

**Zoekacties en gebruikte codes  
binnen  
COREX (CORpusEXploratiesoftware)**

- Vragen, opmerkingen of suggesties: [servicedesk@inl.nl](mailto:servicedesk@inl.nl)
- Informatie over het CGN: <http://www.inl.nl/tst-centrale/nl/producten/corpora/corpus-gesproken-nederlands/6-17>
- CGN-webcursus: [http://www.inl.nl/tst-centrale/images/stories/producten/webcursus/cgn\\_webcursus/](http://www.inl.nl/tst-centrale/images/stories/producten/webcursus/cgn_webcursus/)

# 1 Metadata search

## 1.1 Session Keys (via Session)

<i>CGN.wordCount</i>	number of words in the sample {number}
<i>CGN.recCount</i>	duration of the sample expressed in the total number of seconds {number}
<i>CGN.recDate</i>	recording date {(date or) year}
<i>CGN.fon.available</i>	whether there is a manually verified phonetic annotation available or not {true, false}
<i>CGN.syn.available</i>	whether there is a syntactic annotation available or not {true, false}
<i>CGN.pro.available</i>	whether there is a prosodic annotation available or not {true, false}

## 1.2 Participant Keys (via Session/Participant)

<i>CGN.age</i>	age class to which speaker belonged at the time the sample was recorded;	
	<b>age0</b>	under 18 years of age
	<b>age1</b>	18-24 years
	<b>age2</b>	25-34 years
	<b>age3</b>	35-44 years
	<b>age4</b>	45-55 years
	<b>age5</b>	over 55 years of age
	<b>ageX</b>	age unknown
<i>CGN.birth.year</i>	speaker's year of birth	
<i>CGN.birth.place</i>	speaker's place of birth, represented in terms of the first three digits of the postal code preceded by the country code (BE, NL)	
<i>CGN.birth.region</i>	(geographical) region where the speaker was born. For a list of regions distinguished, see below.	
<i>CGN.residence.place</i>	speaker's place of residence, represented in terms of the first three digits of the postal code preceded by the country code	
<i>CGN.residence.reg</i>	(geographical) region where the speaker resides. For a list of regions distinguished, see below.	
<i>CGN.residence.size</i>	indication of the (present) size of the place where the speaker resided while (s)he was between 4 and 16 years of age	
<i>CGN.education.place</i>	place where speaker lived for the most part between ages 4 and 16 represented in terms of the first three digits of the postal code preceded by the country code (BE, NL)	
<i>CGN.education.opleiding</i>	type of education; e.g. lager onderwijs, mbo, universiteit	
<i>CGN.education.reg</i>	(geographical) region where speaker lived while (s)he attended secondary education (region where speaker lived for the most part between ages 4 and 16). For a list of regions distinguished,	

	see below.	
<i>CGN.education.level</i>	level of education	
	<b>edu1</b>	high
	<b>edu2</b>	middle
	<b>edu3</b>	low
	<b>edu4</b>	unknown

### 1.3 Geographical regions:

<b>regN1a</b>	NL, central region, Zuid-Holland, excl. Goeree Overflake
<b>regN1b</b>	NL, central region, Noord-Holland, excl. West-Friesland
<b>regN1c</b>	NL, central region, West-Utrecht, incl. the city of Utrecht
<b>regN2a</b>	NL, transitional region, Zeeland, incl. Goeree Overflakee and Zeeuws-Vlaanderen
<b>regN2b</b>	NL, transitional region, Oost-Utrecht, excl. stad Utrecht
<b>regN2c</b>	NL, transitional region, Gelders rivierengebied, incl. Arnhem and Nijmegen
<b>regN2d</b>	NL, transitional region, Veluwe up to the river IJssel
<b>regN2e</b>	NL, transitional region, West-Friesland
<b>regN2f</b>	NL, transitional region, Polders
<b>regN3a</b>	NL, peripheral region 1 (north east), "Achterhoek"
<b>regN3b</b>	NL, peripheral region 1 (north east), Overijssel
<b>regN3c</b>	NL, peripheral region 1 (north east), Drenthe
<b>regN3d</b>	NL, peripheral region 1 (north east), Groningen
<b>regN3e</b>	NL, peripheral region 1 (north east), Friesland
<b>regN4a</b>	NL, peripheral region 2 (south), Noord-Brabant
<b>regN4b</b>	NL, peripheral region 2 (south), Limburg
<b>regNx</b>	NL, unknown
<b>regV1</b>	Flanders, central region (Antwerpen and Vlaams-Brabant)
<b>regV2</b>	Flanders, transitional region (Oost-Vlaanderen)
<b>regV3</b>	Flanders, peripheral region 1 (West-Vlaanderen)
<b>regV4</b>	Flanders, peripheral region 2 (Limburg)
<b>regVx</b>	Flanders, unknown
<b>regW</b>	Wallonia
<b>regZ</b>	region known to be outside of The Netherlands and Flanders
<b>regX</b>	region unknown

