



UNIVERSITETET I OSLO  
DET HUMANISTISKE FAKULTET

# Spoken language corpora

Janne Bondi Johannessen

AAU, March 2014

# Spoken corpora at the Text Laboratory, UiO



UNIVERSITETET I OSLO  
DET HUMANISTISKE FAKULTET

- Big Brother corpus
  - 2001. ca. 550 000 words
- Nordic dialect corpus
  - 2012. Ca. 2 mill. words
- NoTa-Oslo
  - 2006. Ca. 900 000 words
- Ruija corpus
  - 2006. Ca. 400 000 words.
- TAUS
  - 1975 / 2006 Old Oslo. Ca. 250 000 words.
- Doctor patient corpus
  - 2014





# Spoken language corpora



- Purposes of a spoken language corpus
  - Study aspects of language as with written corpora
  - In addition, studies of:
    - Discourse
    - Phonology
    - Dialect
    - Diachronic change
    - Agelect
    - ...



# How make spoken data into a corpus?

- Digitise sound/video
- Transcribe speech
  - Choose transcription
    - Phonetic
    - Orthographic
- Time code transcription for alignment with sound file
- Annotate
  - Extra-linguistic features (laughter, coughing)
  - Semi-linguistic features (meaningful sounds)
  - Gestures
  - Grammatically
    - POS tagging
    - Syntactic parsing
    - Other



# Different situation and informants - different speech types



- The Big Brother Corpus
- Pros
  - Lots of spontaneous speech data
  - Lots of dialogue and polylogue
  - Lots of emotional speech in different dialogue situations
    - Conflict, argument, love, irritation etc.
- Cons
  - Not a dialect corpus
  - No representativity w.r.t. age, social class, education etc.
  - Not "controlled" recording situations
  - Small number of informants

- Nordic Dialect Corpus
- Pros
  - Lots of spontaneous speech data
  - Lots of dialogue, **no** polylogue
  - **No** emotional speech
- Cons
  - Is a dialect corpus
  - Partly representative
  - Controlled recording situation
  - Small AND large number of informants





# Big Brother vs Nordic Dialect Corpus



[Trouble viewing video?](#)  
**context±**  
Offset  
Left   
Right   
- Start +  
- Stop +  
>>



[Trouble viewing video?](#)  
**context±**  
Offset  
Left   
Right   
- Start +  
- Stop +  
>>



## Other dialect corpora?

### We know of no comparable resource for any language

- *Sounds familiar? Accents and Dialects of the UK*
  - No grammatical search options
  - No results handling
- *The British National Corpus*
  - No audio
  - Orthographic transcription
  - Unreliable dialect categories
- *The DynaSand dialect database*
  - Few spontaneous utterances
- *The Spoken Dutch Corpus*
  - Not web-based
  - Orthographic transcriptions
  - Not dialect data
- The Scottish Corpus of Text and Speech
  - Not a dialect corpus
  - No searchable linguistic features
- Others under development:
  - Corpus of Estonian Dialects
  - Spoken Japanese Dialect Corpus
- Paul Thompson at the University of Reading: Posting at Corpora List 30 Nov. 2008 about linked audio or video files with transcripts: 15 answers, of which only one on dialects: ours



## Comparison: NoTa-Oslo, Talesøk, BySoc, BNC, SCOTS

	NoTa	Talesøk	GSLC	BySoc	BNC	SCOTS
Transcription linked to audio	Yes	Yes	No	No	No	Yes
Transcription linked to video	Yes	No	No	No	No	Yes
User-friendly search without regular expressions	Yes	Yes	No	No	No	No
Possible to limit informant selection	Yes	Yes	Yes	Yes	Yes	Yes
Overlaps/ turntaking annotated	Yes	Yes	Yes	Yes	No	Yes
Transcription as standard orthography (or slightly modified)	Yes	Yes	Yes	Yes	Yes	Yes
POS tagged	Yes	No	No	No	Yes	No
POS tags can be used as the only search expressions	Yes	–	–	–	No	–



# Credits

- Collaboration between the research network ScanDiaSyn and Nordic Centre of Excellence in Microcomparative Syntax (NORMS)
- Some of the corpus is financed by national research councils in the individual countries
- The technical development has been financed by the University of Oslo, the Norwegian Research Council, and the Nordic research funds NOS-HS and NordForsk



## Why the Nordic Dialect Corpus was developed

- Initiated by members of Nordic Centre of Excellence in Microcomparative Syntax and the ScanDiaSyn network
- Overarching goal: to study the dialects of the North-Germanic dialect continuum
  - The Nordic languages are closely related and have some mutually intelligibility
  - Studying the dialects within each national language is misguided from a theoretical and principled point of view
  - Difficult for each researcher to get hold of relevant data on their own in such a large area.
  - Many different kinds of data needed for syntax research



# Corpus features

- Linguistic contents
  - Dialects from five closely related languages
- Annotation
  - POS tagging and two types of transcription
- Search interface
  - Advanced possibility to combine an array of search criteria and results presentation in an intuitive interface
- Many search variables
  - Linguistics-based, informant-based, time-based
- Multimedia display
  - Linking of sound and video to transcriptions
- Display of informant details
  - Number of words and other informant information
- Advanced results handling
  - Concordances, collocations, counts, statistics ...
- The corpus is available on the web





# Linguistic contents and numbers

- The corpus contains dialect data from the national languages Danish, Faroese, Icelandic, Norwegian, and Swedish
- Contains speech data from app. 821 informants with 2.8 mill. words
- All the recordings represent spontaneous speech
- Differences in data collection due to differences in financing
  - Norwegian, Oevdalian, and some Danish: two kinds of recordings per informant:
    - semi-formal interview
    - informal conversation between two informants
  - Recordings of both young and old informants, both genders
  - Both new and old recordings
  - Audio or both audio and video recordings



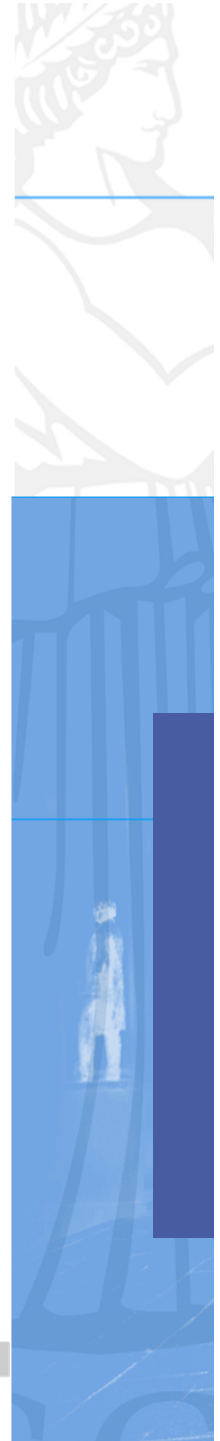




# Annotation: transcription

- Each dialect has been transcribed by the standard official orthography of that country
  - In addition all the Norwegian dialects and some Swedish dialects have also been transcribed phonetically
    - The Norwegian phonetic transcription follows roughly that of Papazian and Helleland (2005). The transcription of the Oevdalian dialect follows the Oevdalian orthography (standardised in 2005 by the *Råðdjärum* – The Oevdalian Language Council).
  - The phonetic transcription is translated to an orthographic transcription via a semi-automatic dialect transliterator





# Annotation: tagging

- The corpus is POS tagged, with selected morpho-syntactic features language by language

# Search interface – Glossa

Scandinavian Dialect  
Corpus

Glossa ( [my results](#) | [my annotations](#) | [statistics](#) | [full query](#) | [help](#) )

criteria»

+  
-



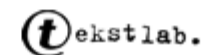
[Recording locations](#)  
[Transcriptions](#)

<b>Regular expressions:</b> <input type="checkbox"/>	<b>Hits per page:</b> <input type="text" value="20"/>	Randomize <input type="checkbox"/>	Orthographic <input checked="" type="radio"/>	<input type="button" value="Search corpus"/>
<b>Search within:</b> <input type="text" value="s"/>	<b>Max results :</b> <input type="text" value="2000"/>	Skip tot. freq. <input checked="" type="checkbox"/>	Phonetic <input type="radio"/>	<input type="button" value="Reset form"/>
<b>informant</b> <sup>+</sup>				<input type="button" value="Show texts"/>
<hr/>				<input type="button" value="Save subcorpus"/>
<b>country</b> <sup>+</sup>	<b>region</b> <sup>+</sup>	<b>area</b> <sup>+</sup>	<b>place</b> <sup>+</sup>	<a href="#">Choose subcorpus</a>
<hr/>				
<b>agegroup</b> <sup>+</sup>	<b>sex</b> <sup>+</sup>	<b>rec (year)</b> <sup>+</sup>	<b>genre</b> <sup>+</sup>	
<hr/>				
<b>Display:</b> <input type="text"/>	<b>Search within:</b> <input type="text"/>			



Denne Glossaversjonen er under utvikling.  
This version of Glossa is undergoing development.

