



Korpus og leksikografi

Leksikografane i
revisjonsprosjektet for *Bokmålsordboka* og *Nynorskordboka*
<https://ordbok.uib.no>





Informasjonskjelder i revisjonsprosjektet

- tekstkorpus (lemmautval, tyding, syntaktisk åtferd, bruksdøme)
- Språkrådet (normering)
- Språksamlingane ved UiB (særleg *Metaordboka*)
- introspeksjon
- andre ordbøker





Relevante verktøy for å søkje i korpus

- CLARINO Bergen Center:
 - Corpuscle (verktøy for søk på ordnivå)
 - INESS (verktøy for søk i setningsanalysar)
- Nasjonalbiblioteket: NB n-gram og ordgalaksar
- Retriever: Atekst (avistekst)
- Universitetet i Oslo: HaBiT (norsk webkorpus)





Corpuscle – søk og analyse (ordnivå)

- <http://clarino.uib.no/korpuskel/>
- Grensesnitt for å tilgjengeleggjere, søkje i og analysere tekstkorpus.
- Inngår i CLARINO Centre Bergen (Universitetsbiblioteket ved UiB).





”Corpuscle-Lex” (eigen server)

Korpusnamn	Språk	Storleik (# ord)	Tidsrom	Sjanger	Lemma + ordklasse	Tilgang
Talk Of Norway	bokmål. nynorsk	63,8 mill	1999-2016	sakprosa	ja	open
Aviskorpus (bokmål)	bokmål	1509,1 mill	1998-2015	sakprosa (avis)		open
NBs frie tekster (bokmål)	bokmål	516,4 mill	1765-2013	blanda: sakprosa, skjønnlitteratur		open
Leksikografisk bokmålskorpus	bokmål	102,3 mill	1985-2013	blanda: sakprosa, skjønnlitteratur	ja	avgrensa
Aviskorpus annotert	bokmål	29,0 mill	2001-2009	sakprosa (avis)	ja	open
Forskning.no (2017)	bokmål	21,5 mill	1998-2017	sakprosa (avis)	ja	avgrensa
Nynorskkorpus	nynorsk	107,8 mill	1866-2012	blanda	ja	avgrensa
NBs frie tekster (Nynorsk)	nynorsk	46,2 mill	1850-2010	blanda		open
Aviskorpus (nynorsk)	nynorsk	16,1 mill	1998-2015	sakprosa (avis)		open



Lemmutval: trunkerte søk

Søk Refine window: 5 tokens | Stop | Lagrede søk ... | Lagre søket som

 felles

Done. Running time: 0.04 sec. (0.05 CPU sec.)

Treff 1 – 30 av 6693 | Forrige Neste | Gå til: | Last ned (Excel-modus) | Copy query URL | Vis linjefilter | Attributter ... | Linjer per side: |