## 1.4 Content Keys (via Session/Content/Keys)

<i>CGN.textclass.target</i>	gives information about four aspects: text type (specifies the component to which a sample belongs), degree of preparedness, mode and domain	
<b>Text types</b>		
<b>tta</b>	spontaneous conversations (face-to-face)	
<b>ttb</b>	interviews with teachers of Dutch	
<b>ttd</b>	spontaneous telephone dialogues (recorded on MD with local interface)	
<b>tte</b>	simulated business negotiations	
<b>ttf</b>	interviews/discussions/debates (broadcast)	
<b>ttg</b>	(political) discussions/debates/meetings (non-broadcast)	
<b>tth</b>	lessons recorded in a classroom	
<b>tti</b>	live (e.g. sport) commentaries (broadcast)	
<b>ttj</b>	news reports/reportages (broadcast)	
<b>ttk</b>	news (broadcast)	
<b>ttl</b>	commentaries/columns/reviews (broadcast)	
<b>ttm</b>	ceremonious speeches/sermons	
<b>ttn</b>	lectures/seminars	
<b>tto</b>	read speech	
<b>Degree of preparedness</b>		
<b>prep1</b>	scripted,	
<b>prep2</b>	unscripted,	
<b>prep3</b>	more-or-less scripted;	
<b>mode</b>		
<b>mod1</b>	broadcast, radio;	
<b>mod2</b>	broadcast, tv;	
<b>mod3</b>	non-broadcast;	
<b>Domain</b>		
<b>dom1</b>	private	
<b>dom2</b>	public	

## 2 Content Search

### 2.0 Reguliere expressies

Het is binnen de Corex-tool “content search” mogelijk om je zoekactie te definiëren met behulp van een reguliere expressie. Een **reguliere expressie** is een manier om patronen te beschrijven waarmee een computer tekst kan herkennen. Dat is handig wanneer je bijvoorbeeld wilt zoeken op alle woorden die met “over” beginnen, zonder dat je al die woorden afzonderlijk moet invoeren.

Variabele	Betekenis
.	matches an arbitrary character
+	matches the preceding pattern element one or more times
?	matches the previous character zero or one times
*	matches the previous character zero or more times
^	matches the beginning of a word
\$	matches the end of a word
[adg]	matches the character a, d or g
[^ad]	matches any character except a and d

**BELANGRIJK:** Wanneer je een zoekstring ingeeft, dan zoekt Corex niet exact op die zoekstring, maar naar alle woorden die die zoekstring bevatten. Het is niet aan te raden om variabelen zoals “.” voor of na je zoekstring te zetten, want dat levert verkeerde resultaten op.

Wanneer je wil aangeven dat je zoekstring aan het begin of aan het einde van een woord moet voorkomen, kun je gebruikmaken van de tekens ^ en \$.

Voorbeelden van zoekacties wanneer “Reg. Exp.” aangevinkt is:

- **bel** vindt bel, gebeld, bellen, babbelen, dubbel,...
- **^bel** vindt bel, bellen, belt, belangrijke,...
- **bel\$** vindt bel, dubbel, heibel,...
- **^b.l\$** vindt bel, bil, bol, bal,...
- **^b.\*l\$** vindt bel, bil, bol, bal, bedoel, bewijsmateriaal, boksschool,...
- **^b[ae]l\$** vindt bel en bal
- **^.it\$** vindt uit, zit, dit,...

### 2.1 Orthography, marked, lemma

Codering	Betekenis
<i>Dialect</i>	Dialectwoord
<i>Foreign</i>	Vreemd/buitenlands woord
<i>Incomplete</i>	Afgebroken woord
<i>Mispr</i>	Verkeerd uitgesproken woord of klanknabootsing
<i>Regionalpr</i>	Zwaar dialectisch uitgesproken woord
<i>uncertain</i>	Woorden waarvan de transcribent niet zeker is of hij/zij ze goed gehoord heeft

