is thine stirma...	LIEF	AA(degree-pos, position-pred, case-unc1, number-unc1, gender-unc1)	LIEF

Before hit -	Hit -	After hit -	Lemma	Part of speech	features
Middelfrankische Reimbibel A ... gesot Thu haues veroren thaz	got	ande sal emer thion thise...	GOED	NOU-C(number-sg, case-acc, gender-n)	GOED
Middelfrankische Reimbibel B ... inden stuzet unde thise werelelich	güt	stuzet Thu scait uon thinen...	GOED	NOU-C(number-pl, case-acc, gender-n)	GOED
(Expositio) Wilrammi Eberspergenis Abbas in Cantu's Cantuarum ... Of thei menisco al sin	guod	hine geguet hz is himo angen...	GOED	NOU-C(number-sg, case-acc, gender-n)	GOED
... machot hin contemprom alles erthiscan	guodes	ande machot hin gregan thes...	GOED	NOU-C(number-sg, case-gen, gender-n)	GOED
Die altnieder- und altniederfränkischen Psalmen und Glossen in H ... ansene bogin that geresoda uerthin	lieua	thina behalden duo mit forboron...	LIEF	NOU-C(number-pl, case-mon, gender-unc)	LIEF
Die altnieder- und altniederfränkischen Psalmen und Glossen in H ... mit cofte mikito cung cofte	lieuis	lieuis in scuonis husis te...	LIEF	NOU-C(number-sg, case-gen, gender-unc)	LIEF
... crethe mikito cung crethe lieuis	lieuis	in scuonis husis te deline...	LIEF	NOU-C(number-sg, case-gen, gender-unc)	LIEF

The tab options

The tab options specifies the contextual properties, such as whether the token occurs at the end of a sentence, and the repetition pattern. However, since most Old Dutch texts lack punctuation, it is not always meaningful to use this function.



Managing sequences of token boxes

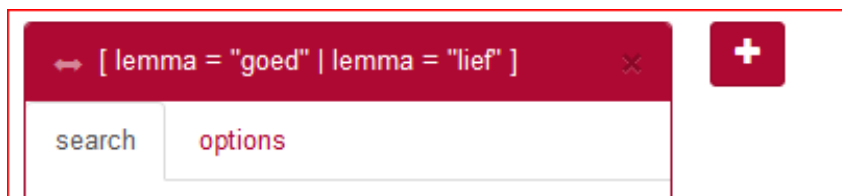
There are three ways to manage the sequence and the number of token boxes:

- *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).

↓ (1)

↓ (2)

↓ (3)



Uploading value lists in the query builder

It's 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 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 query builder* 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"]` (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 35 hits, but the second one in zero!

Some examples:

- Simple: `[word="man"]`, e.g. the attribute word matches the regular expression *man*; `[word!="man"]`, e.g. the attribute word does **not** match the regular expression *man*; `[lemma=". *man"]` matches all lemmata ending with *man*, including *man* itself.
- 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. [\[word or lemma="goed" & pos!="AA"\]](#) or - equivalently - [\[word or lemma="goed" & !pos="AA"\]](#) matches all occurrences of *goed*, not being an adjective.

- Repetition operators: [\[pos="NOU-C"\]{3}](#) matches a sequence of 3 common nouns, [\[pos="NOU-C"\]{2,4}](#) matches a sequence of 2 to 4 common nouns, [\[pos="NOU-C"\]{3,}](#) matches a sequence of 3 or more common nouns.
- 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: ["dijn" "heer" | "zijn" "zoon"](#), or even [\("dijn" "heer" | "zijn" "zoon"\)+](#), matching any sequence of *dijn heer* or *zijn zoon*. 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 [\[lemma=".*man"\]\[pos="nou-c"\]](#) - 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.