Kontekststørrelse: 400px

corpus	treff	year
avis-nno en mellom Voss og Medkila i Prestegardsmoen i dag. Ballen	curla	2008
avis-nno azeem Ahmed det meste sjølv då han skar inn frå høgre og	curla	2011
avis-nno azeem Ahmed det meste sjølv då han skar inn frå høgre og	curla	2011
avis-nno azeem Ahmed det meste sjølv då han skar inn frå høgre og	curla	2011
avis-nno azeem Ahmed det meste sjølv då han skar inn frå høgre og	curla	2011
avis-nno I sjå ut. TV2 har kjøpt OL-rettar. Altså risikerer vi i framtida	curling	2011
avis-nno istoria bak fotball, tennis, basketball, bordtennis, volleyball,	curling	2015
avis-nno bildeCredit> FOTO: ARNE S. GJONE ¶ Les også En skole for	curlingbarna	2005
avis-nno lre tiåringane har det, går det ut på eitt. Nokre har kalla dei	curlingbarna	2006
avis-nno rl- Eirik Kval. Unge menneske i dag er ikkje kalla glasur- og	curlinggenerasjon	2015
avis-nno får dei unge mange kallenamn, som «lydig», «flinkis» og «	curlinggenerasjonen	2014
avis-nno 274481.ece> < F ¶ ~ Norsk seier over Japan ¶ De norske	curlinggutta	2006
avis-nno lekklasse, og at det på havets botn ligg ein tusen år gammal	curlingstein	2015
avis-plain ui " på Hawaii, Paramount går snart i gang med " Girl in the	curl	2001
avis-plain ui " på Hawaii, Paramount går snart i gang med " Girl in the	curl	2001
avis-plain og mars: ¶ Øvelse Sett repetisjoner ¶ Knebøy: 3 x 15 ¶ Leg	curl	2006
avis-plain ai og juni ¶ Øvelse Sett repetisjoner ¶ Knebøy: 3 x10 ¶ Leg	curl	2006
avis-plain imith 3-4 x 15 (på hver fot) ¶ Leg extension 3-4 x 15 ¶ Leg	curl	2006
avis-plain Smith 3-4 x10 (på hver fot) ¶ Leg extension 3-4 x 10 ¶ Leg	curl	2006
avis-plain Nordby synes OL-isen er rask. - Det er ikke så mye skru ("	curl	2006
avis-plain j forbanna på isen og forholdene. Vi er et lag som liker mye	curl	2006



Lemmutval (forts.): ordlister

Documentation
FAQ
Links

Korpusliste

Søk
Konkordans
Kollokasjoner
Distribusjon
Ordlister

Diagram
Tekst
Metadata
Oversikt
Variabler
Korpusdok.

Ordbøker
Ordbank

"curl.*"

Søk

Refine

window:

5 tokens

Stop

Lagrede søk ...

Done. Running time: 0.04 sec. (0.05 CPU sec.)

Antall treff: 6693, unique values: 429. Attributt: word | ignorer storskriving | sorter: etter frekvens relative to: -

Side 1 av 2. Previous Next

1514 (22,62%)	curling	21 (0,31%)	curlingpresident	7 (0,10%)	curling-	4 (0,06%)	curling-sporten
913 (13,64%)	curling-VM	20 (0,30%)	curlinglandslag	7 (0,10%)	curling-gull	4 (0,06%)	curlingeksperten
439 (6,56%)	curlinggutta	20 (0,30%)	curlingspilleren	7 (0,10%)	curlingbarn	4 (0,06%)	curlingfest
301 (4,50%)	curlet	19 (0,28%)	curlingguttene	7 (0,10%)	curlingbarna	4 (0,06%)	curlingforbund
204 (3,05%)	curling-EM	18 (0,27%)	curlingdamene	7 (0,10%)	curlingdamer	4 (0,06%)	curlinggenerasjon
159 (2,38%)	curlingjentene	18 (0,27%)	curlingsteiner	7 (0,10%)	curlingforelder	4 (0,06%)	curlinginnsats
150 (2,24%)	curlinglaget	16 (0,24%)	curlingforbundets	7 (0,10%)	curlingkarrieren	4 (0,06%)	curlingklubben
129 (1,93%)	curlinglandslaget	16 (0,24%)	curlingklovnene	7 (0,10%)	curlingkolleger	4 (0,06%)	curlingkunster
121 (1,81%)	curlinghallen	15 (0,22%)	curling-gutta	7 (0,10%)	curlinglagene	4 (0,06%)	curlinglegenden
116 (1,73%)	curlingherrene	15 (0,22%)	curlingklubb	7 (0,10%)	curlingnasjon	4 (0,06%)	curlingmamma
113 (1,69%)	curlinglag	14 (0,21%)	curlingbuksene	6 (0,09%)	curling-finalen	4 (0,06%)	curlingmesterskap
112 (1,67%)	curler	14 (0,21%)	curlinggull	6 (0,09%)	curlingdronning	4 (0,06%)	curlingnasjonen
111 (1,66%)	curler	14 (0,21%)	curlingkvinner	6 (0,09%)	curlingfans	4 (0,06%)	curlingproff
99 (1,48%)	curlingmillioner	14 (0,21%)	curlingmiljøet	6 (0,09%)	curlingfeberen	4 (0,06%)	curlingsenter
71 (1,06%)	curlingforbundet	14 (0,21%)	curlingseier	6 (0,09%)	curlingfolket	4 (0,06%)	curlingsjef
65 (0,97%)	curlingbanen	13 (0,19%)	curlere	6 (0,09%)	curlinghistorie	4 (0,06%)	curlingskipen
62 (0,93%)	curlingguttas	12 (0,18%)	curla	6 (0,09%)	curlinginteressen	4 (0,06%)	curlingsport
60 (0,90%)	curlingforeldre	12 (0,18%)	curlingbukser	6 (0,09%)	curlingjenter	4 (0,06%)	curlingsportens
54 (0,81%)	curling-seirer	12 (0,18%)	curlingfeber	6 (0,09%)	curlingkampene	4 (0,06%)	curlingsten
54 (0,81%)	curlingherrer	12 (0,18%)	curlingforeldrene	6 (0,09%)	curlingkonkurranse	3 (0,04%)	curlers
53 (0,79%)	curlinghall	12 (0,18%)	curlinglandslaget	6 (0,09%)	curlingleksjon	3 (0,04%)	curling-Norge
47 (0,70%)	curlingbane	12 (0,18%)	curlingmenn	6 (0,09%)	curls	3 (0,04%)	curlinga-gale