[Please report bugs and errors here](#)



# Searching for lemmas



gås +

criteria» -

- word » lemma
- occurrences » start of word
- pos » within word
- num » end of word
- degr » case sensitive
- case » exclude word
- sex » add word form
- nlex » add negated word form
- pers » add lemma
- temp » add negated lemma
- def »
- descr »
- type »
- phonetic »

Re Se

per page: 20

Randomize

Orthographic

Phonetic

results : 2000

Skip tot. freq.

inf

country + region + area + place +

agegroup + sex + rec (year) + genre +

Display:  Search within:

[Recording locations](#)  
[Transcriptions](#)

Search corpus

Reset form

Show texts

Save subcorpus

[Choose subcorpus](#)



Denne Glossaversjonen er under utvikling.  
This version of Glossa is undergoing development.

[Please report bugs and errors here](#)

# Searching for more than one word

ig	interval:		+
criteria»	min	criteria»	-
end of word	3 max	konj	



[Recording locations](#)  
[Transcriptions](#)

Regular expressions: <input type="checkbox"/>	Hits per page: <input type="text" value="20"/>	Randomize <input type="checkbox"/>	Orthographic <input checked="" type="radio"/>
Search within: <input type="text" value="s"/>	Max results: <input type="text" value="2000"/>	Skip tot. freq. <input checked="" type="checkbox"/>	Phonetic <input type="radio"/>

\_\_\_\_\_

\_\_\_\_\_

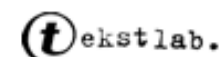
[Choose subcorpus](#)

Display:  Search within:



Denne Glossaversjonen er under utvikling.  
This version of Glossa is undergoing development.

Please report bugs and errors [here](#)



# Search results

CWB expression: "([((word=".\*ig" %c)) [0,3] [((pos="konj"))]) ;"

Informants: 240

scandiasyn:

CWB expression: "([((word=".\*ig" %c)) [0,3] [((pos="konj"))]) ;"

Action :

: 474

Results pages: [1](#) [2](#) [3](#) [4](#) [5](#) [6](#) [7](#) [8](#) [9](#) [10](#) [11](#) [12](#) [13](#) [14](#) [15](#) [16](#) [17](#) [18](#) [19](#) [20](#) [21](#) [22](#) [23](#) [24](#)

- [i](#) [M](#) [K](#) [alvdal\\_02](#) påbygning # eller det er ikke så **veldig herlig det men**
- [i](#) [M](#) [K](#) [alvdal\\_01](#) det er godt å bli **ferdig snart da ja men** (uforståelig) jævlig mye å gjøre i år
- [i](#) [M](#) [K](#) [alvdal\\_01](#) ikke sett n # så nå har jeg jo **plutselig fire sjøl da for** jeg hadde jo tre # tre sjø
- [i](#) [M](#) [K](#) [alvdal\\_02](#) tjener ikke så **jævlig mye men** # det blir jo litt
- [i](#) [M](#) [K](#) [alvdal\\_01](#) ja det s- ja synes jeg ser så **herlig ut \* men** jeg ...
- [i](#) [M](#) [K](#) [alvdal\\_02](#) \* m \* yes \* (latter) \* det ser **herlig ut ja men** det er vel ett nålhue å komme gj

## Some results presented as frequency list

<b>occurences</b>	<b>match</b>
8	artig og
6	rolig og
6	egentlig men
3	mulig og
3	vanskelig og
3	egentlig # men
3	viktig og
3	mulig men
2	tidlig for
2	veldig kjekt og
2	frodig og
2	veldig fin oppvekst og
2	tidlig e begynne og
2	vanlig skole og
2	artig det og
2	veldig men
2	egentlig og
2	egentlig heile # oppveksten og
2	forferdelig r�h�lja heter det og
2	veldig gode venner i_hvert_fall og
2	veldig fin natur og
2	gr�dig n� et men

# Searching for part of speech



criteria»

word »

occurrences »

pos »

num »

degr »

case »

sex »

nlex »

pers »

temp »

def »

descr »

type »

phonetic

(spm)

adj

adv

cbl

det

inf-merke

interj

konj

pause

pause2

prep

pron

pron/det

sbu

subst

subst:adj

sann

ukjent

verb

verb:subst

ambiguous

Recording locations  
Transcriptions

Hits per page:  Randomize

Max results:  Skip tot. freq.

Orthographic

Phonetic

Search corpus

Reset form

Show texts

Save subcorpus

Choose subcorpus

country +

agegroup +

sex +

region +

area +

place +

genre +

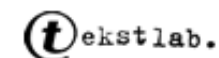
Display:

Search within:



Denne Glossaversjonen er under utvikling.  
This version of Glossa is undergoing development.

Please report bugs and errors [here](#)





# Phonetic querying

æøå...» æøå...»

jeg interval: er

min

criteria» max criteria»

+  
-



- [Transcription guidelines, translation lists,](#)
- [Recording locations](#)
- [Transcriptions](#)
- [News: Search Interface Documentation](#)

[add phrase](#) [delete phrase](#)

**Regular expressions:**  **Hits per page:**  **Randomize**  **Orthographic**   
**Max results:**  **Skip tot. freq.**  **Phonetic**   
**Both**

[Search corpus](#)  
[Reset form](#)

**informant** <sup>+</sup>

---

**country** <sup>+</sup> **region** <sup>+</sup> **area** <sup>+</sup> **place** <sup>+</sup>

---

**agegroup** <sup>+</sup> **sex** <sup>+</sup> **rec (year)** <sup>+</sup> **genre** <sup>+</sup>

[Show informants](#)  
[Save subcorpus](#)  
[Choose subcorpus](#)

# Displaying results

aasen\_35 ja var det inte det ?

aasen\_48 **ojoj och hur fin (uförståelig) tänk att i går # så pratade M1 och jag just om ... # för då hade vi bröllopsdag i går**

aasen\_35 jaha ?



[Trouble viewing video?](#)

context

± - ▾

- Start +

- Stop +

>>

Informants: 240

scandiasyn:

CWB expression: "**(((phon="um.\*" %c)))** ;"




Action :




: 242




Results pages: [1](#) [2](#) [3](#) [4](#) [5](#) [6](#) [7](#) [8](#) [9](#) [10](#) [11](#) [12](#) [13](#)

- i 🗨 🗨 aasen\_48 ojoj och hur fin (uförståelig) tänk att i går # så pratade M1 och jag just **om** ... # för då hade vi bröllopsdag i går
- i 🗨 ls men vad gjorde **om** om rasterna ?
- i 🗨 ls men vad gjorde om **om** rasterna ?
- i 🗨 🗨 aasen\_35 och eh # till slut var det någon som var in där och tit- skulle se efter **om** hon hade haft varit darinne och där stod en slik # stor # brödkartong
- i 🗨 ls var det göra det\_där byta **om** ? lära sig något nytt ?
- i 🗨 🗨 aasen\_35 det det var ju så det var ju många # stockholmare som var här **om** sommaren # över i gårdarna och hade ju lite # lärde ju oss lite av dem
- i 🗨 🗨 aasen\_35 ja det var verkligen riktig bocksvenska **om** de förstod det vet jag inte

- Phonetic and orthographic transcription

   aal\_01um nå skal vi nå er jeg er med i em (front-click) det derre prosjektet som H  
no ska me no e e e me i em \_ de dære prosjekkte somm Hallingdal Gå  
Now we are going now I'm in em (front-click) the derre project Hallingdal G

   aal\_01um jeg er så glad i Ål jeg  
e e så gla i ÅL e  
I'm so happy I Eel (google)

   aal\_01um jeg er så en altså d- # jeg budde jo i Oslo ett år i\_fjor  
e e så n asså d # e budd jo i Ossjlo æitt år i\_fjoL  
I saw an ie d ... I lived in the Oslo one year i\_fjor (google)



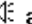


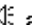


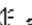


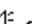


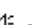


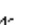





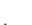



























# Display of transcription and tagging

CWB expression: "(((word="bil.\*" %c))) ;"

Action:

Hits found: 124

Results pages: [1](#) [2](#) [3](#) [4](#) [5](#) [6](#) [7](#)