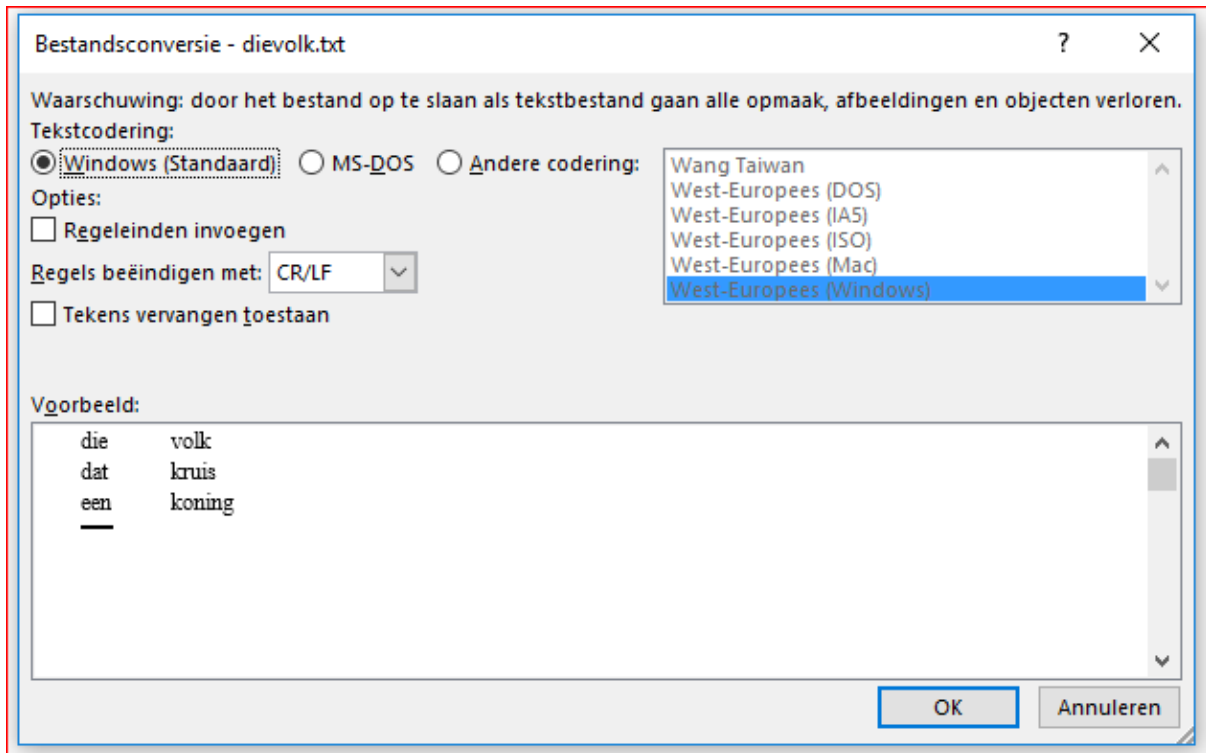
Import query

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, as is shown in the following screenshot:

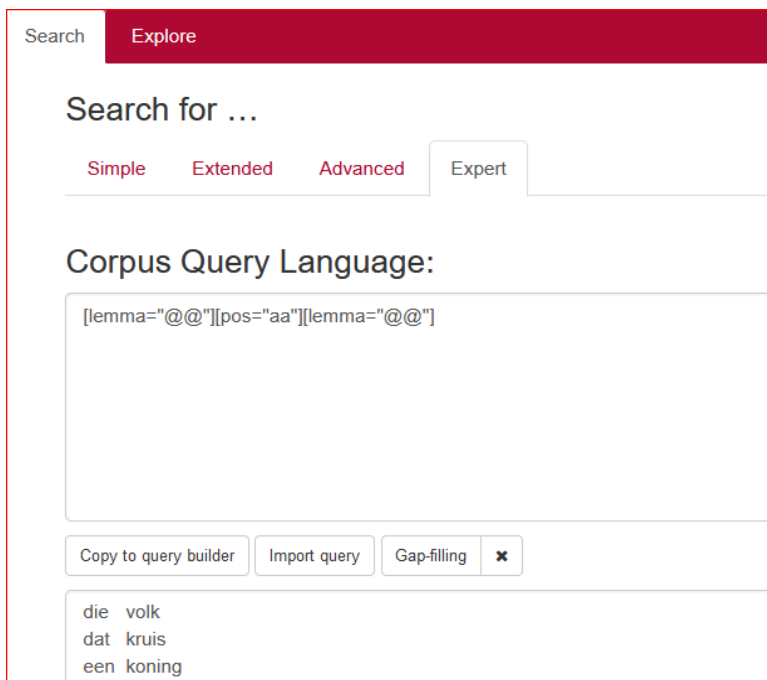


A *.tsv file or a comparable *.txt file enables you to complete a query with marked gaps.

