

CLARIAH Chaining Search: A Platform for Combined Exploitation of Multiple Linguistic Resources

Peter Dekker, Mathieu Fannee, Jesse de Does
Dutch Language Institute (INT)

October 1, 2019

CLARIN Annual Conference 2019
Leipzig, Germany

- Perspective: linguist (with some computational knowledge)
- Goal: combine heterogeneous CLARIN resources: corpora, lexica, treebanks
- Current flawed options
 - Access web interfaces
 - Download entire dataset

- Corpus Gysseling: 13th century Dutch
- Web interface:
<http://gysseling.corpus.taalbanknederlands.inl.nl/gysseling/page/search>

The screenshot displays the Gysseling web interface. At the top, there is a navigation bar with the logo and links for Home, Instituut voor Nederlandse Lexicologie, and CLARIN. The main header features the text 'DE NEDERLANDSE' in a large, light blue font. Below this, there are search filters for 'Simple' and 'OQL query'. The search area is divided into two columns: 'Search for...' and 'Filter search by'. The 'Search for...' column includes fields for 'Wordform' (containing 'man'), 'Lemma', and 'P.o.S.'. The 'Filter search by' column includes fields for 'Title', 'Author', 'Date' (with 'From' and 'To' sub-fields), and 'Source'. A 'Search' button and a 'Reset' button are located below the search fields. A 'Show me:' dropdown menu is set to '50 results'. Below the search area, there is a navigation bar with tabs for 'Per Hit', 'Per Document', 'Hits grouped', and 'Documents grouped'. The 'Per Hit' tab is selected. A pagination bar shows 'Prev' and 'Next' buttons, with a 'Toggle list' button. The main content area displays a table of search results with columns for 'Left context', 'Hit text', 'Right context', 'Lemma', and 'Part of speech'. The table contains several rows of text, each starting with the word 'man' in bold, followed by a period and a description of a context. The 'Lemma' column contains the word 'MAN' and the 'Part of speech' column contains the corresponding grammatical category.

Left context	Hit text	Right context	Lemma	Part of speech
... gheleent waren. Ende vore sine	man	ende sine soepene, die hier na ...	MAN	NOUtype-common,number-pl,inflexion=0
... die oec mins her wouters	man	es. Hier ouer waren oec andre ...	MAN	NOUtype-common,number-sg,inflexion=0
... Hier ouer waren oec andre sine	man	die oec sine scepene sijn. ...	MAN	NOUtype-common,number-pl,inflexion=0
... dese letten ghehanghen. Ende wi	man	ende scepene mins her wouters ...	MAN	NOUtype-common,number-pl,inflexion=0
... alsoelken dienst also onse trouwe	man	haar jan die here van ...	MAN	NOUtype-common,number-sg,inflexion=0
... Onser maeressen ende wijsdome onser	man	Dier scuylech op waren de ...	MAN	NOUtype-common,number-pl,inflexion=0
... ombre de meere Sekerheit, wi	man	hier vore gheruut sijn ombre dat wi ...	MAN	NOUtype-common,number-pl,inflexion=0