- 1    **alvdal\_02** så skal jeg ta førerkort på bil nå # når jeg bli atten # så er jo blir litt godt å ## endelig få kjøre **bil**
- 1    **alvdal\_02** ja jeg # kjøpte meg **bil** i sommer # som jeg har pussa opp att
- 1    **alvdal\_02** har du traktor så # eller **bil** så kjører du heller det enn å sparke # går litt fortere
- 1    **alvdal\_04** nei da så er da mye men nå er det vel så du kommer ikke innover med **bil** heller nå vil jeg tru
- 1    **alvdal\_03** nye veien utover er så stille og så utoverbakke det er så at **bilene** bare trille utover jeg hører dem ikke i det hele tatt
- 1    **alvdal\_04** ja # ja når du skal ut med **bil** så er det vel det
- 1    **alvdal\_03** men sa at jeg tar da ikke **bil** for å reise på skitur # hvis jeg il
- 1    **alvdal\_03** og nå har jeg blitt for lat så jeg gidder ikke å gå opp her heller **sex:** mask kjøre på Kleberswanga det er helst der jeg kjører for jeg tør ikke kjøre opp
- 1    **alvdal\_03** og jeg tenkte at jeg kommer aldri hjem att med denne derre de **num:** fl **type:** appell er
- 1    **alvdal\_03** og e så var jeg en ta disse postfolkene så da # " har du fått pro **defn:** be **bil**en ? "
- 1    **alvdal\_03** " ja " sa " jeg har gjort det " men e ja vi har skulle ringe på n M **descr:** me **nlex:** 100 en det er ikke han som er # e der jeg går med **bil**en " sa jeg # men det v
- 1    **alvdal\_03** ja jeg hadde ikke lovt han noe men han ville gjerne ha **bil**en
- 1    **alvdal\_04** ja dem er aldeles sjuk på gamle **biler** den # er m hun ene der a F9 hun kjøpte denne gamle Ladaen som jeg hadde vet du jeg hadde hadde Lada jeg
- 1    **alvdal\_04** men så byn- ja ja en gikk den men hun var aldeles sjuk på n vet du denne **bil**en så hun # kjøpte nå for en billig penge der
- 1    **alvdal\_04** men så byn- ja ja en gikk den men hun var aldeles sjuk på n vet du denne bilen så hun # kjøpte nå for en **billig** penge der
- 1    **alvdal\_03** hva\_for en **bil** du har nå ?
- 1    **alvdal\_04** han driver på med **bilforretning** inni Oslo # så han skaffa meg denne der er en totusenmodell så den er det er en svær # firehjulstrekker og

# Informant-based querying

  
criteria»

[Recording locations](#)  
[Transcriptions](#)

Regular expressions:   
Search within:

Hits per page:   
Max results :

Randomize   
Skip tot. freq.

Orthographic   
Phonetic

Search corpus

Reset form

informant

- aasen\_35
- aasen\_48
- aeroe2
- aeroe3
- als1
- alvdal\_01
- alvdal\_02

choose

country

- Denmark
- Faroer
- Iceland
- Norway
- Sweden

choose

region

- Gotland
- Götaland
- Norrland
- Svealand

choose

area

- Blekinge
- Bohuslän
- Dalarna
- Dalsland
- Gotland
- Gästrikland
- Halland

choose

place

- als
- alvdal
- andøya
- Ankarsum
- Anundsjö
- Arjeplog
- Asby

choose

agegroup

A  
 B

sex

F  
 M

rec (year)

- 1959
- 2000
- 2007
- 2008

choose

genre

- ikke\_udskrevet
- intervju
- klip
- samtale
- 

choose

Show texts

Save subcorpus

[Choose subcorpus](#)

# Display information on informants 1

CWB expression: "`((word="bil.*" %c))` ;"

Action:

Hits found: 124

Results pages: [1](#) [2](#) [3](#) [4](#) [5](#) [6](#) [7](#)

[i](#) [H](#) [K](#) [alvdal\\_02](#) så skal je

[i](#) [H](#) [K](#) [alvdal\\_02](#) ja jeg # k

[i](#) [H](#) [K](#) [alvdal\\_02](#) har du tra

[i](#) [H](#) [K](#) [alvdal\\_04](#) nei da så

[i](#) [H](#) [K](#) [alvdal\\_03](#) nye veien

[i](#) [H](#) [K](#) [alvdal\\_04](#) ja # ja nå

[i](#) [H](#) [K](#) [alvdal\\_03](#) men sa at

[i](#) [H](#) [K](#) [alvdal\\_03](#) og nå har

[i](#) [H](#) [K](#) [alvdal\\_03](#) og jeg ter

[i](#) [H](#) [K](#) [alvdal\\_03](#) og e så v

[i](#) [H](#) [K](#) [alvdal\\_03](#) " ja " sa "

[i](#) [H](#) [K](#) [alvdal\\_03](#) ja jeg hadde ikke lovt han noe men han ville gjerne ha **bilen**

[i](#) [H](#) [K](#) [alvdal\\_04](#) ja dem er aldeles sjuk på gamle **biler** den # er m hun ene der a F9 hun kjøpte denne gamle Ladaen som jeg hadde vet du jeg

[i](#) [H](#) [K](#) [alvdal\\_04](#) men så byn- ja ja en gikk den men hun var aldeles sjuk på n vet du denne **bilen** så hun # kjøpte nå for en billig penge der

[i](#) [H](#) [K](#) [alvdal\\_04](#) men så byn- ja ja en gikk den men hun var aldeles sjuk på n vet du denne bilen så hun # kjøpte nå for en **billig** penge der

[i](#) [H](#) [K](#) [alvdal\\_03](#) hva\_for en **bil** du har nå ?

Informant details for *alvdal\_03* in the Scandiasyn corpus

```
SELECT tid,sex,agegroup,country,region,area,place,wc,rec FROM SCANDIASYNauthor WHERE tid = 'alvdal_03'
```

Code	Sex	Age group	Country	Region	Area	Place	Word count	Recorded
alvdal_03	F	B	Norway			alvdal	4876	2008

Fullført

ele tatt

t der jeg kjører fo

jeg går med **bile**

## Display information on informants 2

**SELECT SUM(wc) FROM SCANDIASYNauthor WHERE active REGEXP '^ (y)\$' AND active IS NOT NULL**

**Word count for selected informants: 2824585 (total for Scandiasyn corpus: 2824585)**

**Selected informants: 821 from 228 places in 5 countries**

Code	Sex	Age group	Country
ankarsrum_om1	M	B	Sweden
ankarsrum_om3	M	B	Sweden
ankarsrum_ym1	M	A	Sweden
anundsjo_om2	M	B	Sweden
anundsjo_ow3	F	B	Sweden
arjeplog_om1	M	B	Sweden
arsunda_ow1	F	B	Sweden
arsunda_om1	M	B	Sweden
arsunda_om2	M	B	Sweden
arsunda_ow2	F	B	Sweden
asby_om2	M	B	Sweden
asby_om3	M	B	Sweden

# Action menu

CWB expression: "(((word="bil.\*" %c))) ;"

Action :

Hits four

Results p



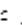
- count
- download
- sort
- collocations
- annotate
- show metadata
- metadata distribution
- delete hits
- save hits



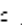
a førerkort på bil nå # når jeg bli atten # så er jo blir litt godt å ## endelig få kjøre **bil**




te meg **bil** i sommer # som jeg har pussa opp att




or så # eller **bil** så kjører du heller det enn å sparke # går litt fortere




da mye men nå er det vel så du kommer ikke innover med **bil** heller nå vil jeg tru




   **alvdal\_03** nye veien utover er så stille og så utoverbakke det er så at **bilene** bare trille utover jeg hører dem ikke i det h




   **alvdal\_04** ja # ja når du skal ut med **bil** så er det vel det




   **alvdal\_03** men sa at jeg tar da ikke **bil** for å reise på skitur # hvis jeg ikke skal inn i S1

   **alvdal\_03** og nå har jeg blitt for lat så jeg gidder ikke å gå opp her heller jeg nå jeg må kjøre på Klebersvanga det er helst

   **alvdal\_03** og jeg tenkte at jeg kommer aldri hjem att med denne derre derre **bilen** der

   **alvdal\_03** og e så var jeg en ta disse postfolkene så da # " har du fått problemer med **bilen** ? "

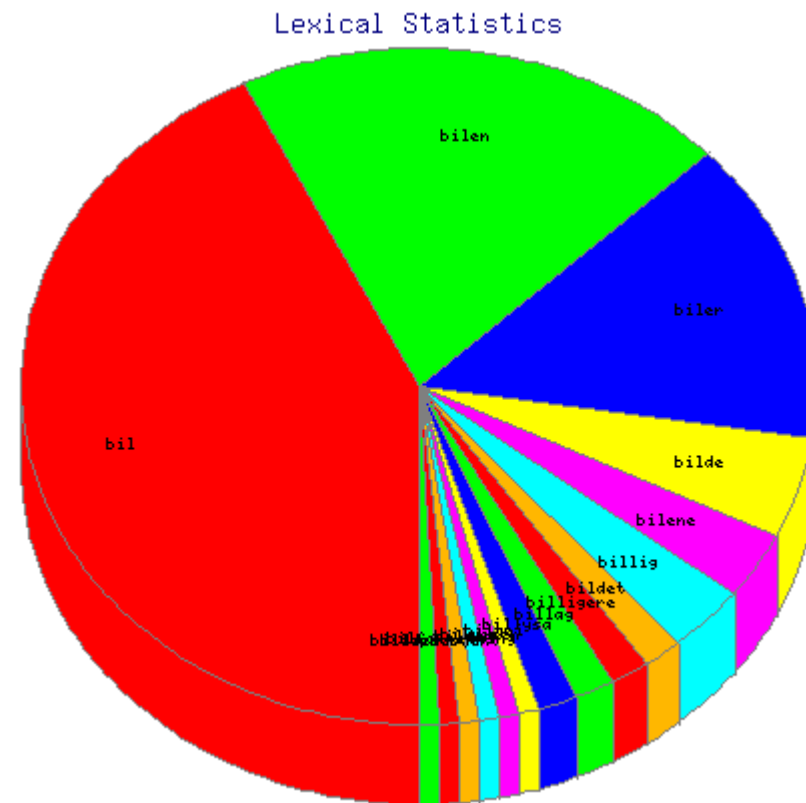
   **alvdal\_03** " ja " sa " jeg har gjort det " men e ja vi har skulle ringe på n M4 han # " ja men det er ikke han som er # e der je