Kollokasjonar: bruksdøme, underoppslag

Documentation
FAQ
Links

Korpusliste

Søk
Konkordans

Kollokasjoner

Distribusjon
Ordliste
Diagram
Tekst
Metadata
Oversikt
Variabler
Korpusdok.

Ordbøker
Ordbank

[lemma="hatt" & morph = ("subst")]

Søk Refine window: 5 tokens

Lagre søket som

Done. Running time: 32.39 sec. (:

Vis kollokasjoner by word, venstrekontekst: 2, høyrekontekst: 2, combine context
MI * log(Freq) | Last ned

11708 collocations calculated; page 1 of 235. Forrige Neste | Show concordance for selection | show
freq. | show MI | show LL

Freq.	Delta	Collocate
74	-2	<input type="checkbox"/> fjør _ hatten
1	2	<input type="checkbox"/> << hatt >> _ halitt
1	-2	<input type="checkbox"/> deflasjonskrise _ Hattene
2	-2	<input type="checkbox"/> høysåter _ Hattane
1	-1	<input type="checkbox"/> SPARKA HATTEN
1	-1	<input type="checkbox"/> DISSE HATTENE
48	-2	<input type="checkbox"/> fjær _ hatten
4	-1	<input type="checkbox"/> bredbremmete hatter
5	-1	<input type="checkbox"/> bredbremmede hatter
2	-1	<input type="checkbox"/> tresnutede hatter
2	-1	<input type="checkbox"/> vidbremmede hatter
1	-2	<input type="checkbox"/> HØY _ HATTEN
1	-1	<input type="checkbox"/> trisnuta hattane
1	2	<input type="checkbox"/> hattane _ utfrunsa
1	-2	<input type="checkbox"/> stiletthæl _ Hattane
7	-1	<input type="checkbox"/> bredbremmete hatten
2	-1	<input type="checkbox"/> SIN HATT



Frekvens fortel ikkje alt

- Ein høg frekvens er ikkje åleine nok:
 - Korleis er ordet distribuert over tid?
 - Korleis er ordet distribuert på tvers av domene?
- Døme: søk på "aor.*" i Aviskorpuset (bokmål)





Korpuskel :: Aviskorpus (Bokmål) :: Ordliste

- Hjem
- Komme i gang
- Dokumentasjon
- FAQ
- Publikasjoner
- Lenker

Korpusliste

- Søk
- Konkordans
- Kollokasjoner
- Distribusjon

Ordliste

- Tekst
- Metadata
- Oversikt
- Variabler
- Korpusdok.