Current option: Download entire dataset

- Corpus Gysseling: 13th century Dutch
- Download dataset: <https://ivdnt.org/taalmaterialen/102-taalmaterialen/2014-tstc-corpus-gysselingh>

```
<f name="type"><symbol value="eigen"/></f><f name="building"><symbol value="met-s-of-th"/></f><f name="pos"><symbol value="N"/></f>
<f name="lemma"><string>HENDRIK</string></f></f>
<f type="001" lemma="ZON" pos="NOUN[types=common,number=sg,inflection=0]" xml:id="w.482761" msd="N(soort,ev,met-e)" n="482761"><seg type="orth">zoeene.</seg><f name="0" type="ml">
<f name="type"><symbol value="soort"/></f><f name="getal"><symbol value="ev"/></f><f name="building"><symbol value="met-e"/></f><f name="pos"><symbol value="N"/></f>
<f name="lemma"><string>ZOOKE</string></f></f>
<f type="020" lemma="MILHELM" pos="NOUN[types=proper,inflection=0]" xml:id="w.482762" msd="N(eigen,zonder)" n="482762"><seg type="orth">wille-expan resp=editor" xmlns="http://www.tei-c.org/ns/1.0">
<f name="type"><symbol value="eigen"/></f><f name="building"><symbol value="zonder"/></f><f name="pos"><symbol value="N"/></f>
<f name="lemma"><string>MILHELM</string></f></f>
<f type="700" lemma="VAN" pos="ADD[types=general,inflection=0]" xml:id="w.482763" msd="VZ(zonder)" n="482763"><seg type="orth">van.</seg><f name="0" type="ml">
<f name="building"><symbol value="zonder"/></f><f name="pos"><symbol value="VZ"/></f>
<f name="lemma"><string>VAN</string></f></f>
<f type="020" lemma="KLEIEM" pos="NOUN[types=proper,inflection=0]" xml:id="w.482764" msd="N(eigen,zonder)" n="482764"><seg type="orth">kleihem.</seg><f name="0" type="ml">
<f name="type"><symbol value="eigen"/></f><f name="building"><symbol value="zonder"/></f><f name="pos"><symbol value="N"/></f>
<f name="lemma"><string>KLEIEM</string></f></f>
<f type="020" lemma="JOHANNES" pos="NOUN[types=proper,inflection=0]" xml:id="w.482765" msd="N(eigen,zonder)" n="482765"><seg type="orth">jan.</seg><f name="0" type="ml">
<f name="type"><symbol value="eigen"/></f><f name="building"><symbol value="zonder"/></f><f name="pos"><symbol value="N"/></f>
<f name="lemma"><string>JOHANNES</string></f></f>
<f type="006" lemma="WOND" pos="NOUN[types=common,number=sg,inflection=0]" xml:id="w.482766" msd="N(soort,ev,zonder)" n="482766"><seg type="orth">wont.</seg><f name="0" type="ml">
<f name="type"><symbol value="soort"/></f><f name="getal"><symbol value="ev"/></f><f name="building"><symbol value="zonder"/></f><f name="pos"><symbol value="N"/></f>
<f name="lemma"><string>WOND</string></f></f>
<f type="020" lemma="GERARD" pos="NOUN[types=proper,inflection=0]" xml:id="w.482767" msd="N(eigen,zonder)" n="482767"><seg type="orth">gherard.</seg><f name="0" type="ml">
<f name="type"><symbol value="eigen"/></f><f name="building"><symbol value="zonder"/></f><f name="pos"><symbol value="N"/></f>
<f name="lemma"><string>GERARD</string></f></f>
<f type="471" lemma="DE" pos="PD[types=art,subtype=def,inflection=0]" xml:id="w.482768" msd="LID(bep,met-e)" n="482768"><seg type="orth">die.</seg><f name="0" type="ml">
<f name="wtype"><symbol value="bep"/></f><f name="building"><symbol value="met-e"/></f><f name="pos"><symbol value="LID"/></f>