   **alvdal\_03** ja jeg hadde ikke lovt han noe men han ville gjerne ha **bilen**



# Count

occurrences	match
53	bil
25	bilen
18	biler
6	bilde
4	bilene
4	billig
2	bildet
2	billigere
2	billag
2	billysa
1	bilvei
1	bilveien
1	bilveg
1	bilforretning
1	bildeprosjekt
1	bildeler



# Deleting or selecting some results

CWB expression: "(((word="bil.\*" %c))) ;"

Action :

Hits found: 124

Results pages: [1](#) [2](#) [3](#) [4](#) [5](#) [6](#) [7](#)

Delete selection

select all

unselect all

Finished deleting

- i** **alvdal\_02** så skal jeg ta førerkort på bil nå # når jeg bli atten # så er jo blir litt godt å ## endelig få kjøre **bil**
- i** **alvdal\_02** ja jeg # kjøpte meg **bil** i sommer # som jeg har pussa opp att
- i** **alvdal\_02** har du traktor så # eller **bil** så kjører du heller det enn å sparke # går litt fortere
- i** **alvdal\_04** nei da så er da mye men nå er det vel så du kommer ikke innover med **bil** heller nå vil jeg tru
- i** **alvdal\_03** nye veien utover er så stille og så utoverbakke det er så at **bilene** bare trille utover jeg hører dem ikke i det hele tatt
- i** **alvdal\_04** ja # ja når du skal ut med **bil** så er det vel det
- i** **alvdal\_03** men sa at jeg tar da ikke **bil** for å reise på skitur # hvis jeg ikke skal inn i S1
- i** **alvdal\_03** og nå har jeg blitt for lat så jeg gidder ikke å gå opp her heller jeg nå jeg må kjøre på Klebersvanga det er helst der jeg (uforståelig)
- i** **alvdal\_03** og jeg tenkte at jeg kommer aldri hjem att med denne derre derre **bilen** der
- i** **alvdal\_03** og e så var jeg en ta disse postfolkene så da # " har du fått problemer med **bilen** ? "

# Annotating results

CWB expression: "`((word="hun" %c))((pos="subst"))` ;"



Action :

Hits found: **96**

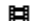

Results pages: [1](#) [2](#) [3](#) [4](#) [5](#)

Save annotations (\* indicates: no reviewed value stored.)



\*

**i**   **mefjordvaer\_19** ja men det hender jo at # at vi slår til og så bruker uttrykk og d- **hun datter** vår som bor i # Oslo hun em # hun jot


\*

**i**   **mefjordvaer\_19** " ja hvor fikk **hun tak** i de jernene ? " " jo det var smeden i Været " sa han # " som laga jernet til hun " sa han

\*

**i**   **stange\_03** og e # ja jeg vet ikke åssen jeg skal si det # nei jeg har akkurat e s- sagt åt **hun dama** de skrev at det jeg kan tenk

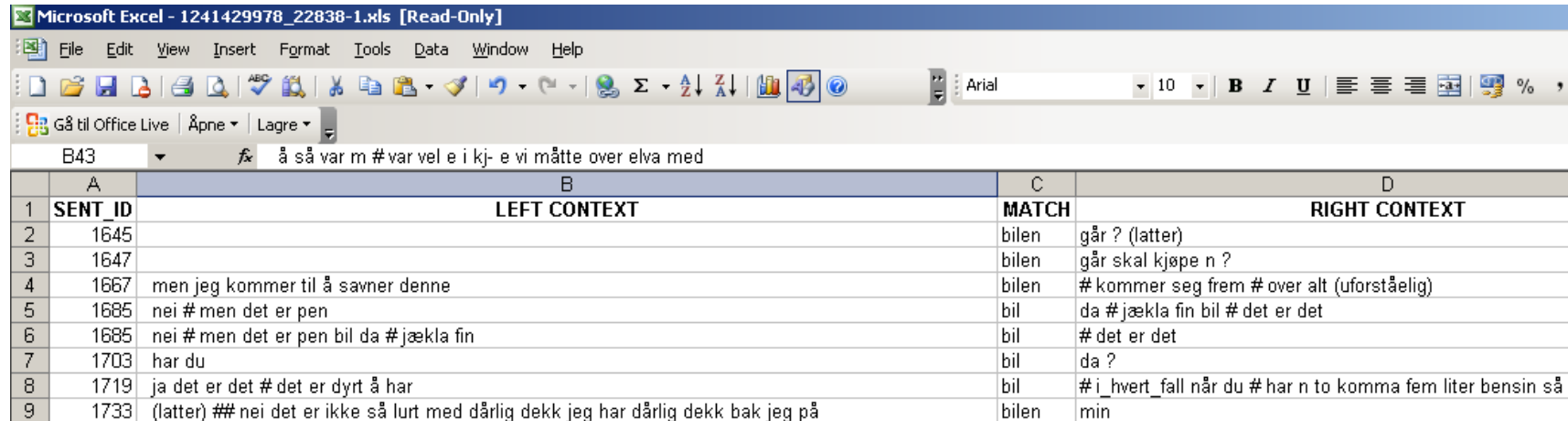
\*

**i**  **stonglandseide\_47** når du kommer dit som du skulle tømme henne så har jo **hun hålt** lekke tom på veien og sånn blir det med skatten veien

\*

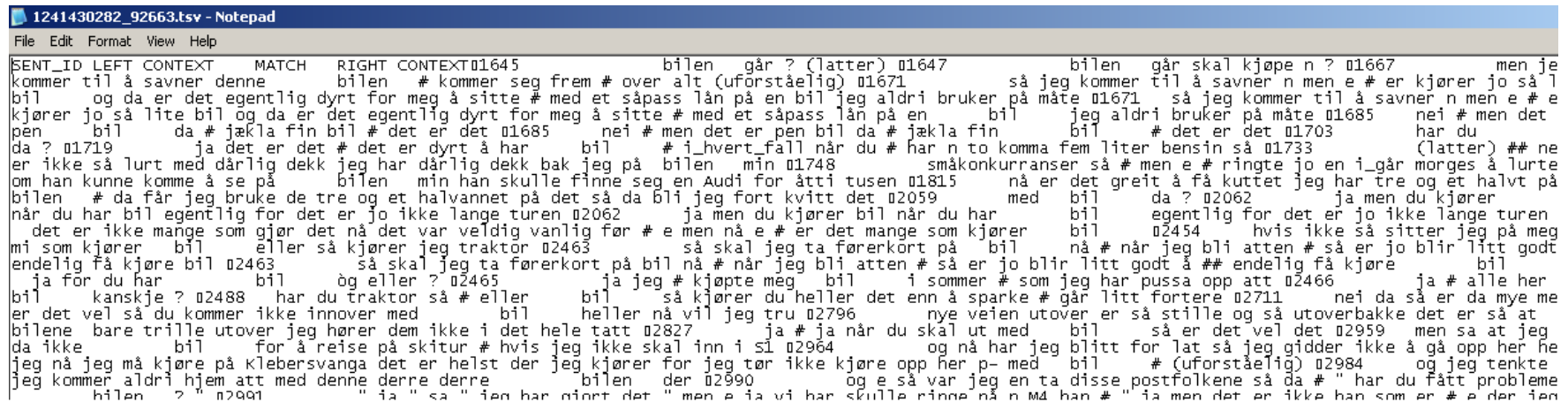
# Downloading results, examples:

## Excel:



	A	B	C	D
	SENT_ID	LEFT CONTEXT	MATCH	RIGHT CONTEXT
1	1645		bilen	går ? (latter)
2	1647		bilen	går skal kjøpe n ?
4	1667	men jeg kommer til å savner denne	bilen	# kommer seg frem # over alt (uforståelig)
5	1685	nei # men det er pen	bil	da # jækla fin bil # det er det
6	1685	nei # men det er pen bil da # jækla fin	bil	# det er det
7	1703	har du	bil	da ?
8	1719	ja det er det # det er dyrt å har	bil	# i hvert fall når du # har n to komma fem liter bensin så
9	1733	(latter) ## nei det er ikke så lurt med dårlig dekk jeg har dårlig dekk bak jeg på	bilen	min