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:



The values in the first column - *die, dat, een* - will be entered at the position of the first gap (@@) and the values in the second column - *volk, kruis, koning* - 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
Mittelfränkische Reimbibel A				
...herren	ein	the hiez Cosdras.	EEN GRIM	PD(type=indef,subtype=art,position=prenom,case=nom,number=sg,gender=m)
Inperside tho	grimme	Ther tho...	KONING	AA(degree=pos,position=prenom postnom pred,case=nom,number=sg,gender=m) NOU-
geweldig was	kunig			C(number=sg,case=nom,gender=m)
...cristenheit her	that	geweldeliche	DAT	PD(type=dem,subtype=art,position=prenom,case=acc,number=sg,gender=n)
ce stordo Ande	heilige	ande uorde iz	HEILIG	AA(degree=pos,position=prenom postnom pred,case=acc,number=sg,gender=n) NOU-
nam	cruce	insin...	KRUIS	C(number=sg,case=acc,gender=n)
Mittelfränkische Reimbibel B				
...zû ierusalem	thaz	quam inthaz lant	DAT	PD(type=dem,subtype=art,position=prenom,case=acc,number=sg,gender=n)
uant thie turch	heilige	Then sie...	HEILIG	AA(degree=pos,position=prenom postnom pred,case=acc,number=sg,gender=n) NOU-
	cruce		KRUIS	C(number=sg,case=acc,gender=n)
...sie ire sageten	thaz	uerborgen	DAT	PD(type=dem,subtype=art,position=prenom,case=acc,number=sg,gender=n)
war sie	heilige	habeten Sie	HEILIG	AA(degree=pos,position=prenom postnom pred,case=acc,number=sg,gender=n) NOU-
	cruce	sageten thaz...	KRUIS	C(number=sg,case=acc,gender=n)

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

Before hit -	Hit -	After hit -	Lemma	Part of speech + features
Mittelfränkische Reimbibel A				
...jüthen uiengen ande an ein	cruce	gehiengen Wie solden the juthen...	KRUIS	NOU-C(number=sg,case=acc,gender=n)
<p>oug lugeliche thing her were rachelis kint Ene unt fênge uan gode rachel. the bodescaf brehte ere gabriel. Thaz ist sagode her gelogen that petrus ande paulus sagon Wande the zuene herren tho zo Roma waren Christus sagode her ere herro was ein drugenere Then thie jüthen uiengen ande an ein cruce gehiengen Wie solden the juthen that gedon mugen Of christus, ere herre godes sun were Woldon mig sagodo her mine uiande uan ik wolde in allon unt gan Ik mohte sagodo her uer suinden under eren handen Hedde ik auar thes geren ik mohte ouer se uêren Wie solden se mir gedâren</p>				
Property		Value		
Word		cruce		
Lemma		KRUIS		
Part of speech + features		NOU-C(number=sg,case=acc,gender=n)		
Word id		w.25269		

