

Brieven als Buit-2 application manual

Table of contents

| | |
|--------------------------------|----------|
| Introduction | 3 |
| Information about the corpus | 3 |
| GiGaNT Lexicon service | 3 |
| Metadata categories | 4 |
| Letter | 4 |
| Year | 4 |
| Text type | 4 |
| Autograph | 4 |
| Signature | 4 |
| Sender | 4 |
| Name | 4 |
| Gender | 4 |
| Class | 5 |
| Age | 5 |
| Region of residence | 5 |
| Relationship to addressee | 5 |
| Addressee | 5 |
| Name | 5 |
| Place | 5 |
| Country | 6 |
| Region | 6 |
| Ship | 6 |
| Sent from | 6 |
| Application user manual | 7 |
| Getting started | 7 |
| Searching the corpus | 8 |
| Simple search | 8 |
| Search | 8 |
| Wildcards | 9 |
| Reset | 9 |
| History | 9 |
| Global settings | 10 |
| Extended search | 10 |
| Filter search by | 11 |
| Advanced search | 13 |
| The query builder | 13 |

| | |
|--|-----------|
| The tab Search | 13 |
| Token attributes | 14 |
| Adding attributes to a token box | 14 |
| Function of the two +-buttons in a token box | 15 |
| The tab Context | 15 |
| Managing sequences of token boxes | 16 |
| Uploading value lists in the query builder | 16 |
| Copy to CQL editor | 16 |
| Expert search | 16 |
| Copy to query builder | 17 |
| Import query | 17 |
| Gap filling | 17 |
| Viewing results | 19 |
| Per Hit view | 19 |
| Sorting results | 19 |
| Grouping results | 20 |
| Per Document view | 21 |
| Sorting results | 21 |
| Grouping results | 22 |
| Exporting results | 23 |
| Information about a document | 23 |
| Content | 23 |
| Metadata | 23 |
| Statistics | 24 |
| Images | 24 |
| Exploring the corpus | 24 |
| Documents | 24 |
| N-grams | 25 |
| Options | 25 |
| Example | 26 |
| Statistics (frequency lists) | 27 |
| Options | 27 |
| Example | 27 |
| Appendix: Corpus Query Language | 29 |
| CQL support | 29 |
| Supported features | 29 |
| Differences from CWB | 30 |
| (Currently) unsupported features | 31 |
| Using Corpus Query Language | 31 |
| Matching tokens | 31 |
| Sequences | 32 |
| Regular expression operators on tokens | 32 |

| | |
|-----------------------------------|----|
| Punctuation | 32 |
| Case- and diacritics-sensitivity | 33 |
| Matching XML elements | 33 |
| Labeling tokens, capturing groups | 34 |
| Global constraints | 34 |

Introduction

This manual describes the corpus exploitation environment for the *Brieven als Buit-2* ('Letters as loot-2') corpus. The corpus application is developed by the Dutch Language Institute (Instituut voor de Nederlandse Taal or INT). The backend of the application is the BlackLab Lucene based search engine developed for corpora with token-based annotation (<https://blacklab.ivdnt.org/>). The web-based frontend is a further development of the corpus-frontend application developed by INT (<https://github.com/instituutnederlandsetaal/blacklab-frontend>). Its design is inspired by the first version of the OpenSoNaR user interface by Tilburg University and Radboud University (<https://github.com/Taalmonsters/WhiteLab2.0>).

Information about the corpus

Brieven als Buit-2 is a spin-off of the research programme *Brieven als Buit / Letters as Loot. Towards a non-standard view on the history of Dutch* (see www.brievenalsbuit.nl). This programme, initiated and directed by prof. dr. Marijke van der Wal (Leiden University) and funded by the Netherlands Organisation for Scientific Research (NWO), successfully ran from 1 September 2008 till 1 September 2013.

Just as the earlier *Brieven als Buit* internet application, the present *Brieven als Buit-2* comprises Dutch letters which were taken as loot by privateers and confiscated by the High Court of Admiralty during the wars fought between The Netherlands and England from the second half of the 17th to the early 19th centuries. The National Archives (Kew, UK), the current location of the confiscated documents, offered the opportunity to photograph selected letters, all kept in the High Court of Admiralty (HCA) Archives.

The letters were sent home by sailors and others from abroad but also vice versa by those staying behind who needed to keep in touch with their loved ones. These confiscated letters of men, women and even children represent priceless material for historical linguists. They allow us to gain access to the everyday Dutch of the past, the colloquial Dutch of people from the middle and lower classes.

All texts in *Brieven als Buit-2* were tokenised, but unlike the earlier *Brieven als Buit* corpus, the tokens are not tagged with Part of Speech and are not lemmatized. It is therefore only possible to search for word forms.

This first online accessible version of the corpus *Brieven als Buit-2* was released on 26 February 2021.

GiGaNT Lexicon service

Contrary to the earlier *Brieven als Buit* internet application, *Brieven als Buit-2* has not yet been annotated with Part of Speech and Lemma. To make the corpus more accessible, suggestions for query expansion are given, using the INT lexicon service with the historical computational lexicon [GiGaNT-HILEX](#).

The current version of GiGaNT-HILEX in the lexicon service contains the lexicon modules based on the *Dictionary of the Dutch Language (Woordenboek der Nederlandsche Taal, WNT)* and the *Dictionary of Middle Dutch (Middelnederlandsch Woordenboek, MNW)*.

If you want to make use of this service, please contact Katrien Depuydt (katrien.depuydt@ivdnt.org).

Metadata categories

The *Brieven als Buit-2* corpus has been enriched with an elaborate set of metadata categories. These metadata are described below. In the corpus application it is possible to limit a search by filtering on metadata categories.

Letter

A letter from the *Brieven als Buit-2* corpus.

Year

The year(s) in which the letter(s) was (were) written, as evidenced by the date of the letter.

Text type

The type of text to which the letter belongs (business, private, combination). Compared to the original internet application, *Brieven als Buit-2* contains more business letters or letters of a mixed private-business character.

Autograph

In the case of autographs (autograph), it has been established that the sender actually wrote the letter. In the case of non-autographs (non-autograph) the sender did not write the letter himself, but had it done by someone else. In uncertain cases it was not completely certain whereas in unknown cases it could not be determined at all whether the letter is an autograph or not.