## Tab separated values:



```
1241430282_92663.tsv - Notepad
File Edit Format View Help
SENT_ID LEFT CONTEXT MATCH RIGHT CONTEXT01645 bilen går ? (latter) 01647 bilen går skal kjøpe n ? 01667 men je
kommer til å savner denne bilen # kommer seg frem # over alt (uforståelig) 01671 så jeg kommer til å savner n men e # er kjører jo så l
bil og da er det egentlig dyrt for meg å sitte # med et såpass lån på en bil jeg aldri bruker på måte 01671 så jeg kommer til å savner n men e # e
kjører jo så lite bil og da er det egentlig dyrt for meg å sitte # med et såpass lån på en bil jeg aldri bruker på måte 01685 nei # men det
pen bil da # jækla fin bil # det er det 01685 nei # men det er pen bil da # jækla fin bil # det er det 01703 har du
da ? 01719 ja det er det # det er dyrt å har bil # i hvert fall når du # har n to komma fem liter bensin så 01733 (latter) ## ne
er ikke så lurt med dårlig dekk jeg har dårlig dekk bak jeg på bilen min 01748 småkonkurranser så # men e # ringte jo en i går morges å lurte
om han kunne komme å se på bilen min han skulle finne seg en Audi for åtti tusen 01815 nå er det greit å få kuttet jeg har tre og et halvt på
bilen # da får jeg bruke de tre og et halvannet på det så da bli jeg fort kvitt det 02059 med bil da ? 02062 ja men du kjører
når du har bil egentlig for det er jo ikke lange turen 02062 ja men du kjører bil når du har bil egentlig for det er jo ikke lange turen
det er ikke mange som gjør det nå det var veldig vanlig for # e men nå e # er det mange som kjører bil 02454 hvis ikke så sitter jeg på meg
mi som kjører bil eller så kjører jeg traktor 02463 så skal jeg ta førerkort på bil nå # når jeg bli atten # så er jo blir litt godt
endelig få kjøre bil 02463 så skal jeg ta førerkort på bil nå # når jeg bli atten # så er jo blir litt godt å ## endelig få kjøre bil
ja for du har bil og eller ? 02465 ja jeg # kjøpte meg bil i sommer # som jeg har pussa opp att 02466 ja # alle her
bil kanskje ? 02488 har du traktor så # eller bil så kjører du heller det enn å sparke # går litt fortere 02711 nei da så er da mye me
er det vel så du kommer ikke innover med bil heller nå vil jeg tru 02796 nye veien utover er så stille og så utoverbakke det er så at
bilene bare trille utover jeg hører dem ikke i det hele tatt 02827 ja # ja når du skal ut med bil så er det vel det 02959 men sa at jeg
da ikke bil for å reise på skitur # hvis jeg ikke skal inn i sl 02964 og nå har jeg blitt for lat så jeg gidder ikke å gå opp her he
jeg nå jeg må kjøre på Klebersvanga det er helst der jeg kjører for jeg tør ikke kjøre opp her p- med bil # (uforståelig) 02984 og jeg tenkte
jeg kommer aldri hjem att med denne derre bilen der 02990 og e så var jeg en ta disse postfolkene så da # " har du fått probleme
bilen ? " 02991 " ja " sa " jeg har gjort det " men e ia vi har skulle ringe nå n M4 han # " ia men det er ikke han som er # e der ien
```















# Saving results

CWB expression: "`[((lemma="bil" %c))]` ;"

Action : 

Hits four

Results 1



- count
- download
-   sort
-   collocations
-   annotate
-   show metadata
-   metadata distribution
-   delete hits
-   **save hits**



? (latter)



mer til å savner denne **bilen** # kommer seg frem # over alt (uforståelig)



mer til å savner n men e # er kjører jo så lite **bil** og da er det egentlig dyrt for meg å sitte # med et såpass lån på en bil jeg aldri bruker på måte



mer til å savner n men e # er kjører jo så lite bil og da er det egentlig dyrt for meg å sitte # med et såpass lån på en **bil** jeg aldri bruker på måte



  **alvdal\_02** ja det er det # det er dyrt å har **bil** # i\_hvert\_fall når du # har n to komma fem liter bensin så



  **alvdal\_01** småkonkurranser så # men e # ringte jo en i\_går morges å lurte på om han kunne komme å se på **bilen** min han skulle finne seg en Audi for åtti tusen



  **alvdal\_01** ja men du kjører **bil** når du har bil egentlig for det er jo ikke lange turen

  **alvdal\_01** ja men du kjører bil når du har **bil** egentlig for det er jo ikke lange turen

  **alvdal\_01** det er ikke mange som gjør det nå det var veldig vanlig før # e men nå e # er det mange som kjører **bil**

  **alvdal\_02** hvis ikke så sitter jeg på meg mor mi som kjører **bil** eller så kjører jeg traktor

  **alvdal\_02** så skal jeg ta førerkort på **bil** nå # når jeg bli atten # så er jo blir litt godt å ## endelig få kjøre bil

  **alvdal\_02** så skal jeg ta førerkort på bil nå # når jeg bli atten # så er jo blir litt godt å ## endelig få kjøre **bil**

# How to get several transcriptions quickly

Use a semi-automatic transliteration program

Input: a phonetically transcribed file

Processing:

Pick a similar dialect as a starting point

Transliterate to standard orthography, then train. Continue to transliterate, train again.

Output: aligned phonetic and orthographic transcriptions

# Dialekttranslitterator

## Transkripsjoner

+ Ny transkripsjon
+ Ny dialekt
- Slett valgte

Dialekt	Transkripsjonsfil	Dato	Last ned
roeros	roeros_04gk-amr_ifg_kk.trs	04.02.2011	<a href="#">A</a>
kirkenes	kirkenes_01um-kb_hl_kk.trs	17.02.2011	<a href="#">A</a>
kirkenes	kirkenes_02uk-kb_hl_kk.trs	17.02.2011	<a href="#">A</a>
bjugn	bjugn_15-19_sb_ifg.trs	17.02.2011	<a href="#">A</a>
bjugn	bjugn_15-aml_sb_ifg.trs	22.02.2011	<a href="#">A</a>
bjugn	bjugn_16-ma_sr_kk.trs	22.02.2011	<a href="#">A</a>
bjugn	bjugn_19-aml_sb_ifg.trs	22.02.2011	<a href="#">A</a>
kirkenes	kirkenes_03gm-kb_kk_lh.trs	24.02.2011	<a href="#">A</a>
kirkenes	kirkenes_04gk-kb_kk_lh.trs	24.02.2011	<a href="#">A</a>
bjugn	bjugn_23-pmk_hl_kk.trs	24.02.2011	<a href="#">A</a>
roeros	roeros_03gm-04gk_ifg.trs	24.02.2011	<a href="#">A</a>
kjoellefjord	kjoellefjord_03gm-04gk_sb_amg_lh.trs	28.02.2011	<a href="#">A</a>
beiarn	beiarn_01um-kb_kk_sb.trs	01.03.2011	<a href="#">A</a>
beiarn	beiarn_02uk-kb_kk_sb.trs	03.03.2011	<a href="#">A</a>
beiarn	beiarn_01um-02uk_kk_sb.trs	03.03.2011	<a href="#">A</a>
Oppdal	oppdal_02-31_eo_hl_ifg.trs	11.03.2011	<a href="#">A</a>
beiarn	beiarn_03gm-kb_sb_kk.trs	17.03.2011	<a href="#">A</a>
beiarn	beiarn_04gk-kb_sb_kk.trs	17.03.2011	<a href="#">A</a>
beiarn	beiarn_03gm-04gk_kk_sb.trs	17.03.2011	<a href="#">A</a>
lavangen	lavangen_03gm-ov_sb_eo.trs	17.03.2011	<a href="#">A</a>
lavangen	lavangen_04gk-ov_sb_eo.trs	19.03.2011	<a href="#">A</a>
lavangen	lavangen_03gm-04gk_sb_eo.trs	19.03.2011	<a href="#">A</a>
skaugdalen	skaugdalen_36-43_hl_kk_ifg.trs	21.03.2011	<a href="#">A</a>
vardoe	vardoe_01um-kb_eo_kk.trs	23.03.2011	<a href="#">A</a>
vardoe	vardoe_02uk-kb_eo_kk.trs	23.03.2011	<a href="#">A</a>
vardoe	vardoe_01um-02uk_eo_kk.trs	24.03.2011	<a href="#">A</a>
heroeyMR	heroeyMR_01um-ta_hl_lh.trs	28.03.2011	<a href="#">A</a>
heroeyMR	heroeyMR_02uk-ta_hl_lh.trs	29.03.2011	<a href="#">A</a>
heroeyMR	heroeyMR_01um-02uk_hl_lh.trs	29.03.2011	<a href="#">A</a>
Stranda	stranda_01um-ta_hl_kk.trs	30.03.2011	<a href="#">A</a>

## Fildeler

Filnavn	Sjekket?
kirkenes_01um-kb_hl_kk.trs.01	
kirkenes_01um-kb_hl_kk.trs.02	
kirkenes_01um-kb_hl_kk.trs.03	
kirkenes_01um-kb_hl_kk.trs.04	
kirkenes_01um-kb_hl_kk.trs.05	
kirkenes_01um-kb_hl_kk.trs.06	