Lokalisering

Avansert søk | bytt til [Enkelt søk](#)

[Query history ...](#)

"aor.*"

Søk Refine window: 5 tokens | Stop

som felles

Done. Running time: 0.03 sec. (0.03 CPU s)

Antall treff: 113, unique values: 21. Attributt: word | ignorer storskriving | sorter: etter frekvens relative to:

[Last ned](#) include counts include fractions

Side 1 av 1.

58 (51,33%)	aortastenose
22 (19,47%)	aorta
8 (7,08%)	aortaklaffen
4 (3,54%)	aorta-ventiler
3 (2,65%)	aorta-disseksjon
2 (1,77%)	aortaklaff
2 (1,77%)	aortic
1 (0,88%)	aorakelsteiner
1 (0,88%)	aorhhh
1 (0,88%)	aorin
1 (0,88%)	aorta-aneurisme
1 (0,88%)	aorta-blodåren
1 (0,88%)	aortaaneurisme
1 (0,88%)	aortabuen
1 (0,88%)	aortaklaffene
1 (0,88%)	aortaklaffeprotese
1 (0,88%)	aortaklaffer



Dokumentasjon

FAQ

Publikasjoner

Lenker

Korpusliste

Søk

Konkordans

Kollokasjoner

Distribusjon

Ordlister

Tekst

Metadata

Oversikt

Variabler

Vis distribusjon type: absolutt | counts only | include structures

av word | ignorer storskriving, Δ : 0
relative to year | ignorer storskriving, Δ : 0
grupper etter - | ignorer storskriving, Δ : 0
og - | ignorer storskriving, Δ : 0

Fractions sum up to 1.0 in each row. Fractions in blue are unweighted means of group fractions.

Side 1/1 av 1x1. | Last ned

(sum)

aortastenose

(sum)	58 (100,0)	58 (100,0)
	58 (100,0)	58 (100,0)
2005	57 (100,0)	57 (100,0)
2014	1 (100,0)	1 (100,0)



Corpuscle-Lex :: Aviskorpus (Nynorsk), Aviskorpus (Bokmål), Forskning.no (2017), Leksikalsk bokmålskorpus, NBs frie tekster (Nynorsk), NBs frie tekster (Bokmål), Nynorsk-korpus, Talk Of Norway :: Distribusjon

Avansert søk | bytt til [Enkelt søk](#)

[Query history ...](#)

"assistenttren[a|e]r"

Søk Refine window: 5 tokens | Stop | [Lagrede søk](#)

... | [Lagre søket](#) som felles

Done. Running time: 0.21 sec. (0.34 CPU sec.)

Vis distribusjon type: | counts only | include structures

av ignorer storskriving, Δ : filter:

relative to ignorer storskriving, Δ : filter:

grupper etter ignorer storskriving, Δ : filter:

og ignorer storskriving, Δ : filter:

Fractions sum up to 1.0 in each row. Fractions in blue are unweighted means of group fractions. Fractions in green are distributions of total numbers.

Side 1/1 av 1x1. | [Last ned](#)

	(sum)	<input checked="" type="checkbox"/> assistenttrenar	<input checked="" type="checkbox"/> assistenttrener
(sum)	13264 (100,0)	163 (1,2)	13101 (98,8)
	13264 (100,0)	163 (31,3)	13101 (68,7)
avis-nno	29 (100,0)	29 (100,0)	
avis-plain	12977 (100,0)		12977 (100,0)
fn-new	2 (100,0)		2 (100,0)
lbk	18 (100,0)		18 (100,0)
nnk	238 (100,0)	134 (56,3)	104 (43,7)