Signature

The signature of the box containing letters in the archives of the High Court of Admiralty (HCA) in the National Archives in Kew, United Kingdom. A signature always starts with the letters HCA followed by a series of numbers, for example HCA 30-223.

Sender

The person who sent the letter. Please note that this is not always the actual scribe of the letter (see under Autograph).

Name

The name of the sender.

Gender

The gender of the sender (female, male, unknown).

Class

Four social layers are distinguished, based on the stratification that is common among historians (see Willem Frijhoff & Marijke Spies, 1650. *Bevochten eendracht*. Den Haag: Sdu, 1999, pp. 188-190). The four social segments of the application are: low class, middle-low class, middle-high class and high class.

The low class includes, for example, seafarers from the lowest ranks, servants, soldiers and the poor (11 documents). The middle-low class includes small shopkeepers, small farmers, seafarers from lower ranks and craftsmen (78 documents). The middle-high class consists of, for example, small businessmen, wealthy farmers, master craftsmen, captains and lower ranking officers such as mates (333 documents). The high class includes wealthy merchants, ship owners, academics, senior officials and senior officers in the army and navy (277 documents). It should be noted that high class does not refer to the highest social class of nobility and non-noble ruling classes. That upper class does not occur in the corpus. There are 687 documents of which it is not known to which class the sender belonged.

Age

Three age groups are distinguished, namely <30 (under 30), 30-50 (30 to 50 years) and >50 (over 50). There are 623 letters for which the age of the sender has not been determined (unknown).

Region of residence

This designation refers to the region where a sender grew up or where he or she spent most of his or her life. The region of residence is usually a Dutch province such as Zeeland or Zuid-Holland. The corpus contains many letters from the Caribbean, sent by people who temporarily or for a longer period of time stayed in the Caribbean region, but who originally came from Zeeland or Zuid-Holland. The region of residence is then Zeeland or Zuid-Holland and not, for example, Curaçao. This characteristic is important for linguistic research into dialect differences.

Note that for Noord-Holland, a distinction is made between Noord-Holland, Amsterdam - a city that was a metropolis at the time - and Noord-Holland (excluding Amsterdam).

Relationship to addressee

The relationship, personal and/or professional, that the sender has with the addressee, such as acquaintance, employer, friend, grandson, mother, nephew/cousin or sister.

Addressee

The person to whom the letter was sent.

Name

The name of the consignee.

Place

The place to which a letter was sent, e.g. Enkhuizen or Middelburg.

Country

The country to which a letter has been sent. Note that contemporary names are used, for example Sri Lanka (and not Ceylon) and Saint Kitts (and not St. Christopher).

Region

Within the Dutch language area, the term region refers to a province or dialect area, such as Noord-Brabant, West-Vlaanderen and Zeeland. Outside the Dutch language area, it indicates a geographical region: Azië (Asia), Caraïbisch gebied (Caribbean), Noord-Europa (Northern Europe), West-Afrika (West Africa) and Zuid-Europa (Southern Europe).

Ship

The ship to which a letter was sent, e.g. the *Spiegel*.

Sent from

The four place designations that are distinguished in this tab (Place, Country, Region, Ship) have already been described above at Addressee. Please note that the names of places, countries, regions and ships are different from those found at Addressee.

Application user manual

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

Getting started

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

- To search for a word literally in the form you specify, use Simple search or Extended search.
 - Simple Search for Word [scepe](#)
 - Extended Search for Word [dochter](#)
- To search for words satisfying a certain pattern, use *wildcards* in Simple Search or Extended Search, or *regular expressions* in Advanced Search and Expert Search.
 - words starting with *ver* and ending with *len* in [Simple Search](#)
 - words starting with *ver* and ending with *len* in [Extended Search](#)
 - lemmata starting with *ver* and ending in *eren* with (mostly) one syllable in between in [Advanced Search](#)
 - lemmata starting with *ver* and ending in *eren* with (mostly) one syllable in between in [Expert Search](#)
- To see which unique forms occur as a result of your search, use Group Results.
 - example Group by Annotation: [different words following lieve](#)
 - example Group by Annotation: [different words preceding the word huis](#)
- To explore the distribution of document properties in the corpus, use the Explore feature.
 - example: [characteristics of the signature](#)
 - example: [speaker age distribution](#)

Searching the corpus

Simple search

Search

The Simple Search allows you to quickly search for specific word forms (e.g. *huys*). After entering a search term, a spinner briefly appears on the right side of the search bar. Based on the keyed in word, suggestions are given of possible variants of spelling and/or form from the [GiGaNT-lexicon](#). If you know exactly which word you are looking for, you can also – while the wheel is spinning – press Enter directly. The search will then start immediately.

Based on the information in this lexicon all spelling variants of the search term found are suggested (see the screenshot below). You can then choose from the presented suggestions or select all at the same time (Select all). To make your search even more targeted, it is also possible to limit the search to the parts of speech that were found in the historic component of the GiGaNT-lexicon in connection to the search term.

The screenshot shows the 'Simple Search' interface. At the top, there are tabs for 'Search' and 'Explore'. Below that is a search bar with the text 'Search for ...'. Underneath the search bar are four tabs: 'Simple' (selected), 'Extended', 'Advanced', and 'Expert'. A 'Word' input field contains the text 'huys'. Below the input field are two buttons: 'Select all' and 'Deselect all'. A grid of checkboxes lists various spelling variants of 'huys': hues, huis, huise, huisen, huize, huizen, hus, huse, husen, huss, huus, huys, huysel, huysen, and huyze. Below the grid is a section titled 'Limit to Part of Speech' with two checked options: 'huis (NOU-C)' and 'huizen (VRB)'. At the bottom of the interface are four buttons: 'Search', 'Reset', 'History', and a settings gear icon.

It is also possible to enter a phrase: *voor eeniege dagen* or *watt myn aangaat*. You will then find all occurrences of that exact phrase. Furthermore, you can search for different values simultaneously by separating them without spaces by a vertical line, e.g. *god|man|lief* or - with the use of wildcards - *god|aan*|hond*.

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

Wildcards

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

- * The asterisk matches any character zero or more times. Therefore, searching for *a*n* matches all word forms that start with an *a* and end with a *n*, e.g. *aan*, *adverteeren* and *alschoon*.
- ? The question mark matches a single character once. Therefore, searching for *a?n* matches *only* three-letter values starting with an *a* and ending with a *n*, e.g. *aan*, *aen*, *ahn*, *ann*, *adn* and *aen*.