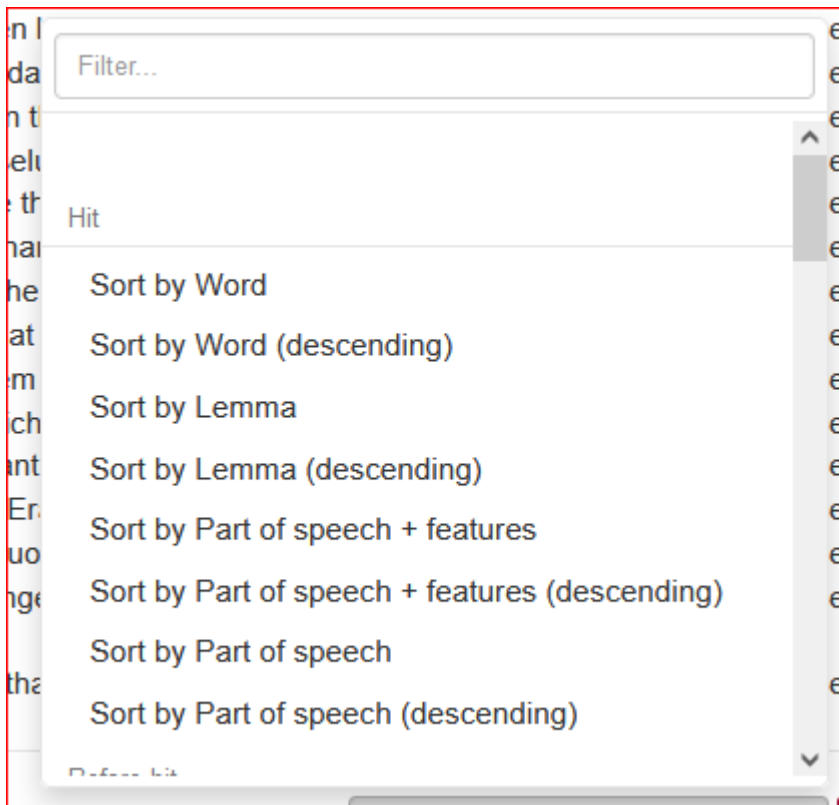
Hit rows are always preceded by a row containing the document title in which those hits occurred, in this case “Mittelfränkische Reimbibel A”. Document titles can be toggled on or off by using the Hide Titles (or Show Titles when titles are hidden) button at the bottom of the page.

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.

Per Hit		Per Document			
Hits				Total hits:	31 (0.0503%)
Group hits by...				Search time:	0s
« 1 2 »					
	Before hit	Hit	After hit	Lemma	Part of speech + features
Mittelfränkische Reimbibel A					
... jüthen uiegen	Word		en Wie solden die juthen...	KRUIS	NOU-C(number=sg,case=acc,gender=n)
... he lieze p	Lemma		paulo that houuet...	KRUIS	NOU-C(number=sg,case=acc,gender=n)
... ce stordo Ande nan	Part of speech + features		iche ande uorde iz insin...	KRUIS	NOU-C(number=sg,case=acc,gender=n)
... imo sagode that	Part of speech		fhat her iz wither...	KRUIS	NOU-C(number=sg,case=acc,gender=n)
... hus her zostorc			usalem uorde [T]ho her...	KRUIS	NOU-C(number=sg,case=acc,gender=n)
... burg porta ze samene slog want her that		cruce	niet otmutliche ne drög Ein...	KRUIS	NOU-C(number=sg,case=acc,gender=n)
... burg porten stont ande hauodo ein		cruce	an siner hant Ande sprag...	KRUIS	NOU-C(number=sg,case=acc,gender=n)
... the sco Ande drog that		cruce	mit grozer uorhten tho offonodo...	KRUIS	NOU-C(number=sg,case=acc,gender=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 (Hit, Before hit, After hit, Date, Localization, Text type and Title and author).



Grouping results

Results Per Hit can be grouped by properties of Hit, Before hit, After hit and of 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. (In the example below we searched for the **lemma** *god* and grouped the hits on **word form**.)