<f name="lemma"><string>DE</string></f></f>
<f type="006" lemma="MEVEL" pos="NOUN[types=common,number=sg,inflection=0]" xml:id="w.482769" msd="N(soort,ev,zonder)" n="482769"><seg type="orth">wueel.</seg><f name="0" type="ml">
<f name="type"><symbol value="soort"/></f><f name="getal"><symbol value="ev"/></f><f name="building"><symbol value="zonder"/></f><f name="pos"><symbol value="N"/></f>
<f name="lemma"><string>MEVEL</string></f></f>
<f type="020" lemma="IMETH" pos="NOUN[types=proper,inflection=0]" xml:id="w.482770" msd="N(eigen,zonder)" n="482770"><seg type="orth">weinin.</seg><f name="0" type="ml">
<f name="type"><symbol value="eigen"/></f><f name="building"><symbol value="zonder"/></f><f name="pos"><symbol value="N"/></f>
<f name="lemma"><string>IMETHIN</string></f></f>
<f type="700" lemma="VAN" pos="ADD[types=general,inflection=0]" xml:id="w.482771" msd="VZ(zonder)" n="482771"><seg type="orth">van.</seg><f name="0" type="ml">
<f name="building"><symbol value="zonder"/></f><f name="pos"><symbol value="VZ"/></f>
<f name="lemma"><string>VAN</string></f></f>
<f type="020" lemma="VARSNARE" pos="NOUN[types=proper,inflection=0]" xml:id="w.482772" msd="N(eigen,met-a)" n="482772"><seg type="orth">versen-expan resp=editor" xmlns="http://www.tei-c.org/ns/1.0">
<f name="type"><symbol value="eigen"/></f><f name="building"><symbol value="met-a"/></f><f name="pos"><symbol value="N"/></f>
<f name="lemma"><string>VARSNARE</string></f></f>
<f type="010" lemma="SCHEPEN" pos="NOUN[types=common,number=pl,inflection=0]" xml:id="w.482773" msd="N(soort,pl,zonder)" n="482773"><seg type="orth">scepenen.</seg><f name="0" type="ml">
<f name="type"><symbol value="soort"/></f><f name="getal"><symbol value="pl"/></f><f name="building"><symbol value="zonder"/></f><f name="pos"><symbol value="N"/></f>
<f name="lemma"><string>SCHEPEN</string></f></f>
<f type="700" lemma="VAN" pos="ADD[types=general,inflection=0]" xml:id="w.482774" msd="VZ(zonder)" n="482774"><seg type="orth">van.</seg><f name="0" type="ml">
<f name="building"><symbol value="zonder"/></f><f name="pos"><symbol value="VZ"/></f>
<f name="lemma"><string>VAN</string></f></f>
<f type="474" lemma="DE/HEI" pos="PD[types=art,subtype=def,inflection=0]" xml:id="w.482775" msd="LID(bep,met-n)" n="482775"><seg type="orth">den.</seg><f name="0" type="ml">
<f name="wtype"><symbol value="bep"/></f><f name="building"><symbol value="met-n"/></f><f name="pos"><symbol value="LID"/></f>
<f name="lemma"><string>DE/HET</string></f></f>
<f type="104" lemma="VRIJ" pos="ADJ[number=sg,inflection=0]" xml:id="w.482776" msd="ADJ(ev,met-n)" n="482776"><seg type="orth">vriec.</seg><f name="0" type="ml">
<f name="getal"><symbol value="ev"/></f><f name="building"><symbol value="met-n"/></f><f name="pos"><symbol value="ADJ"/></f>
<f name="lemma"><string>VRIJ</string></f></f>
<f type="204" lemma="DOEN" pos="VB[types=main, finiteness=finite,tense=present,inflection=0]" xml:id="w.482777" msd="WI(hoofd,pl,tw,met-n)" n="482777"><seg type="orth">doen.</seg><f name="0" type="ml">
<f name="wtype"><symbol value="hoofd"/></f><f name="worm"><symbol value="ov"/></f><f name="building"><symbol value="tw"/></f><f name="pos"><symbol value="N"/></f>
```