This wildcard can be used more than once. Thus *a????n* matches words like *armen*, *allen*, *abben* and *agten*.

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.

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).

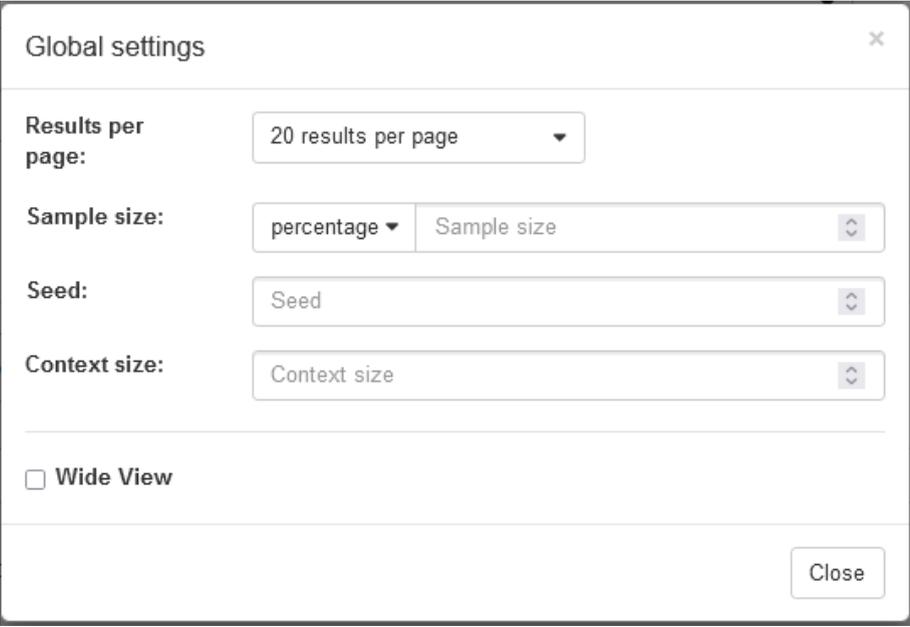


Every search query has its own url. If you copy this url via History (Copy as link) or directly from the address bar of your browser, you can send it to someone else who can import this link via Import from a link. It offers that person the possibility to run the search on his or her own computer.

Global settings

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

- *Results per page*: you can choose whether you want 20, 50, 100 or 200 results to be shown;
- *Sample size*: selecting a value here will instruct the search engine to return a random sample drawn from the complete result set. The sample size can be limited by
 - a percentage of the total number of search results (percentage);
 - the number of results displayed (count).
- *Seed*: a ‘random seed’ is a number used to initialize a so-called pseudo-random number generator. Keeping the same seed will ensure that two samples drawn from the same result set are identical. A new seed will most likely result in a different sample;
- *Context size*: by entering a number you can determine the number of words Before hit and After hit;
- *Wide View*: the default setting is ‘small view’; you can change to Wide View by ticking the checkbox.



The image shows a dialog box titled "Global settings" with a close button (X) in the top right corner. The dialog contains the following settings:

- Results per page:** A dropdown menu set to "20 results per page".
- Sample size:** A dropdown menu set to "percentage" and a text input field labeled "Sample size" with a spinner.
- Seed:** A text input field labeled "Seed" with a spinner.
- Context size:** A text input field labeled "Context size" with a spinner.
- Wide View:** A checkbox that is currently unchecked.
- Close:** A button in the bottom right corner.

Extended search

Like in Simple search, Extended Search allows you to quickly search for specific word forms or phrases. The search is performed in the same way as described for Simple Search.

After entering a search term in the search field Word, a spinner briefly appears on the right side of the search bar. Based on the keyed in word, suggestions are given of possible variants of spelling and/or forms from the [GiGaNt-lexicon](#). If you know exactly which word you're looking for, you can also – while the wheel is spinning – press Enter directly. The search will then start immediately. Like in Simple Search, you can also enter a phrase here.

Based on the information in this lexicon all spelling variants of the search term found are suggested. To make your search even more targeted, it is also possible to limit the search to certain parts of

speech in connection to the search term. You can then choose from the presented suggestions or select all at the same time (Select all). In the screenshot below, all options have been selected.

Search **Explore**

Search for ...

Simple **Extended** Advanced Expert

Word

Select all Deselect all

camer camers kaamer
 kaamers kamer kameren
 kamers

Limit to Part of Speech

kamer (NOU-C)

Case- and diacritics-sensitive

Extended Search allows to search case- and diacritics-sensitive. Note that the default setting for search is case- and diacritics-insensitive. For example, searching for the Word *jan* (& Jan (NOU-P NOU-C)) will result in 1026 occurrences of this name. By ticking the box Case- and diacritics-sensitive you will only find the Word *jan* (149x), but not the variant *Jan*. In order to directly find only occurrences of the Word (form) *Jan* (877x), tick the box Case- and diacritics-sensitive under the search field Word (as shown below).

Search **Explore**

Search for ...

Simple Extended **Advanced** Expert

Word

Select all Deselect all

jan ian jans

Limit to Part of Speech

jan (NOU-C) gunnen (VRB) Jan (NOU-P NOU-C)

Case- and diacritics-sensitive

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

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 Letter (*Year, Text type, Autograph, Signature*), Sender (*Name, Gender, Class, Age, Region of residence, Relationship to addressee*), Addressee

(*Name, Place, Country, Region, Ship*) and Sent from (*Place, Country, Region, Ship*). To view the results for all documents, simply leave the attributes in the filtering form empty.

There are two different ways to specify a filter, depending on the field type. You can either fill in a value yourself - for instance Sender Name - or choose one or more values from a drop-down list - for instance Class. The drop-down list has been applied especially when the number of values to choose from is relatively small. Sender Gender for instance has only three possibilities (female, male and unknown). 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.

The screenshot shows a web interface for filtering search results. At the top, there's a header "Filter search by ...". Below it are four filter categories: "Letter", "Sender" (with a small circle containing the number "1"), "Addressee", and "Sent from". The "Sender" category is expanded, showing a dropdown menu with three options: "female" (which has a checkmark), "male", and "unknown". Below the dropdown, there's a section for "Age" with a search input field.