## Sett backoff

0.01	aal
0.01	aamot
0.01	aaseral
0.01	aasnes
0.01	alvdal
0.2	andoeya
0.01	aremark
0.01	aukra
0.01	aure
0.01	aurland
0.2	Ballangen
0.01	bardu
0.01	bergen
0.01	bjugn
0.5	bodoe
0.01	boe
0.01	boemlo
0.01	botnhavn
0.01	brandbu
0.01	brekkom
0.01	brunlanes
0.01	bud
0.01	bykle
0.01	dalsbygda



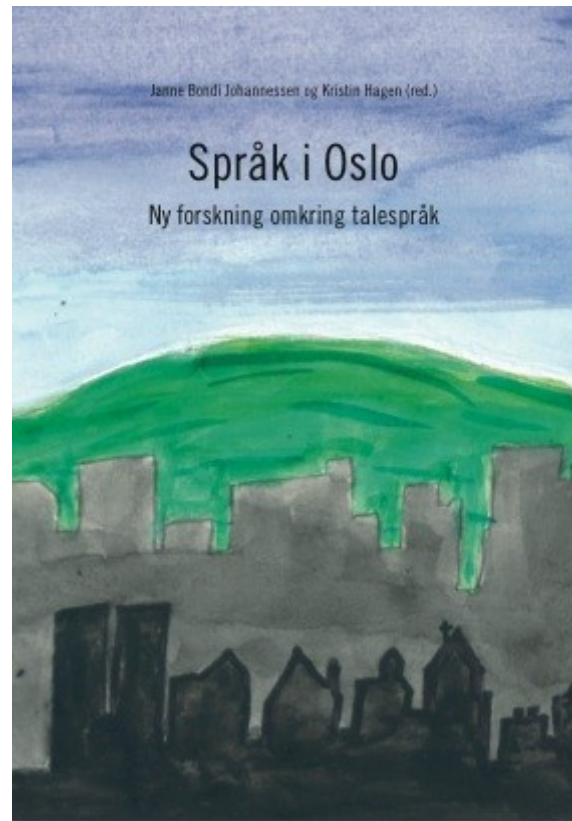
"ja"	"ja"	1
"je	"jeg	1
"je"	"jeg"	1
"jentene"	"jentene"	1
"jo	"jo	1
"jæi"	"jeg"	1
"ka	"hva	1
"ker	"hvor_	2
"kerr	"hvor_	4
"kerr"	"hvor_ "	4
"kær	"hvor_	1
"kærr	"hvor_	2
"kærr"	"hvor_ "	1
"kæssj"	"hvor_ "	1

Result of transliteration:  
a bilingual dictionary



- Nordic dialect corpus:
- [http://omilia.uio.no/glossa/html/index\\_dev.php?  
corpus=](http://omilia.uio.no/glossa/html/index_dev.php?corpus=scandiasyn)
- scandiasyn

# Book:



UNIVERSITETET I OSLO  
DET HUMANISTISKE FAKULTET

- Janne Bondi Johannessen and Kristin Hagen (eds.) 2008. *Language in Oslo – New research on spoken language*. Oslo: Novus.



## Ruth Vatvedt Fjeld (UiO):

### *Talespråksforskningens betydning for leksikografien*

- o Talespråksord mangler i ordbøkene:
  - o *Herre*, men ikke *derre*
  - o Ikke *loka*, *lættis*
  - o toalettpapir, klosettpapir, men ikke *dass*papir, *dass*rull
  
- o Betydningsspesialisering
  - --- det gjør jo automatisk at hvis alt går i *dass* da så --
  - --- oi jeg sølte (uforståelig) # (stønning) # jeg er helt *dass* i øya
  
- o Betydningsutvidelse
  - *skikkelig a2 (fra lty) bra, god, ordentlig ha s-e klær / få s- betaling / s-e folk / adv: oppføre seg s-*
  - den smakte *skikkelig* # forferdelig





## Toril Opsahl, Unn Røyneland og Bente Ailin Svendsen (UiO):

*”Syns du jallanorsk er lættis, eller?” – om taggen [lang=X] i NoTa-Oslo-korpuset*



- o Hva slags ord finnes i korpuset som ikke finnes i ordbøkene?
- o Hvilke befolkningsgrupper bruker mest unormerte ord?
- o Er det forskjeller?
- o Gutter fra østre deler av Oslo er mest kreative i dette nyordsperspektivet.





# Kjell Ivar Vannebo (UiO): *NoTa-informantene og tellemåten*

- Gammel tellemåte: *treogtredve*
- Ny tellemåte: *trettitre*
- Det er over femti år siden den nye tellemåten ble innført.
  - 1950, enstemmig i Stortinget, initiativ fra Telegrafverket
- Eneste gang en normering for talemål har funnet sted for norsk
- Hva kan man konkludere med grunnlag i NoTa-korpuset?
- Reformen har vært vellykket, selv om andre har hevdet det motsatte!
  - 18-25 år: 99 % ny tellemåte
  - 26-50 år: 49 %
  - 51+ år: 32 %





# Øystein Alexander Vangsnes (UiT): *Omkring adnominalt åssen/hvordan i Oslo-målet*



- *hvordan* poker var det du spilte da ?
- *åssen* bikkje er det men- ... ?
- *Åssen* er mer brukt i Oslo Rest enn i Oslo Vest
- Adnominalt *åssen* er vanligere enn adnominalt *hvordan*
- Adnominalt *åssen* er uvanlig blant folk med høy utdanning
- Betydningen til adnominalt *åssen*:
- Egenskap
  - *Åssen* bil kjøpte'ru'a? En rød en.
- Eksemplar
  - *åssen* topp ? ja e på ballen ... på toppen av ballen



# Janne Bondi Johannessen (UiO): *Psykologiske demonstrativer*



- o jo jeg var jo på, på  
Ammerudhjemmet, og da så satt  
*hun dama* som jeg snakka om  
som var snart hundre år (Dame,  
80, Grorud, Oslo)
- o men hva med han derre m leste  
du om han derre *han tyskeren*  
som hadde kuttet av utstyret på en  
fyr og spist det (Mann, 18, Vestre  
Aker)
- En ny type demonstrativ
  - De vanlige (den, denne)
    - Geografisk avstand – nær, fjern
  - De nye (han, hun)
    - psykologisk distanse





# Svein Lie (UiO):

## *Veldig sånn festejente*

- Nye betydninger av *veldig* og *sånn*
- *veldig* brukt med substantiv
  - man oppfatter et bestemt trekk eller egenskap ved substantivet, som så kan graderes.
- *Sånn* – et ekstremt vanlig talespråksord, har mange viktige bruksbetingelser
  - å signalisere høflighet!





# Gjert Kristoffersen (UiB) og Hanne Gram Simonsen (UiO):

## *Oslo! En undersøkelse av sl-sekvensen i NoTa-korpuset*



- Apikal, bakre, uttale eller fremre uttale av <sl>?
- 3477 belegg av <sl>
  - Fremlyd: over 99,6 % bakre uttale
  - Innlyd: 68 % bakre uttale
- Ordet *Oslo*:
  - 84,4 % bakre uttale (inkludert alle de unge)



# Elisabet Engdahl (GU): *Frågor i NoTa*



- o Spørsmålsformuleringer
  - o ja/nei-spørsmål
  - o Ekkospørsmål
    - o han sa måtte pappa måtte fikse det # for det han har installert det
    - o installert hva da ? (Ekkospørsmål)
  - o Spørsmål med deklarativ ordstilling
    - o du har ikke vært der du heller nei ?
- o bruksbetingelser
- o NoTa-korpuset som forskningsverktøy
  - o hvordan transkribørene har annotert spørsmålene
  - o hvordan korpuset fungerer.



## Inger Margrethe Hvenekilde Seim (UiO): *Innhold og struktur i en samtale mellom to ungdommer i et flerkulturelt miljø i Oslo*

- To gutter med innvandrerbakgrunn i en timelange samtale.
  - Hvordan forløper samtalen?
  - Hvordan tas emner opp?
- Emnet etnisitet
  - Egen gruppetilhørighet
  - Hverandres gruppetilhørighet
  - Andres gruppetilhørighet
  
  - Akkurat dette emnet kommer guttene inn på flere ganger gjennom hele samtalen.







# Jan Svennevig (BI/UiO):

## *”Ikke sant” som respons i samtale*



- Sammenligning av flere små, eldre samtalekorpus og NoTa-korpuset
- *ikke sant* brukt som respons
  - har økt fra null i 1994 til en tredel av alle forekomster i 2005.
  - mest vanlig blant yngre mennesker.
- flere betydninger for uttrykket brukt som respons
  - felles for dem er at det brukes som bekreftelse på noe som er sagt i samtalen.
- Kan virke ganske irriterende på samtalepartneren, og Svennevig forklarer hvorfor.