- CLARIAH: digital research infrastructure for arts and humanities in The Netherlands
- Levels of searchability:
 - Local search: resource available as web service
 - Federated search: resources of same type queried as single resource
 - Chaining search: heterogeneous sources combined, sequential search workflows

A Python library and Jupyter web interface to easily combine exploration of linguistic resources published in the CLARIN/CLARIAH infrastructure, such as corpora, lexica and treebanks.

- Large SPARQL query to query multiple resources
- But:
 - Not all resources available as Linked Open Data, queried with SPARQL
 - Query becomes highly complex

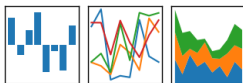
```
values ?n_ontolex_writtenRep { "{word}" } .
?n_entry a ontolex:LexicalEntry .
?n_entry ud:pos ?udobj . ?udobj rdfs:label ?udpos .
?n_form a ontolex:Form .
?n_sense a ontolex:LexicalSense .
?n_syndef a diamant:SynonymDefinition .
?n_sensedef a lemon:SenseDefinition .
?n_syndef diamant:definitionText ?n_syndef_definitionText .
?n_sensedef diamant:definitionText ?n_sensedef_definitionText .
?n_entry ontolex:canonicalForm ?n_form .
?n_entry ontolex:sense ?n_sense .
?n_sense lemon:definition ?n_syndef .
?n_sense lemon:definition ?n_sensedef .
?n_sense diamant:attestation ?n_attest_show .
?n_sense diamant:attestation ?n_attest_filter .
?n_attest_show diamant:text ?n_q_show .
?n_attest_filter diamant:text ?n_q_filter .
?n_attest_show a diamant:Attestation .
?n_attest_filter a diamant:Attestation .
?n_q_filter a diamant:Quotation .
?n_q_show a diamant:Quotation .
?n_q_filter diamant:witnessYearFrom ?wy_f_filter .
?n_q_filter diamant:witnessYearTo ?wy_t_filter .
?n_q_show diamant:witnessYearFrom ?wy_f_show .
?n_q_show diamant:witnessYearTo ?wy_t_show .
FILTER (xsd:integer(?wy_f_show) >= {{timelineStart}})
FILTER (xsd:integer(?wy_t_show) >= {{timelineStart}})
FILTER (xsd:integer(?wy_f_show) <= {{timelineEnd}})
FILTER (xsd:integer(?wy_t_show) <= {{timelineEnd}})
bind("searchLemma" as ?resultType) .
} UNION
```

- Python library
 - Based on pandas DataFrames [McKinney \(2011\)](#)
 - Four modules: search, process, ui and utils
- Jupyter notebooks
 - Examples notebook
 - Sandbox: start coding yourself
 - Sailing letters case study



pandas

$$y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$$



- Corpora
 - CLARIN Federated Content Search endpoint (Stehouwer et al., 2012)
 - Queried via CQL
 - More than 1200 freely available corpora
 - BlackLab corpus search engine (De Does et al., 2017)
 - Queried via CQL
 - More search options, e.g. metadata filtering
- Lexica
 - Linked Open Data, *Ontolex* lexicon model (McCrae et al., 2017)
 - Queried via SPARQL
 - Internal INT *lexicon service* API
- Treebanks
 - CGN (corpus of spoken Dutch) and Lassy (written Dutch)
 - Queried via XPath query

Search in linguistic resource:

```
results = create_lexicon(lexicon_name).lemma(  
    word).search()
```

```
results = create_corpus(corpus_name).pattern(  
    pattern).search()
```

```
results = create_treebank(treebank_name).  
    pattern(pattern).search()
```

Create results table (Pandas DataFrame) with keywords-in-context:

```
df = results.kwic()
```

`utils` provides functions general operations applied to search results tables.

```
diff = column_difference( df1["lemma_0"  
                             ], df2["lemma_0"] )
```

`process` handles linguistic processing of data.

```
df_lexicon = extract_lexicon(df_corpus  
)
```

Operations regarding showing the user interface and loading/saving data

```
save_dataframe(df, "test.csv")  
df = load_dataframe("test.csv")
```

Look for adjectives which should have ending -e, but miss it with determiner *een*.