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

When on the other hand the set of possible values is rather large (e.g. Addressee Name), you have to type a specific value in the 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 *abraham*. 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 (Sender) Name "*abraham bor*" will result in five letters sent by him.

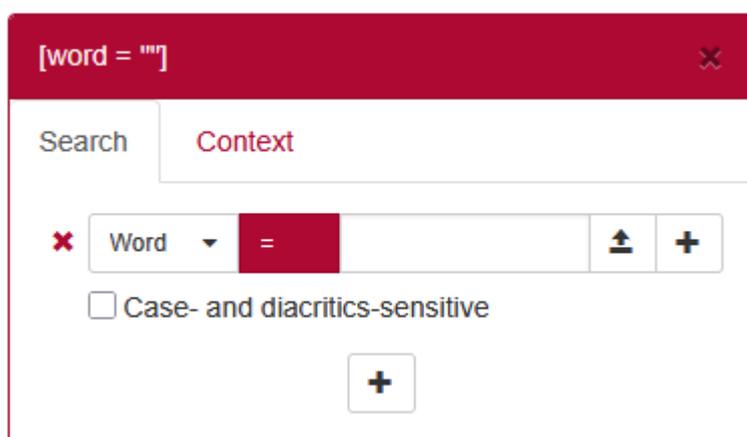
For a detailed description of the metadata, see the section [Metadata categories](#).

Advanced search

The query builder

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

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



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

The tab Search

The tab Search contains a set of attributes a token in the corpus must have to be matched by the query. By clicking the + -button on the right hand side of this token, you can add new attributes (see below). Then enter a value that the attribute must have for the token to be found. The search command Word ‘starts with’ *ge* and Word ‘ends with’ *den* for example results in both verbal forms (*gesonden*, *geworden*, *gelden*) and plural nouns (*geaffectioneerden*, *geliefden*).

It is only possible to search by word forms. However, you can specify whether that word form should be equal or not equal to the entered search term. You can also specify whether or not a word should begin or end with a particular character combination.

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:

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 past participles of strong verbs. First, fill in the attribute Word ‘starts with’ *ge*, then click +, choose *AND*, and choose Word ‘ends with’ *en*.

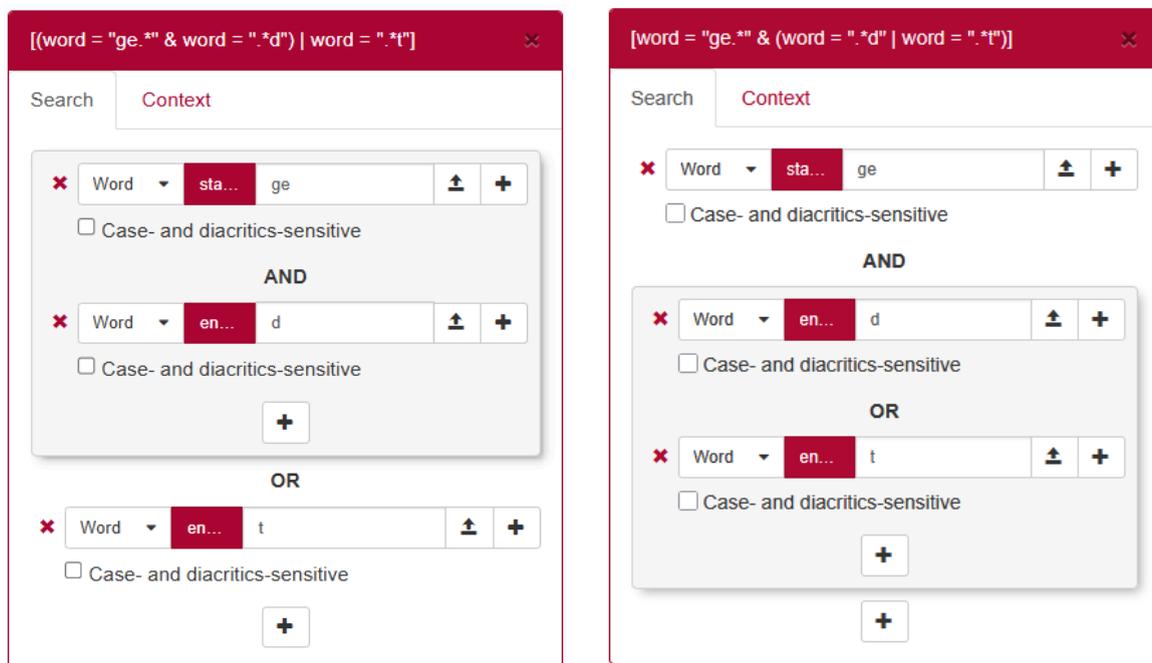
Similarly, creating a new attribute using *OR* will create a token query matching tokens that have the original attribute *or* the new attribute. For instance, enter Word ‘starts with’ *ge*, add a new attribute with the *OR* option and enter Word ‘ends with’ *en* to match tokens as *gezegd*, *gezelschap* and *hebben, ontfangen*.

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.

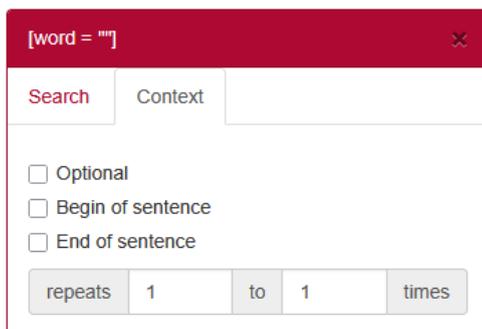
Suppose we want to look for past participles of weak verbs, i.e. verbs that end in an *-d* or a *-t*. If we add the attributes in the order Word 'starts with' *ge* AND Word 'ends with' *d*, OR Words 'ends with' *t* using the +-signs **below** the attributes, as in the left screenshot below, we get the token query [(word = "ge.*" & word = ".*d") | word = ".*t"]. This will also match forms such as *stadt*, *niet*, so this is not what we were after.

If, on the other hand, we add OR Word 'ends with' *t* with the +-sign to the **right** of the attribute Word 'ends with' *d*, it will be inserted in a subclause, thus resulting in the correct query [word = "ge.*" & (word = ".*d" | word = ".*t")], as shown in the right screenshot below.



The tab Context

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



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.