## 2.2 Fonetische transcriptie

### 2.2.1 CGN-symbolen

Klasse	Voorbeeld	CGN-symbool	Klasse	Voorbeeld	CGN-symbool
Plosieven	<u>p</u> ut	p	Korte vocalen	li <u>p</u>	l
	<u>b</u> ad	b		le <u>g</u>	E
	<u>t</u> ak	t		la <u>t</u>	A
	<u>d</u> ak	d		bo <u>m</u>	O
	<u>k</u> at	k		pu <u>t</u>	Y
	<u>g</u> oal	g		Lange vocalen	lie <u>p</u>
Fricatieven	<u>f</u> iets	f	<u>b</u> uur	y	
	<u>v</u> at	v	le <u>g</u>	e	
	<u>s</u> ap	s	de <u>u</u> k	2	
	<u>z</u> at	z	la <u>t</u>	a	
	<u>s</u> jaal	S	bo <u>o</u> m	o	
	<u>r</u> avage	Z	bo <u>o</u> k	u	
	<u>x</u> icht	x	Sjwa	ge <u>l</u> ijk	@
	<u>g</u> egen	G	Diftongen	wi <u>js</u>	E+
Sonoranten	<u>h</u> eel	h	<u>h</u> uis	Y+	
	<u>n</u> ang	N	<u>k</u> oud	A+	
	<u>m</u> at	m	Leenvocalen	sc <u>è</u> ne	E:
	<u>n</u> at	n	<u>f</u> reule	Y:	
	<u>j</u> oranje	J	<u>z</u> one	O:	
	<u>l</u> at	l	Nasale vocalen	vacc <u>in</u>	E~
	<u>r</u> at	r	<u>c</u> roiss <u>ant</u>	A~	
	<u>w</u> at	w	<u>c</u> ong <u>é</u>	O~	
	<u>j</u> as	j	<u>p</u> arfum	Y~	

### 2.2.2 Zoekcodes voor de fonetische transcripties

Code	Betekenis
<i>SEP</i>	gescheiden (er wordt geen foneem gedeeld over de woordgrenzen heen)
<i>SHARE-P</i>	een gedeeld plosief (plofklank) over woordgrenzen heen
<i>SHARE-NP</i>	een gedeeld non-plosief over woordgrenzen heen
<i>SHARE-W</i>	uitzondering: het woord/de woorden 'da's'
<i>INSERT</i>	een geïnserteerd foneem tussen twee woorden

## 2.3 Part of speech

Afkorting woordsoort	Woordsoort
ADJ	Adjectief
BW	Bijwoord
LET	Leesteken
LID	Lidwoord
N	Substantief
SPEC	speciale woorden (vb. delen van eigennamen, woorden in zelfnoemfunctie,...)
TSW	Tussenwerpsel
TW	Telwoord
VG	Voegwoord
VNW	Voornaamwoord
VZ	Voorzetsel
WW	Werkwoord

Het is mogelijk om binnen de content search te zoeken op woordsoort (vb. WW) of op een volledige tag (vb. N(soort,mv,basis), maar het is ook mogelijk om een zoekactie te doen op basis van gedeeltelijke POS-tags.

*Voorbeeld:* er zijn meerdere tags voor onbepaalde lidwoorden, nl. LID(onbep, dial), LID(onbep,gen,evf) en LID(onbep,stan,agr).

Om alle uitingen te vinden waarin onbepaalde lidwoorden voorkomen, kun je drie afzonderlijke zoekacties doen, maar het is ook mogelijk om met één zoekopdracht alle uitingen met onbepaalde lidwoorden te vinden.

Dit kan door met de rechtermuisknop te klikken op het vak na LID waarin je reeds een van de tags voor onbepaalde lidwoorden geselecteerd hebt, bijvoorbeeld (onbep,stan,agr). Na de klik wordt het vak wit en kan je de cursor in het vakje plaatsen. Plaats de cursor vóór dat deel van de tag dat je wil deleten, bv. voor "stan". Vervolgens druk je op de deleteknop en dan wordt de tag vanaf "stan" vervangen door een sterretje: \*. Er staat nu als tag LID(onbep,\*). Wanneer je een zoekactie start worden alle onbepaalde lidwoorden gezocht.

The screenshot shows a search filter configuration area. It includes a label 'Annotation on tier' followed by a dropdown menu set to 'part of speech'. Next is the label 'matches' followed by a dropdown menu set to 'LID'. To the right of 'LID' is a text input field containing '(onbep,\*)' and a small dropdown arrow. Below these are two checkboxes: 'Reg. Exp.' and 'Case Sens.', both of which are unchecked. At the bottom of the configuration area are two buttons: 'Add new constraint' and 'Delete last constraint'.

**Let op:** alleen het laatste deel van een tag kan verwijderd worden.

## 2.4 Prosodische annotatie

De volgende prosodische fenomenen maken deel uit van de prosodische annotatie:

- **prosodische grenzen** (een onderbreking in het spraaksignaal):
  - o zwakke grenzen zijn per definitie aangeduid in een manuele annotatie en binnen Corex kregen die de naam **weak**
  - o sterke grenzen werden aangeduid in een manuele annotatie (strong) of in een automatische annotatie (strong)
- **klemtoon**
- **abnormale klankverlenging**

De volgende zoekacties zijn mogelijk binnen de prosodische annotatie (zie ook hieronder):

- je kan bepalen of de grens die je zoekt links of rechts van een woord moet staan (leftb en rightb)
- je kan op een bepaald aantal grenzen zoeken (nweakb en nstrongb)
- je kan zoeken op een bepaald aantal klemtonen (nprom) al dan niet in combinatie met een woord
- je kan zoeken op een bepaald aantal klankverlengingen (nlength) al dan niet in combinatie met een woord