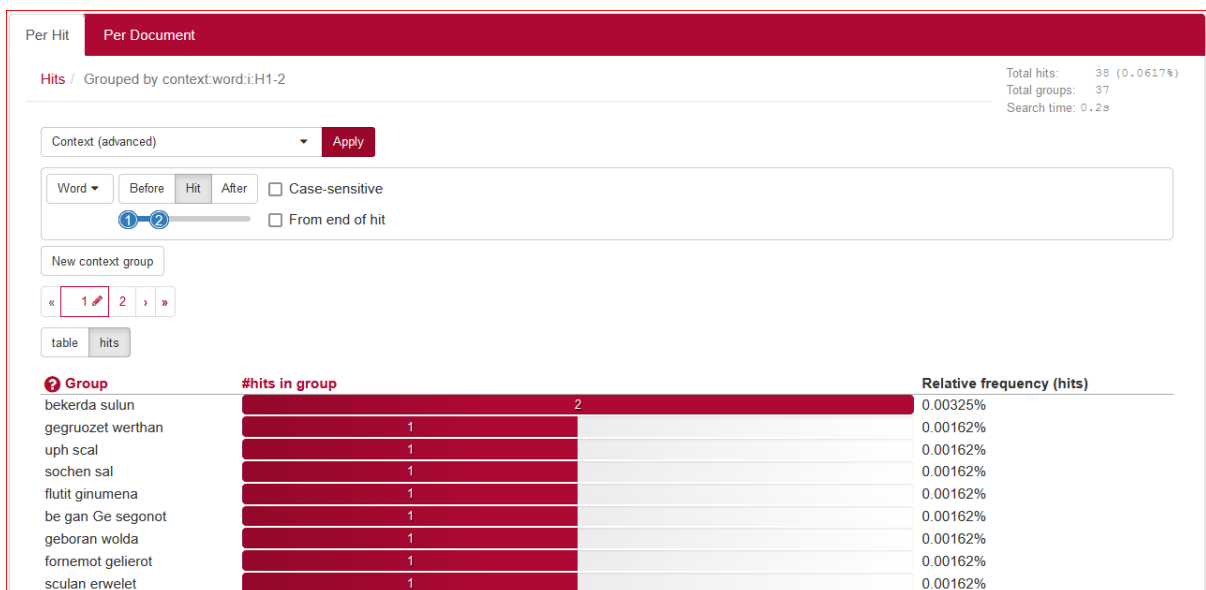


In the Per hit view, advanced grouping options are available by selecting the option Context (advanced). This option allows you to group the results by up to 5 tokens before or after the hit. It also allows you to group the results based on (parts of) the hits. By pressing New context group you can group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

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

Before hit	Hit	After hit	Lemma	Part of speech + features
Mittelfränkische Reimbibel A				
...herren	up wolde	Tho her thaz	OPKEREN	VRB(type=uncl, finiteness=inf, valency=uncl, conjugation=weak)
thaz her	kieren	ge dan	WILLEN	VRB(type=uncl, finiteness=fin, tense=past, mood=conj, number=sg, person=3, valency=uncl, conjugation=irreg)
then buch		hâudo...	OPKEREN	VRB(type=uncl, finiteness=inf, valency=uncl, conjugation=weak)
...hie� her	ge saget	then léuon	ZEGGEN	VRB(type=uncl, finiteness=pastpart, valency=uncl, conjugation=weak irreg)
pilato also	hâuen	Ther grimmo	HEBBEN	VRB(type=uncl, finiteness=fin, tense=pres, mood=ind, number=pl, person=1, valency=uncl, conjugation=weak irreg)
wir	aua	kuning...	AFNEMEN	VRB(type=uncl, finiteness=inf, valency=uncl, conjugation=strong, verbclass=4)
	nemon			
...aua slan	gesaget	her thon in	ZEGGEN	VRB(type=uncl, finiteness=pastpart, valency=uncl, conjugation=weak irreg)
Thit gerihete	haue	themo	HEBBEN	VRB(type=aux, finiteness=fin, tense=pres, mood=ind, number=sg, person=1, valency=uncl, conjugation=weak irreg)
that ik	gebot	nasten...	GEBIEDEN	VRB(type=uncl, finiteness=fin, tense=past, mood=ind, number=sg, person=3, valency=uncl, conjugation=strong, verbclass=2)

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



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 (not visible in the example below).

Group	#hits in group	Relative frequency (hits)
bekerda sulun	2	0.00325%
gegruozet werthan	1	0.00162%
uph scal	1	0.00162%
sochen sal	1	0.00162%
flutit ginumena	1	0.00162%

◀ View detailed concordances

Before	Hit	After
... ummethiga uuerthin also uuahs that	flutit ginumena uuerthunt	ourfiel fuir in ne gesagon ...