In CQL, spaces only affect a search if they are included in quotes. Whether the search command is [word="schip"] or [word = "schip"] (or just "schip") does not make any difference to the result. However, there is a difference between the queries [word="schip"] and [word=" schip"]. The first search results in exactly 13 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*; [\[word=".*man"\]](#) matches all words ending with *man*, including *man* itself. (Note that [\[word="*man"\]](#) will not give any results, because in Expert Search an asterisk is not a wildcard but a repetition operator.)
- Combination of attributes (combining operators are &, |, !), e.g. [\[word="hoop|geloof|liefde"\]](#) matches either the word *geloof*, the word *hoop* or the word *liefde*.
- The empty [] matches any token, e.g. [\[word="man"\]\[\]{}3\[word="god"\]](#) matches a sequence of *man* followed by *god* with three arbitrary tokens in between.
- Operators |, & and parentheses () and the repetition operators (+, *, ? and {}) can be used to build complex sequence queries. Example: ["lieue" "man" | "almagtige" "god"](#), or even [\("lieue" "man" | "almagtige" "god"\)+](#), matching any sequence of *lieue man* or *almagtige god*. Note that, while most queries up to this point could also have been constructed with the query builder, we really need the power of CQL from here on.

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

Copy to query builder

When the query is relatively simple – like [\[word="schip"\]\[word="den"\]](#) – 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) [word="schip"][[word="den"]that has the same properties.

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

If, for instance, you are interested in the distribution of words that can be placed between two specific words you can create this query in the Corpus Query Language field:

```
[word="@@"][][word="@@"]
```

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: ⓘ

```
[word="@@"][][word="@@"]
```

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 (titles are hidden):

Per Hit Per Document

Hits Total hits: 108 (0.0151%)
Search time: 0.02s

Group Results

« 1 2 3 6 »

| Before Hit | Hit | After Hit |
|--|--|--|
| ...vreesen is ik wens Dat ...verder Zo wens ik dat ...de Westjse als Wybren nannes ...Tjebbe Jans de Groot voerende ...coomen sullen ende soo dra ...te deelen ik heb met ...berigten Als Dat ik in ...het Doen van Assurantie in | de verdraagsam god de Almagtige God het t Schip het Coffe schip het een schip het laaste Schip het voorsz schip het meergemelde schip | niet ut stort syn gramschyp... uw wil Seegnen met uw... Verlooren dog tSy binne daar... De twee Gesussters tot La... of het ander lukt om... communucacie brief weegens het huwelijk... De Leendert Matthys voor Deszelfs... De Leendert Matthijs heb ik... |

This mimics the functionality to upload a list of values in the Advanced Search interface.

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

Viewing results

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

Per Hit view

Click a hit - i.e. a line with the bold words in the column Hit - to display the properties and values of the hit (in the following example **het laeste schip**). Click the hit again to close.

Document id: bab2121

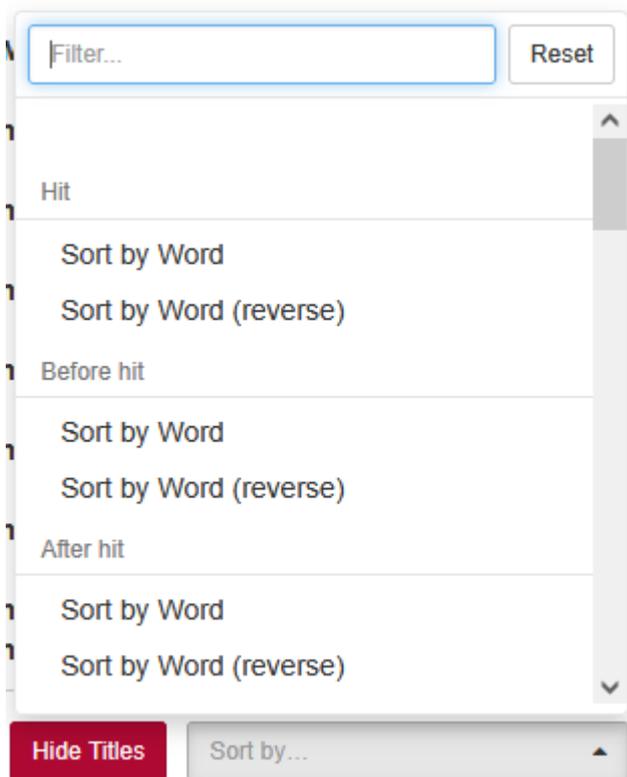
| Before Hit | Hit | After Hit | | | | |
|--|-------------------------|---------------------------------|----------|-------|------|------------------|
| ...suijker te kookken sijn met | het laeste schip | uijt selant Gekoomen soo dat... | | | | |
| <p>...senden alsoo het schip op morgen a oovermorgen sal vertrecken maer hoope met het eerste schip dat naer dit vertreckt een paer duijsent pont oover te senden ende hier hebb ick toe Coomende 23 feberwaeri drie duijsent pont aen de plantaesi tegoet de ketels om suiijker te kookken sijn met het laeste schip uijt selant Gekoomen soo dat ick met den eersten mede u sijlen maeken sal versoeck V E seer vriendelijck van dit boovenstaende en wijnich te assisteeren ende ick hoope met den eersten het u e weder te restituereeren ende versoecke V E mijne stoutich heijt te vereckskuseereen ende mede den...</p> <table border="1"> <thead> <tr> <th>Property</th> <th>value</th> </tr> </thead> <tbody> <tr> <td>Word</td> <td>het laeste schip</td> </tr> </tbody> </table> | | | Property | value | Word | het laeste schip |
| Property | value | | | | | |
| Word | het laeste schip | | | | | |

Hit rows are always preceded by a row containing the document title in which those hits occurred, in this case *To Cornelis van Bullaart, 12 januari 1672 by Johannes Stratius*. The document titles can be toggled on or off by using the Hide Titles (or Show Titles when titles are hidden) button at the bottom of the page. If you hover the mouse over the title, the identification number of the document appears, in this case: bab2121.

Sorting results

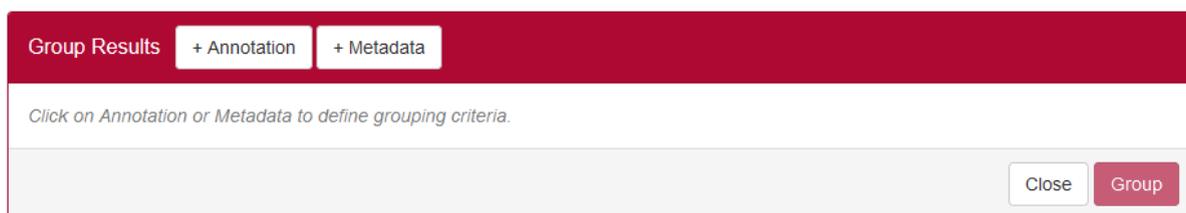
Click on any of the column headings to sort the hits on Words within that column, clicking again inverts the sorting.