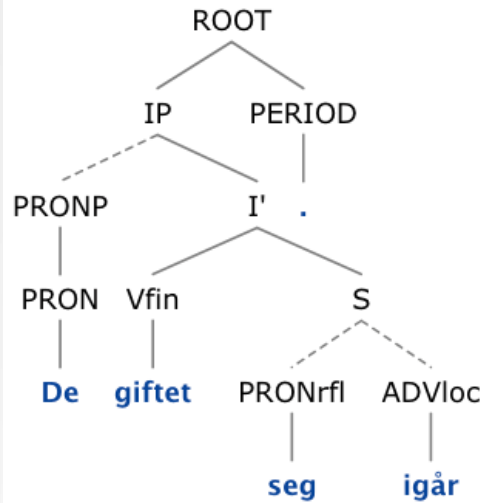
INESS og NorGramBank: søk og analyse (setningsnivå)

- **Trebank**: eit syntaktisk analysert tekstkorpus der kvar setning har ein detaljert syntaktisk analyse.
- **INESS**: ein infrastruktur for å bevare og gjere trebankar tilgjengelege (søk og analyse). Del av CLARINO Bergen Centre.
- **NorGramBank**: ein trebank for norsk, utvikla i prosjektet INESS (2010-2017).





C-structure



F-structure

PRED	'gifte*seg<[8:de]>[4:seg]'
TNS-ASP	10 TENSE past, MOOD indicative
TOPIC	PRED 'de' NTYPE 9 NSYN pronoun REF +, PRON-TYPE pers, PRON-FORM de, PERS 3, NUM pl, DEF +, CASE nom 8
ADJUNCT	1 { 2 PRED 'i går' ADV-TYPE temp }
OBJ	PRED 'seg' NTYPE 7 NSYN pronoun REF -, PRON-TYPE refl, PERS 3, NUM pl, CASE obl 4
SUBJ	[8]
0	VTYP main, VFORM fin, STMT-TYPE decl





Page 1 of 2 | Go to page: | public private

Filter by name: | by description: | by template:

Click on a name to choose a template. (31 stored templates shown.) | Show template expansion

[Store marked templates as sketch ...](#)

	<i>Name</i>	<i>Description</i>
<input type="checkbox"/>	* ADJ-coord(@ADJ)	Adjectives coordinated with an adjective
<input type="checkbox"/>	* ADJ-degreadv(@ADJ)	Degree adverbs modifying an adjective
<input type="checkbox"/>	* ADJ-modifies(@ADJ)	Nouns modified by an adjective
<input type="checkbox"/>	* ADJ-modnominadj(@ADJ)	Adjectives modifying a nominal head adjective
<input type="checkbox"/>	* ADJ-suff(@SUFF)	Adjectives derived with a suffix
<input type="checkbox"/>	* ADV-degmodifies(@ADV)	Adjectives modified by a degree adverb
<input type="checkbox"/>	* ADV-types(@ADV)	The types of an adverb
<input type="checkbox"/>	* N-adjmod(@N)	The adjectives modifying a noun
<input type="checkbox"/>	* N-defmascorfem(@N)	Feminine vs. masculine inflection of a noun





InESS

Sketch

Main Page
Knowledge center
The project
Documentation
FAQ
Publications
Links
Resources

Treebanks

Treebank
Selection
Treebank Details
Sentence
Overview
Sentence
Query
Sketch

Search in: nob-avis, nob-child, nob-fn, nob-lbk-av, nob-lbk-sa, nob-lbk-tv, nob-ndt-lfg, nob-novel, nob-novel_1, nob-novel_2, nob-novel_3, nob-novel_4, nob-novel_5, nob-novel_6, nob-novel_7, nob-sofie

max #: | fragments: none only | fully disamb.: none only

disambiguated: none only | unambiguous: none only

[Select query templates ...](#) | [Select sketch ...](#) | [Query history ...](#)

Template: * V-argframes(@V)

Description: Argument frames of a verb

Parameters:

@V:

Run query



Click on a row to see the matching sentences. | Copy format: plain BRO

UNIVERSITET I BERGEN



<i>Count</i>	<i>#a: atom</i>	<i>#arg1: value</i>	<i>#arg2: value</i>	<i>#arg3: value</i>
1545	gifte*seg	pronoun		
866	gifte*seg*med	pronoun	pronoun	
409	gifte*seg*med	pronoun	common	
346	gifte*seg*med	pronoun	proper	
201	gifte*seg	common		
128	gifte			
119	gifte*seg	proper		
96	gifte*seg*med	common	common	
63	gifte*seg*med	common	pronoun	
59	gifte*seg*med	proper	common	
54	gifte*seg*med	proper	proper	
52	gifte*seg*med	proper	pronoun	
26	gifte*seg*med	common	proper	
25	gifte*bort		pronoun	
20	gifte*bort	pronoun	pronoun	
18	gifte*bort		common	
15	gifte*seg*til	pronoun	common	
10	gifte*bort	pronoun	common	
8	gifte*bort	common	pronoun	
5	gifte*bort	common	common	
5	gifte*seg*til	pronoun	pronoun	
4	gifte*bort		proper	
3	gifte*seg*til	common	common	
2	gifte*bort	proper	pronoun	
1	gifte*bort	pronoun	proper	
1	gifte*seg*til	common	proper	



Count	#a: <i>atom</i>	#arg1: <i>value</i>	#arg2: <i>value</i>	#arg3: <i>value</i>
1545	gifte*seg	pronoun		
866	gifte*seg*med	pronoun	pronoun	
409	gifte*seg*med	pronoun	common	
346	gifte*seg*med	pronoun	proper	
201	gifte*seg	common		

UNIVERSITETET I BERGEN

Page 1 of 11 | Go to page: | [Download](#)

Click on a row to go to the sentence. Mouse over a row to see the structures.

<i>Treebank</i>	<i>Document</i>	<i>Trans.</i>	<i>Id</i>	<i>Sentence</i>	
nob-novel_7	oai:bibsys.no:biblio...	no	1739	En mann som vil gifte seg, er også i stand til slike handlinger, akkurat som sine forfedre.	<input type="button" value="Copy"/>
nob-novel_7	oai:bibsys.no:biblio...	no	2135	- Saken er den at mor har tenkt å gifte seg!	<input type="button" value="Copy"/>
nob-novel_7	oai:bibsys.no:biblio...	no	2194	Hun er en kvinne som gifter seg, akkurat som andre kvinner gifter seg hver dag, hver time på dagen!	<input type="button" value="Copy"/>
nob-novel_7	oai:bibsys.no:biblio...	no	2416	Det var ikke noe galt i at en «kvinne» giftet seg igjen etter skilsmissen.	<input type="button" value="Copy"/>
nob-novel_7	oai:bibsys.no:biblio...	no	2620	- Far giftet seg da han var på min alder.	<input type="button" value="Copy"/>
nob-novel_7	oai:bibsys.no:biblio...	no	5153	Hvis du vil, skal jeg gi deg flere titalls eksempler på yngre søstrer som har giftet seg før eldre, og det har ikke vært til hinder for at de eldste giftet seg med de beste ektemennene som tenkes kan.	<input type="button" value="Copy"/>
nob-novel_7	oai:bibsys.no:biblio...	no	5892	Tanken på at datteren skulle gifte seg, gav ham en merkelig, ubehagelig følelse, enda det stred mot både fornuft og moral.	<input type="button" value="Copy"/>



Oppsummering

Korpus er til god hjelp for leksikografen når det gjeld:

- lemmatilfang (ord inn og ord ut av ordboka)
 - underoppslag
 - ordtydingar
 - syntaktisk åtferd
 - bruksdøme
-
- ... Men korpus må brukast klokt, og kan ikkje erstatte leksikografiske vurderingar
 - Behov for vidareutvikling av tekniske løysingar







UNIVERSITETET I BERGEN