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.

Total documents: 7 (7.78%)

Documents

- Sort by hits
- Sort by hits (ascending)

Date

- Sort by Witness Year Date
- Sort by Witness Year (descending) Date
- Sort by Witness Decade Date
- Sort by Witness Decade (descending) Date

Localization

Date	Hits
51-1200	7
51-1200	4
00-1100	8
01-1000	1
01-1000	12
01-1000	5
01-1000	1

Sort by...

Show Hits

Export CSV

Grouping results

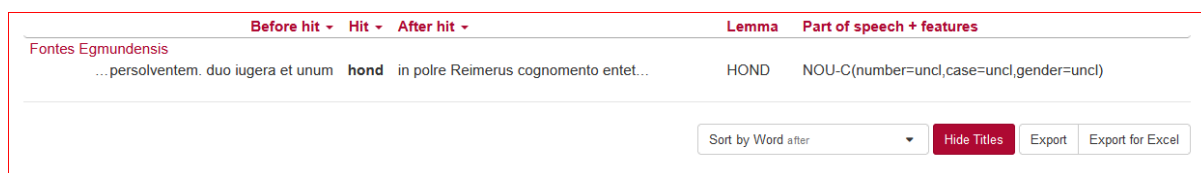
Results Per Document can be grouped by the metadata of the documents in which those hits occur (Date, Localization, Text type, and Title and author). Here, grouping is facilitated by the drop-down menu Group docs by...

Exporting results

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

Information about a document

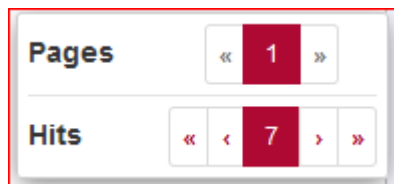
Click on a document title to open the document in a new window - for instance *Fontes Egumndensis* in the screenshot below



The screenshot shows a search results table with columns: 'Before hit', 'Hit', 'After hit', 'Lemma', and 'Part of speech + features'. The first row contains the text 'Fontes Egumndensis' and a snippet of Latin text: '...persolventem. duo iugera et unum **hond** in polre Reimerus cognomento entet...'. The lemma is 'HOND' and the part of speech is 'NOU-C(number=uncl,case=uncl,gender=uncl)'. At the bottom right, there are buttons for 'Sort by Word after', 'Hide Titles', 'Export', and 'Export for Excel'.

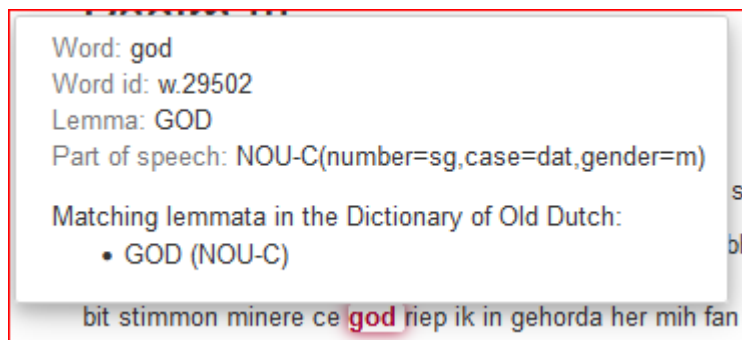
Content

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:



The screenshot shows two navigation controls. The 'Pages' control has a central button with the number '1' and arrows on either side. The 'Hits' control has a central button with the number '7' and arrows on either side.

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:



The pop-up displays the following information:
Word: god
Word id: w.29502
Lemma: GOD
Part of speech: NOU-C(number=sg,case=dat,gender=m)
Matching lemmata in the Dictionary of Old Dutch:
• GOD (NOU-C)
At the bottom, a snippet of text is shown: 'bit stimmon minere ce **god** riep ik in gehorda her mih fan'.

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 corpus. 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.