You can also sort the results by means of the drop-down menu at the bottom of the page (Sort by ...), which offers you the possibility to sort by various attributes for Hit, Before hit, After hit, Letter, Sender, Addressee, Sent from.



Grouping results

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



Results can be grouped by Annotation and by Metadata.

By clicking +Annotation you can group by the first word, by all words or by specific words, whether before the hit, within the hit or after the hit, and based on the annotation Word. When grouping by the first word or specific words, you can also group from the end of the hit. The default grouping is grouping all words within the hit using annotation Word. Clicking +Metadata allows you to group by metadata assigned to the document (Letter, Sender, Addressee, Sent from).

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

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

Group Results + Annotation + Metadata

first Word before hit ✕

I want to group on the first word ▾ before the hit ▾ using annotation Word ▾

Case sensitive:

o dat daar nog **Een** | Schip | vertrekt dan vind ik mijn
e dal daar nog Een Schip vertrekt dan vind ik mijn

Clear Group

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

| Group | #hits in group | Relative frequency (hits) |
|-------|----------------|---------------------------|
| het | 377 | 0.0525% |
| t | 179 | 0.025% |

◀ View detailed concordances Load more concordances

| Before Hit | Hit | After Hit |
|------------------------------------|--------------|--|
| ...t'advertieren dat dan in t | schip | d'hoop Capt Ian mejer heb... |
| ...Eijlandts tabacq blaeden pr t | schip | de hoop schippr jan mejer... |
| ...den 18 nouemb pr t | schip | van Jan bosman waer Jn... |
| ...Middelburgh in Zeelandt Pr t | Schip | De Negotie Zeevaard Schipper Pieter... |
| ...niet min fatiguant geweest t | schip | was in een Gestadige beweeging... |
| ...Van UEd ontffange Per t | Schip | het Witte Paart Capt Obet... |
| ...Armaud a Amsterdam Pr t | Schip | Koningin Esther Capn Jan Nosten... |
| ...en wynig verversching van t | Schip | Soo alls tik selfs koop... |
| ...als Wybren nannes het t | Schip | Verlooren dog tSy binne daar... |
| ...3 voet waatter in t | schip | haden en tot meer ongeluk... |
| ...Coopliede a Dordregt pr t | Schip | de Goede Vrienden Captn Dirk... |
| ...de brieve zak by t | schip | ongeopend moet blyven dus zult... |
| ...misnoeging gegeeven heeft met t | schip | aan fesquet te adresseeren zal... |
| ...risico afgelaaden hebt per t | schip | de Neutraliteit capn Jan HK... |
| ...Capn Jan Hansen voerende t | Schip | d Hoopende Zeeman abzent aan... |
| ...domse Coffijbe gelaaden in t | Schip | Hoofdenburg Capn S:J Lourens voor... |
| ...dus na te gaan t | Schip | zwaar heeft geleeden hebbe vandaag... |
| ...zelfde Schip laaden want t | Schip | is goed de Capn is... |
| ...Amstel tot Amsterdam Pr t | schip | Negotie en Zeevaard Schippr Pieter... |
| ...Staat in handen Per t | Schip | de Negotie en Zeevaard Schippr... |

◀ View detailed concordances Load more concordances

| | | |
|-----|----|----------|
| een | 78 | 0.0109% |
| ons | 45 | 0.00627% |
| int | 29 | 0.00404% |

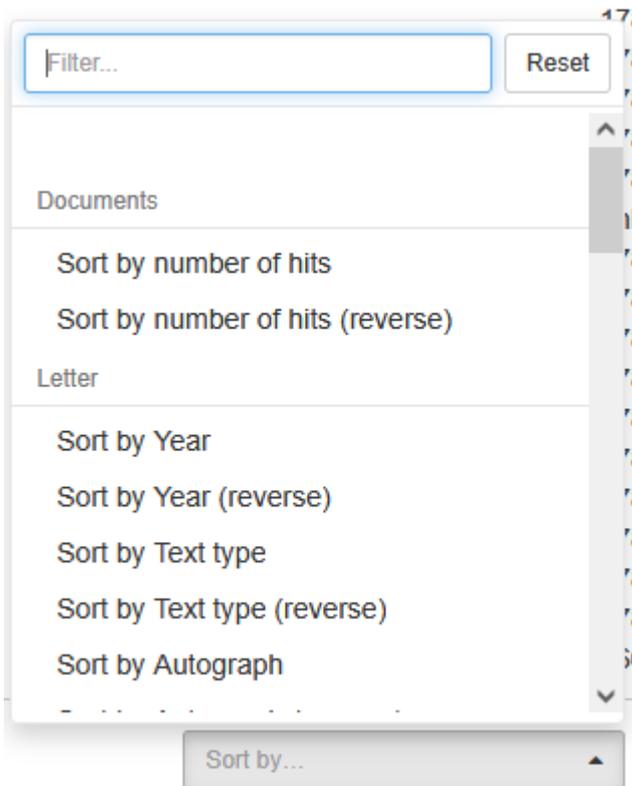
Click on View detailed concordances to go back to the normal hits view to see more detailed information for the hits in this group. The button Go back to grouped results brings you back to the list of groups.

Per Document view

Sorting results

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

You can also sort the results by means of the drop-down menu at the bottom of the page (Sort by...), which offers you the possibility to sort by various attributes such as Hit (Documents), Signature (Letter) and Class (Sender).



Grouping results

Results Per Document can be grouped by metadata assigned to the document (Letter, Sender, Addressee and Sent from). The example below shows all documents in which the Word *engels* occurs grouped by year.

Results for: [word="engels"] within all documents

Per Hit | Per Document

Documents / Grouped by Document Year Total documents: 29 (2.09%)
Total groups: 6
Search time: 0.01s

Group Results + Metadata

Document Year ✕

Select the metadata to group on.