First, search in corpus.

```
df_corp = create_corpus("openchn").  
    pattern(' [ pos="DET"&lemma="een" ] [  
    word = ". * [ ^ e ] $ " _ & _ pos = "AA . * degree =  
    pos . * " ] [ pos = "NOU . * gender = [ fm ] . * " ] ' )  
    . search () . kwic ()
```

Search in lexicon.

```
df_lex = create_lexicon("molex").lemma  
        ( '(.+)[^e]$' ).pos( 'ADJ(degree=pos)'  
        ).search().kwic()  
final_e_condition = df_filter(df_lex["  
        wordform"], 'e$')  
df_lexicon_form_e = df_lex[  
        final_e_condition ]
```

Filter corpus using results from lexicon.

```
e_forms = set(df_lexicon_form_e.lemma)
cond = df_filter(df_corp["word_1"],
                pattern=e_forms, method="isin")
result_df = df_corp[ cond ]
```


Jupyter Examples Last Checkpoint: Yesterday at 11:58 AM (unsaved changes) Logout

File Edit View Insert Cell Kernel Widgets Help Trusted | env

```
In [15]: from chaininglib.search.CorpusQuery import *

query= corpusQueryField.value
corpus_name = corpusField.value

# BEWARE: we limit the results to 500 records here; not limiting may cause a search to take quite some time
df_corpus = create_corpus(corpus_name).pattern(query).max_results(500).search().kwic()
display_df(df_corpus, labels="Results")
```

Results

	left context	lemma 0	pos 0	word 0	right context
0	lijeye man stelt alles te	boek	NOU	boek	waet ghij uijt gef dat
1	alzoo hij niet op de	boek	NOU	bouck	en stondt en hij heeft
2	Schrijfpampier a 6 Srt	boek	NOU	boek.	50 ditto ongsneeden ditto a
3	rollen matten MR & Eenige	boek	NOU	boeken	Maar also die twee Scheepen
4	te Senden, & wat de	boek	NOU	boeken	aan betreft Zo Gelleff maar
5	laaten wagten tot ik mijn	boek	NOU	boeken	heeft, want weet anders niet
6	& houw nu alle de	boek	NOU	boeken	& bonte & die sijn
7	gantsche regel die in myn	boek	NOU	Boek	niet te vinden is, en

Letters as Loot corpus (Van der Wal et al., 2012):
17th and 18th century letters from Dutch
sailors, annotated with metadata



- Research question: *Which vocabulary is specific for:*
 - *Social class* (low or high)
 - *Gender* (male or female)
 - *Time period* (17th or 18th century)
- Use CLARIAH chaining search to:
 1. Retrieve data from corpus
 2. Filter and split data on metadata (class/gender/era)
 3. Compute relative frequencies of words per metadata category
 4. Compute differences (ratios) between relative frequencies
- Notebook [Case_study_paper.ipynb](#) in GitHub repository

Lemmata most specific for time period (highest diff in rel frequency)

17th century				18th century			
lemma	relative frequency			lemma	relative frequency		
	17th	18th	diff		17th	18th	diff
huisvrouw	0.016	0.002	1.013	heer	0.006	0.021	0.985
vriend	0.023	0.011	1.012	mijnheer	0.009	0.018	0.991
man	0.020	0.010	1.010	edele	0.000	0.008	0.992
gezondheid	0.016	0.007	1.010	jaar	0.005	0.012	0.993
goedenacht	0.009	0.000	1.009	mejuffrouw	0.001	0.006	0.995
schipper	0.009	0.000	1.009	kapitein	0.012	0.017	0.995
brief	0.022	0.014	1.008	zuster	0.008	0.013	0.995
suiker	0.010	0.002	1.008	achting	0.000	0.005	0.995
monsieur	0.008	0.001	1.006	familie	0.002	0.006	0.996
schip	0.023	0.017	1.006	liefde	0.001	0.006	0.996

Lemmata most specific for social class (highest diff in rel frequency)

lemma	low		
	low	high	diff
goedenacht	0.026	0.009	1.017
hart	0.027	0.012	1.015
man	0.033	0.020	1.013
brief	0.034	0.022	1.012
Heer	0.022	0.011	1.011
gezondheid	0.026	0.016	1.009
zuster	0.016	0.008	1.008
huisvrouw	0.024	0.016	1.008
kind	0.020	0.013	1.007
zoon	0.016	0.011	1.005

lemma	high		
	low	high	diff
kapitein	0.005	0.012	0.993
monsieur	0.000	0.008	0.993
suiker	0.002	0.010	0.993
sinjeur	0.001	0.006	0.995
mijnheer	0.004	0.009	0.995
schipper	0.005	0.009	0.996
heer	0.002	0.006	0.996
pond	0.000	0.004	0.997
december	0.002	0.005	0.997
rekening	0.001	0.004	0.997