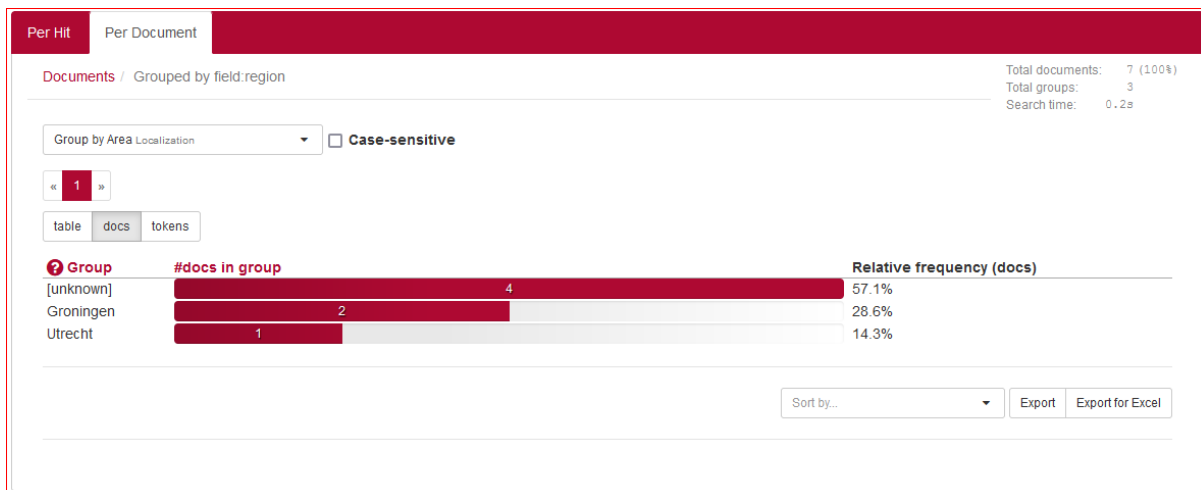
A simple example: suppose we want to obtain information about the Localization of Old Dutch documents dating from 700 to 800 within the *Corpus Oudnederlands*.

- In the Group documents by metadata drop-down menu, choose Group by Area
- In Show groups as, select docs
- In the metadata search form (Filter search by), fill in at Witness Year 700 and 800
- Press 'Search'

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

- Navigation:** 'Search' and 'Explore' tabs at the top.
- Subdivisions:** 'Documents', 'N-grams', and 'Statistics' tabs.
- Grouping:** 'Group documents by metadata' dropdown set to 'Group by Area Localization'.
- Display:** 'Show groups as' dropdown set to 'docs'.
- Filtering:** 'Filter search by ...' dropdown set to 'Date'.
- Witness Year:** Two input fields for '700' and '800'.
- Search Mode:** 'Permissive' and 'Strict' buttons, with 'Strict' selected.
- Summary:** 'Witness Year (Date): 700-800', 'Selected subcorpus: Total documents: 7 (7.78%)', 'Total tokens: 352 (0.571%)'.
- Actions:** 'Search', 'Reset', 'History', and a settings gear icon.

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



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 Part 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 Oudnederlands* for specific metadata.

Example

Within all the documents of the *Corpus Oudnederlands*, you will find two occurrences of this so-called 5-gram:

Per Hit | Per Document

Hits / Grouped by hit:word Total hits: 2 (0.00325%)
Total groups: 2
Search time: 0.2s

Group by Word Case-sensitive

« 1 »