Group by Year ▼

Case sensitive:

Clear Group

« 1 » table docs

| Group | #docs in group | Relative frequency (docs) |
|-------|----------------|---------------------------|
| 1780 | 9 | 31% |
| 1664 | 6 | 20.7% |
| 1781 | 6 | 20.7% |
| 1665 | 5 | 17.2% |
| 1782 | 2 | 6.9% |
| 1783 | 1 | 3.45% |

Exporting results

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

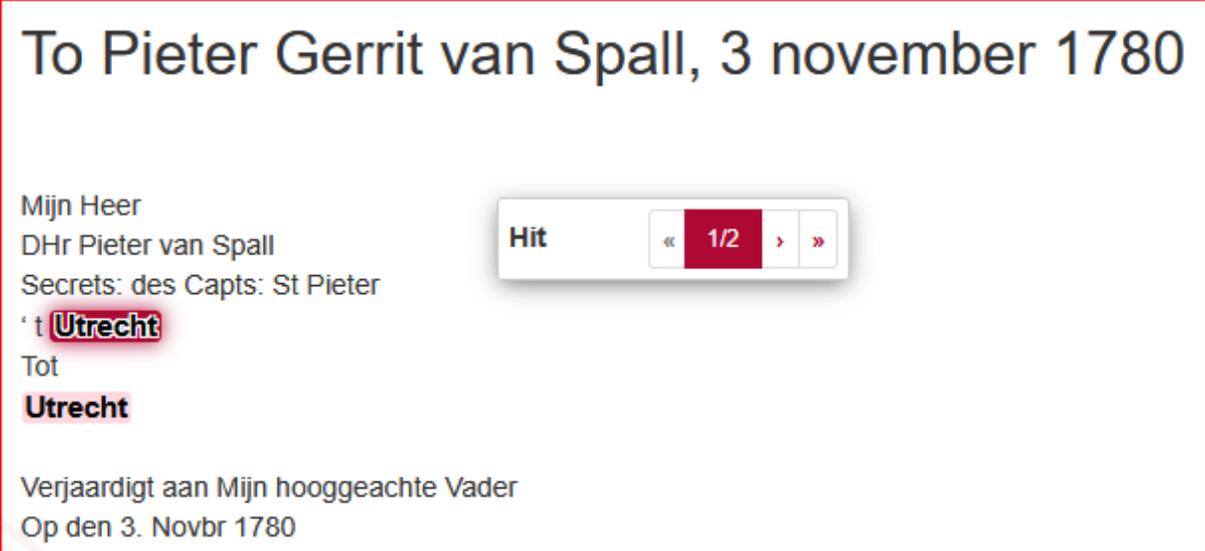
Grouped results can be exported in the same way. However, if you would like to have the metadata with each concordance of a group, you must first click on the red bar of a specific group and then on View detailed concordances. The results you then see can be exported by the use of the Export buttons.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

Information about a document

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

Content

Hits from the current query will be highlighted in bold in the opened document. In the case of several hits only the current hit will also appear in shadow (such as *Utrecht* in the example below). You can navigate from one hit to another by using the arrows at the Hits button (this button can be dragged around):



To Pieter Gerrit van Spall, 3 november 1780

Mijn Heer
DHr Pieter van Spall
Secrets: des Capts: St Pieter
t **Utrecht**
Tot
Utrecht

Verjaardigt aan Mijn hooggeachte Vader
Op den 3. Novbr 1780

Hit « 1/2 » »

Metadata

In the Metadata tab all metadata properties of the document are displayed. They provide information about the Letter, the Sender, the Addressee and Sent from, as well as the Document length (tokens).

Statistics

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

Images

Under images you can find a photo of the original letter, kept in the National Archives (Kew, UK).

Exploring the corpus

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

Documents

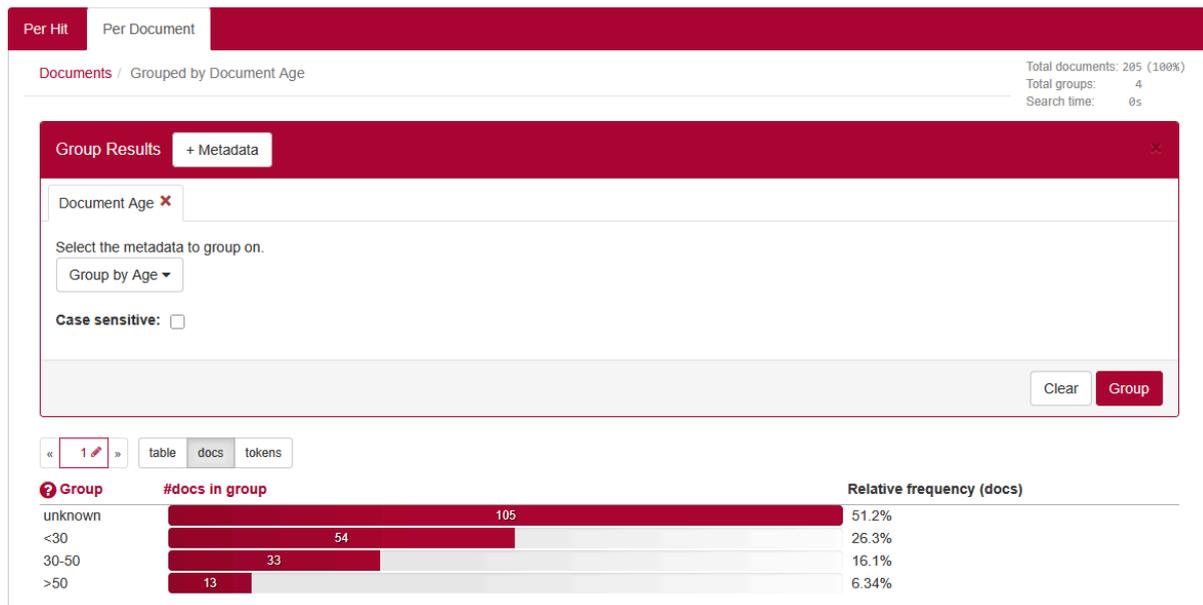
This subtab allows you to investigate the documents. It consists of two drop-down menus to specify the grouping of the metadata and to specify the way the groups are to be shown.

A simple example: suppose we want to obtain information about the age distribution of female senders within the *Brieven als Buit-2* corpus.