Lemmata most specific for gender (highest diff in rel frequency)

lemma	male		
	relative frequency		
	male	female	diff
suiker	0.002	0.000	1.002
schip	0.024	0.022	1.002
oom	0.003	0.002	1.001
vriend	0.022	0.021	1.001
december	0.002	0.001	1.001
maand	0.003	0.002	1.001
huis	0.004	0.003	1.001
port	0.002	0.001	1.001
vaderland	0.002	0.001	1.001
stuk	0.001	0.000	1.001

lemma	female		
	relative frequency		
	male	female	diff
man	0.033	0.041	0.992
brief	0.034	0.038	0.996
hart	0.027	0.031	0.996
Heer	0.022	0.023	0.998
zoon	0.016	0.018	0.998
allerliefste	0.005	0.006	0.998
mens	0.005	0.006	0.998
gezondheid	0.026	0.027	0.999
huisvrouw	0.024	0.025	0.999
reis	0.010	0.011	0.999

Lemmata most specific for gender (highest diff in rel frequency)

lemma	male		
	relative frequency		
	male	female	diff
suiker	0.010	0.001	1.009
monsieur	0.008	0.001	1.006
vriend	0.023	0.019	1.005
sinjeur	0.006	0.002	1.004
december	0.005	0.001	1.004
goed	0.011	0.007	1.004
vracht	0.003	0.000	1.003
mr.	0.003	0.000	1.003
cargasoen	0.003	0.000	1.003
dienaar	0.002	0.000	1.000

lemma	female		
	relative frequency		
	male	female	diff
man	0.020	0.046	0.975
brief	0.022	0.034	0.988
kind	0.013	0.022	0.991
zoon	0.011	0.018	0.993
goedenacht	0.009	0.016	0.993
zuster	0.008	0.015	0.994
gezondheid	0.016	0.022	0.994
dochter	0.004	0.010	0.995
genade	0.008	0.013	0.995
hart	0.012	0.017	0.995

- Case study shows difference in vocabulary across sociolinguistic variables, especially for gender in higher social class
- Case study proves: CLARIAH Chaining search facilitates and accelerates linguistic research
- Discussion
 - Minimum of programming knowledge is needed
 - Long computations: invoke library directly, without notebook
 - Future work: optimization for very large data sets

- Demonstration during poster session

- Try it yourself:

<https://github.com/INL/chaining-search>

- Local install
 - Cloud instance on Azure
- Full API reference in documentation:

<https://chaining-search.readthedocs.io/en/latest/>

References

- De Does, J., Niestadt, J., and Depuydt, K. (2017). Creating Research Environments with BlackLab. In Utrecht University, NL and Odijk, J., editors, *CLARIN in the Low Countries*, pages 245–257. Ubiquity Press.
- McCrae, J. P., Bosque-Gil, J., Gracia, J., Buitelaar, P., and Cimiano, P. (2017). The ontolex-lemon model: development and applications. In *Proceedings of eLex 2017 conference*, pages 19–21.
- McKinney, W. (2011). pandas: a Foundational Python Library for Data Analysis and Statistics. *Python for High Performance and Scientific Computing*, 14:9.

Stehouwer, H., Durco, M., Auer, E., and Broeder, D. (2012). Federated Search: Towards a Common Search Infrastructure. *LREC 2012*, page 5.

Van der Wal, M. J., Rutten, G., and Simons, T. (2012). Letters as loot: Confiscated Letters filling major gaps in the History of Dutch. In Dossena, M. and Del Lungo Camiciotti, G., editors, *Pragmatics & Beyond New Series*, volume 218, pages 139–162. John Benjamins Publishing Company, Amsterdam.