table hits

Group	#hits in group	Relative frequency (hits)						
in fan mannon bluodo behalt	1	0.00162%						
<p>« View detailed concordances »</p> <table border="1"> <thead> <tr> <th>Before</th> <th>Hit</th> <th>After</th> </tr> </thead> <tbody> <tr> <td>... genere mi fan uirkindon unreht</td> <td>in fan mannon bluodo behalt</td> <td>mi vuanda ecco fiengon sela ...</td> </tr> </tbody> </table>			Before	Hit	After	... genere mi fan uirkindon unreht	in fan mannon bluodo behalt	mi vuanda ecco fiengon sela ...
Before	Hit	After						
... genere mi fan uirkindon unreht	in fan mannon bluodo behalt	mi vuanda ecco fiengon sela ...						
in mit mannon ne sulun	1	0.00162%						
<p>« View detailed concordances »</p> <table border="1"> <thead> <tr> <th>Before</th> <th>Hit</th> <th>After</th> </tr> </thead> <tbody> <tr> <td>... an arbeithe manno ne sint</td> <td>in mit mannon ne sulun</td> <td>befilloda uuerthan bethiu hatta sia ...</td> </tr> </tbody> </table>			Before	Hit	After	... an arbeithe manno ne sint	in mit mannon ne sulun	befilloda uuerthan bethiu hatta sia ...
Before	Hit	After						
... an arbeithe manno ne sint	in mit mannon ne sulun	befilloda uuerthan bethiu hatta sia ...						

Sort by...

Statistics (frequency lists)

Here, you can produce frequency lists for a subcorpus. 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 Oudnederlands* for specific metadata

Example

It is possible to determine the use of the most frequently used Old Dutch words in Utrecht in the *Corpus Oudnederlands* by searching for Frequency list type Word and by filtering search by Area (localization):

Search Explore

Explore ...

Documents N-grams Statistics

Frequency list type
Word

Filter search by ...

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

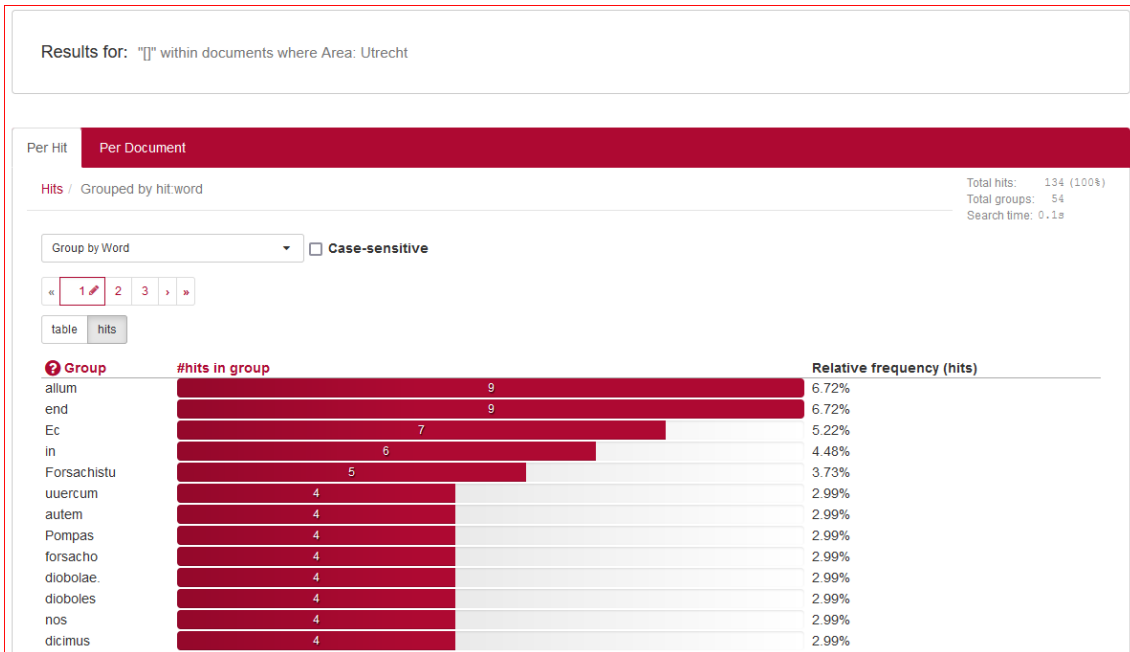
Country
Country

Area
Utrecht

Place
Place

Kloeke location code
Kloeke location code

This results in:



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/accnt-sensitive matching. Note that, unlike in CWB, case-INsensitive matching is currently the default. To explicitly match case/accnt-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!" (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 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.