- In the Group documents by metadata drop-down menu, choose Group by Age (Sender)
- In Show groups as, select *Docs*
- In the metadata search form (Filter search by), select in Gender (Sender) *female*
- Press Search

The screenshot shows the 'Explore' interface with a dark red header containing 'Search' and 'Explore' tabs. Below the header, the 'Explore ...' section has three sub-tabs: 'Documents' (selected), 'N-grams', and 'Statistics'. Under 'Documents', there are two dropdown menus: 'Group documents by metadata' set to 'Group by Age' and 'Show groups as' set to 'Docs'. Below this is the 'Filter search by ...' section with four sub-tabs: 'Letter', 'Sender' (selected and marked with a '1'), 'Addressee', and 'Sent from'. Under 'Sender', there are two dropdown menus: 'Name' set to 'Name' and 'Gender' set to 'female'. A large, faint watermark 'Prison' is visible across the bottom half of the interface.

You will get this result:



N-grams

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

Options

- *N-gram size*: the length of the sequence (a number from 1 to 5; default setting is 5).
- *N-gram type*: the attribute to search for. You can choose: Word (i.e. word form). If you do not specify the search term a series of arbitrary words equal to the n-gram size will be searched for.
- It is also possible to restrict to, for instance, n-grams with some slots already specified, as is shown in the following example. After entering a search term, a spinner briefly appears on the right side of the search bar. Based on the keyed in word, suggestions are given of possible variants of spelling and/or form from the GiGaNT-lexicon and of parts of speech. By clicking on 'Select all' all forms belonging to a GiGaNT lemma are added.
- By using the Filter search by ... you can create a subcorpus within the corpus *Brieven als Buit-2* for specific metadata.

Example

The screenshot shows the 'Explore ...' interface with the following settings and results:

- Search: Explore
- Documents | N-grams | Statistics
- N-gram size: 5
- N-gram type: Word
- Word 1: het
- Word 2: Word
- Word 3: kint|kind|kijnder|kijndt|kijnt|ki
- Word 4: Word
- Word 5: Word
- Select all | Deselect all (for both words)
- Limit to Part of Speech:
 - het
 - hebben (VRB)
 - het (INT)
 - het (PD ART)
 - het (PD)
 - heten (VRB)
- Results list:
 - kint
 - kind
 - kijnder
 - kijndt
 - kijnt
 - kinde
 - kinden
 - kinder
 - kindere
 - kinderen
 - kinders
 - kindre
 - kindren
- Limit to Part of Speech (for results):
 - kind (AA)
 - kind (NOU-C)

Within all the documents of the *Brieven als Buit-2* corpus, you will find 12 occurrences of this so-called 5-gram (choose the option ‘Select all’ for both words and limit your search to Part of Speech *het* (ART) respectively *kind* (NOU-C)).

Per Hit Per Document

Hits / Grouped by Word within hit Total hits: 12 (0.00167%)
Total groups: 12
Search time: 0.003s

Group Results + Annotation + Metadata

Word within hit ✕

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

Case sensitive:

het Zelve Bij UED teffens |
 het Lieve kind Mag Continueeren |
 Was het anders het Zoude
het Zelve Bij UED teffens het Lieve kind Mag Continueeren Was het anders het Zoude

Clear Group

« 1 » table hits

| Group | #hits in group | Relative frequency (hits) |
|------------------------------------|----------------|---------------------------|
| het nablyvende kind te meer | 1 | 0.000139% |
| het lieve kind die grootste | 1 | 0.000139% |
| het Lieve kind Mag Continueeren | 1 | 0.000139% |
| het lieven kint met openen | 1 | 0.000139% |
| het een kint van een | 1 | 0.000139% |
| het Lieve Kint Nagt Dag | 1 | 0.000139% |
| het kleinste kijnt tomas van | 1 | 0.000139% |
| het Lieve kind dat jk | 1 | 0.000139% |
| het lieve kind dirk Marinus | 1 | 0.000139% |
| het jonggeboore kind Waarlijk voor | 1 | 0.000139% |
| het een kint de groetenis | 1 | 0.000139% |
| het lieve kint in de | 1 | 0.000139% |

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*: in this corpus, it is only possible to create frequency lists of Words (i.e. word forms).
- By using the Filter search by, you can create a subcorpus within the corpus *Brieven als Buit-2* for specific metadata.

Example

It is possible to determine the use of the most frequent words by female writers younger than 30 from the high class by searching for Frequency list type Word and by filtering search by Gender: *female*, Class: *high* and Age: *<30*. This results in:

Results for: Word frequency within documents where Age: <30, Gender: female, Class: high

Per Hit Per Document

Hits / Grouped by Word within hit

Total hits: 10.240 (100%)
Total groups: 2.754
Search time: 0.008s

Group Results + Annotation + Metadata

Word within hit

I want to group on using annotation

Case sensitive:

 Mejuffrouw E M Franssen Ten
Mejuffrouw E M Franssen Ten

« 1 2 3 4 6 11 »

| Group | #hits in group | Relative frequency (hits) |
|-------|----------------|---------------------------|
| en | 311 | 3.04% |
| de | 251 | 2.45% |
| dat | 230 | 2.25% |
| ik | 230 | 2.25% |
| van | 204 | 1.99% |
| het | 161 | 1.57% |
| mijn | 136 | 1.33% |
| te | 130 | 1.27% |
| voor | 129 | 1.26% |
| is | 123 | 1.2% |
| als | 119 | 1.16% |

Appendix: Corpus Query Language

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

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

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

CQL support

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

Supported features

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

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

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

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

Differences from CWB

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

For now, here's what you should know:

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

(Currently) unsupported features

The following features are not (yet) supported:

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

Using Corpus Query Language

Matching tokens

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

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

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

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

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

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

```
"man"
```

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

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

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

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

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

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

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

Sequences

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

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

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

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

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

Regular expression operators on tokens

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

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

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

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

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

Or, for two or more adjectives:

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

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

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

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

Punctuation

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

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

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

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

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

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

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

Case- and diacritics-sensitivity

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

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

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

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

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

Matching XML elements

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

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

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

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

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

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

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

The above query will just match the word "baker" as part of a person's name. But you're likely more interested in the entire name that contains the word "baker". So, to find those full names, use:

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

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

```
<person/>
```

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

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

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

Labeling tokens, capturing groups

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

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

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

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

BlackLab also supports numbered groups:

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

Global constraints

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

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

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