# Lars-Olof Delsing (LU):

## *Viskningar och rop – eller hur vi undrar och förundras*

- o Ordstillingen og ordene som er brukt i utrops- og undringskonstruksjoner.
  - o Så du har vackra rosor! (f.svensk)
  - o Så vakre roser du har! (norsk)
  
  - o Vad/så många slipsar du har! (svensk)
  - o Så mange slips du har! (norsk)
  
- o En enorm variasjon mellom de skandinaviske språkene. Men likevel en del generaliseringer.





## Gunnar Hrafn Hrafnbjargarson (UiO/UiT):

### *Substantiverte adjektiv: Det er verste jeg har hørt*

- o Nakne superlativer : En forholdsvis ukjent konstruksjonstype
- o Svakt bøyd, men mangler artikkel,
- o Ligner på utbryting, men er det ikke
  - o det er bare fugler som kan fly
  - o det er eneste som er minuset syns jeg da
- o Dette er ikke utbrytinger, slik de tilsynelatende kan se ut
  - o Det som kan fly er bare fugler
  - o \* Det som er minuset synes jeg er eneste
- o men er faktisk en type ikke-referensiell, substantivert predikativ.







## Marit Julien (LU):

*Så vanleg at det kan ikkje avfeiast – om  
V2 i innføyde setningar.*



- Ordstillingen i underordnede setninger
  - Men det som er er at *han kan ikke* lage sanger.
  - Så jeg bare sier at det kan jeg ikke gjøre.
- Hva slags betingelser er mulige?
- Det ikke spiller noen direkte rolle hva slags verb som innleder at-setningen.
- Det som teller, er utelukkende semantikk: at-setningen må være hevdet.
  - \*Men jeg tviler på at slike konserter hjelper faktisk mot volden.



# Mari Nygård, Kristin M. Eide og Tor A. Åfarli (NTNU): *Ellipsens syntaktiske struktur*



UNIVERSITETET I OSLO  
DET HUMANISTISKE FAKULTET



- Lydløse ord. Uttalte, lydløse eller stumme ord. Ellipse.
  - \_\_ fikk ikke lov å herje som vi ville.
  - nei jeg var vel en ti-tolv år tenker jeg, så \_\_ var ikke så gammel ...
- Ikke noen allmenn forståelse blant språkforskerne om fenomenet ellipse.
  - Noen: det som ikke er uttalt, finnes heller ikke.
  - Nygård, Eide og Åfarli: ellipse er like eksisterende i språket som de ordene vi faktisk hører.
  - de tilsynelatende stumme ordene kan bare forekomme i klart definerte kontekster, og ikke er tilfeldig spredt rundt omkring.



# Fredrik Jørgensen (UiO): *Automatisk gjenkjenning av ytringsgrenser i talespråk*



- Talespråket er en utfordring som må løses for at man siden skal kunne analysere talespråkets grammatikk automatisk.
- Skriftspråk er tydelig avgrenset ved tegnsetting,
- Transkribert talespråk er ikke inndelt på noen enhetlig måte, så her må det nytenkning til.
- Viktig å utvikle automatiske metoder for å finne ytringsgrenser i NoTa-korpuset,
- Noen foreløpige tall viser hvor vellykket det har blitt.



# Peter Juel Henriksen (KU):

## *NoTa – nu med lydskrift*

- Å lage norsk lydskrift av norsk ortografisk transkripsjon på, ved å gå via dansk!
- Systemet NoTaFon
- Utnytter
  - at NoTa-korpuset er grammatisk tagget med ordklasser og bøyning
  - noen enkle ordlister over de vanligste ordene som er mest forskjellige på dansk og norsk.
- Får automatisk laget norske tekst om til dansk tekst
- I dansk kan han så bruke et system som finnes fra før for å gå fra ortografi til lydskrift.
- Setter opp noen enkle regler for hvordan dansk og norsk lyd skiller seg fra hverandre (bløte og harde konsonanter, for eksempel),
- Får et automatiske system fra dansk lydskrift til norsk lydskrift!







## Åshild Søfteland og Anders Nøklestad (UiO):

*Manuell morfologisk tagging av NoTa-materialet med støtte fra en statistisk tagger*



- Talespråk er fullt av gjentakelser, avbrudd og pauser
- Derfor vanskelig å bruke eksisterende analysemetoder, f.eks. Tagger
- Alternativ: Lage en statistisk tagger for talespråk:
  - Begynn med å tagge en del av teksten manuelt; dette blir det første treningskorpuset
  - Tren taggeren på treningskorpuset
  - Kjør taggeren på en ny del av teksten
  - Rett opp feila som taggeren har gjort
  - Legg den oppretta teksten til treningskorpuset
  - Gå til punkt 2
- Metode: Treetagger. Resultat: 96,9 % korrekt.



# Victoria Rosén (UiB):

## *Mot en trebank for talespråk*

- Diskuterer viktigheten av å ha en samling av syntaktisk annoterte setninger
  - Viktig for språkteknologisk utvikling
  - Viktig for språkteknologisk evaluering
- Vanskelig med syntaktisk annotering av talespråk
- En løsning: å foreta manuell for-redigering for senere automatisk analyse



# Janne Bondi Johannessen: *Oslospråket i tall*



UNIVERSITETET I OSLO  
DET HUMANISTISKE FAKULTET



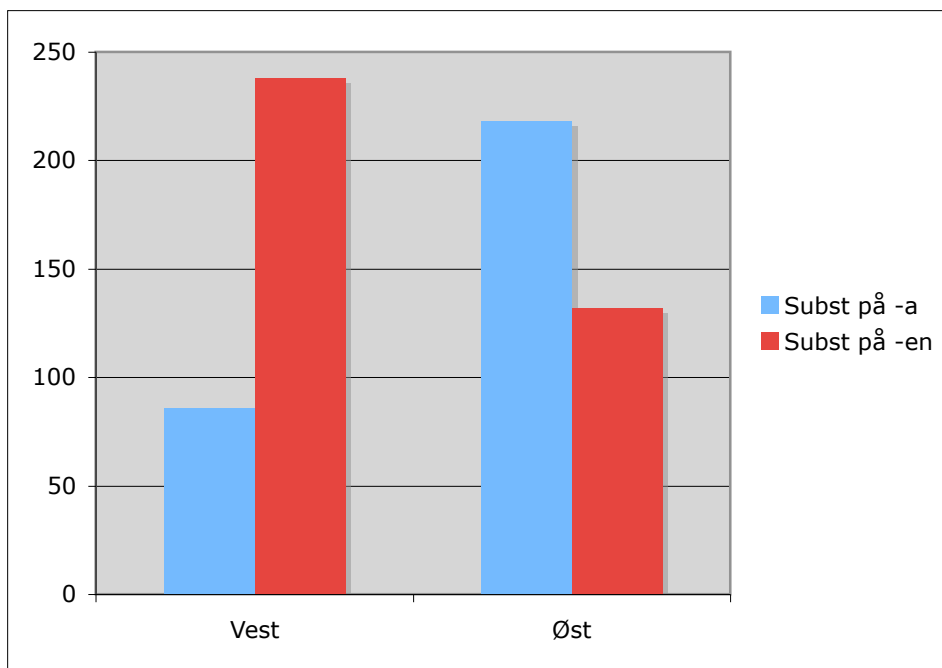
photo by Vic Alexander ©1998

Uekstlab.



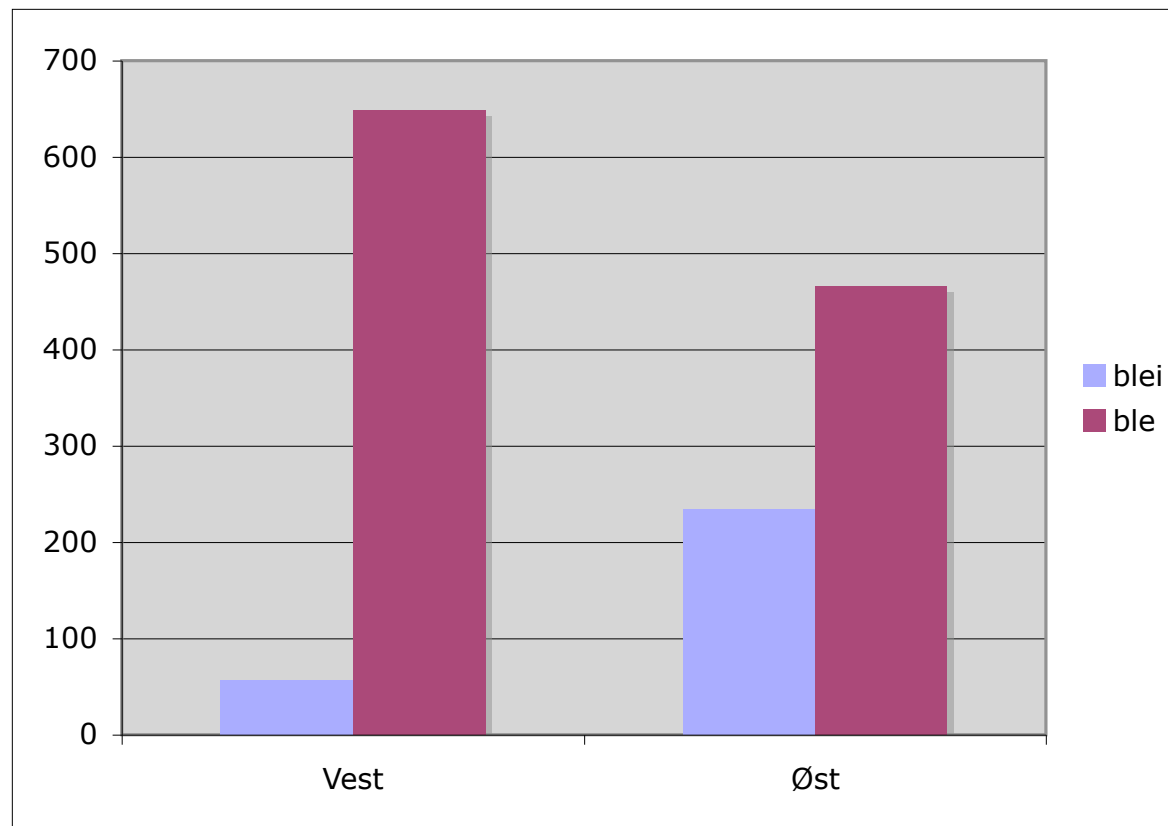
## Øst og vest:

***Elva eller elven, avisa eller avisen, gata eller gaten, boka eller boken, tida eller tiden, sola eller solen?***





# Øst og vest: Diftonger – *blei* eller *ble*?





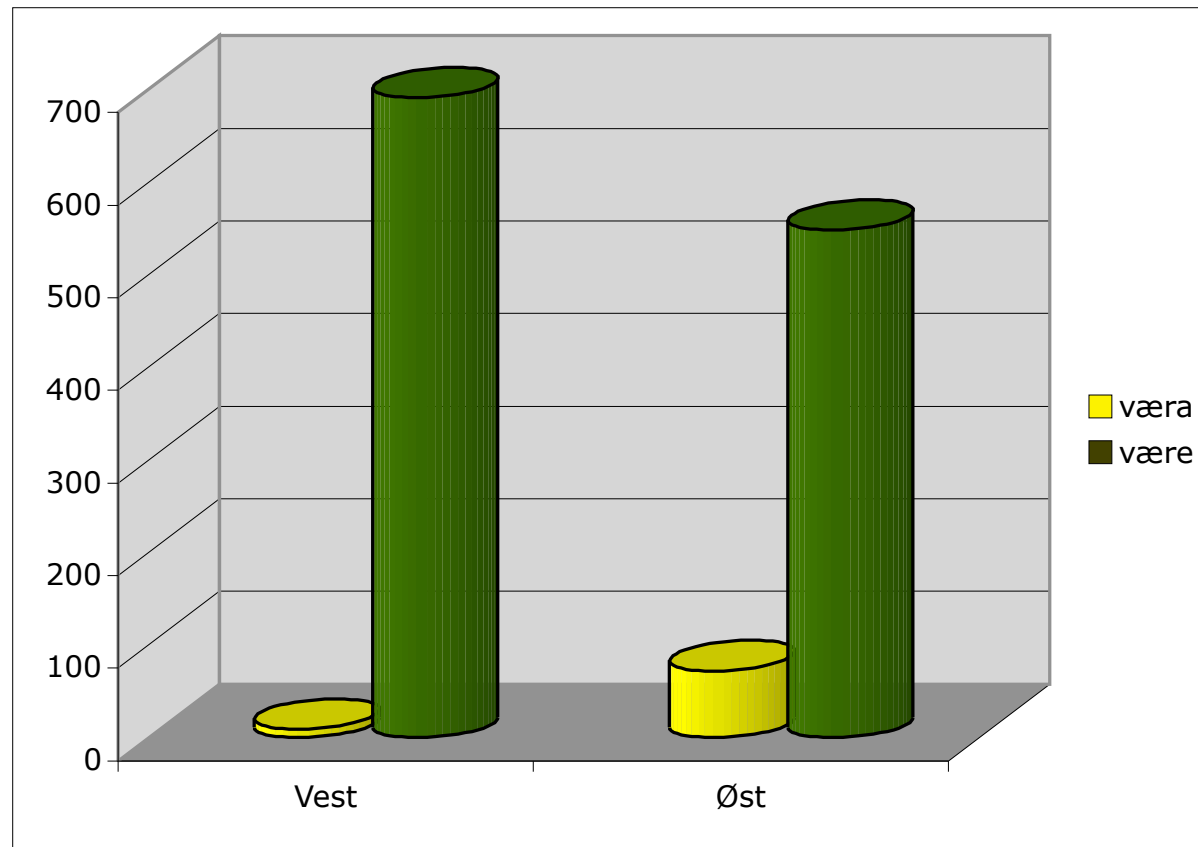
# Øst og vest:

## Kløyvd infinitiv – *Væra* eller *være*?

(antall forekomster)



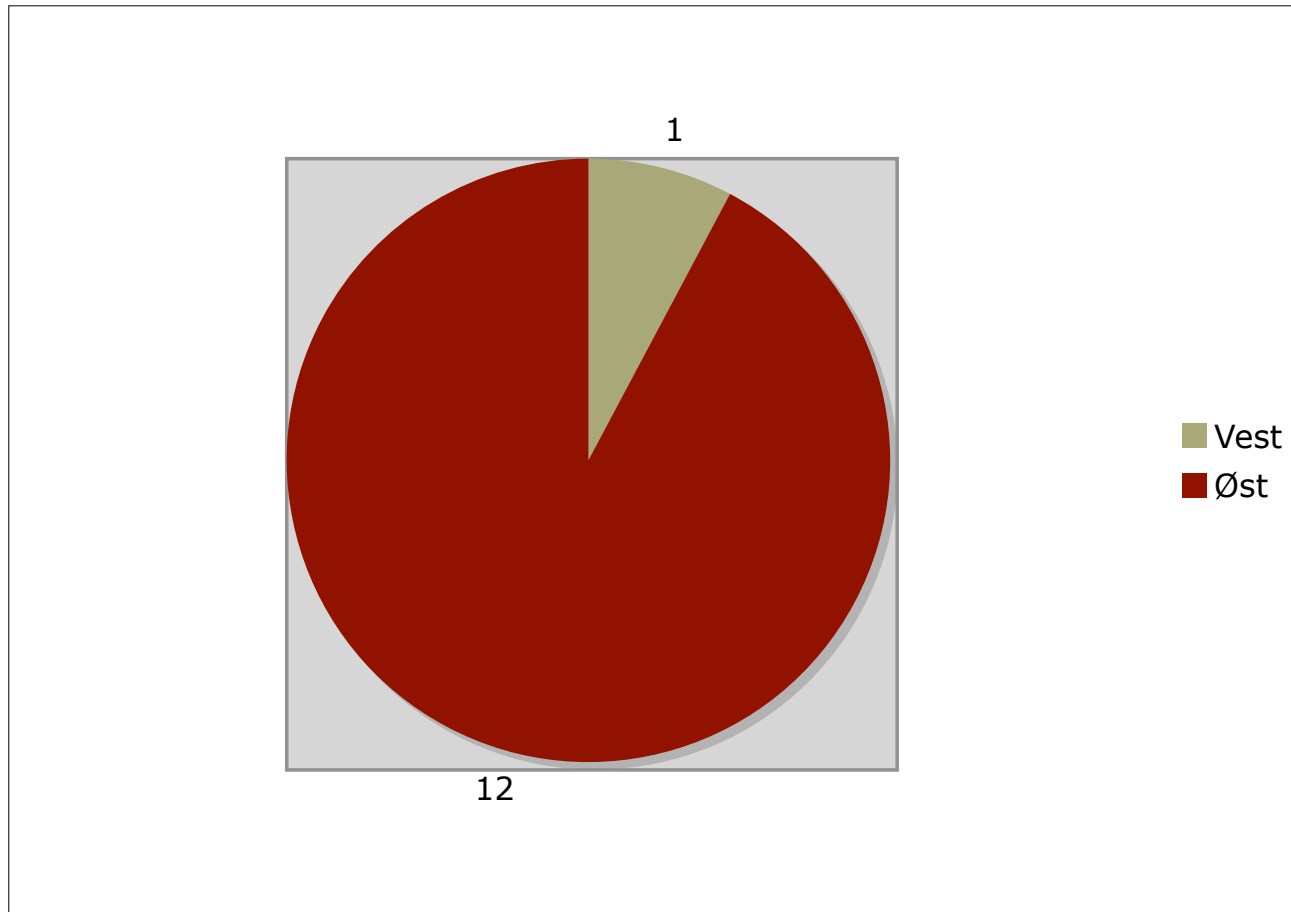
UNIVERSITETET I OSLO  
DET HUMANISTISKE FAKULTET





# Kløyvd infinitiv i prosent: *væra*

i prosent av alle ganger verbet brukes i begge former i øst og vest

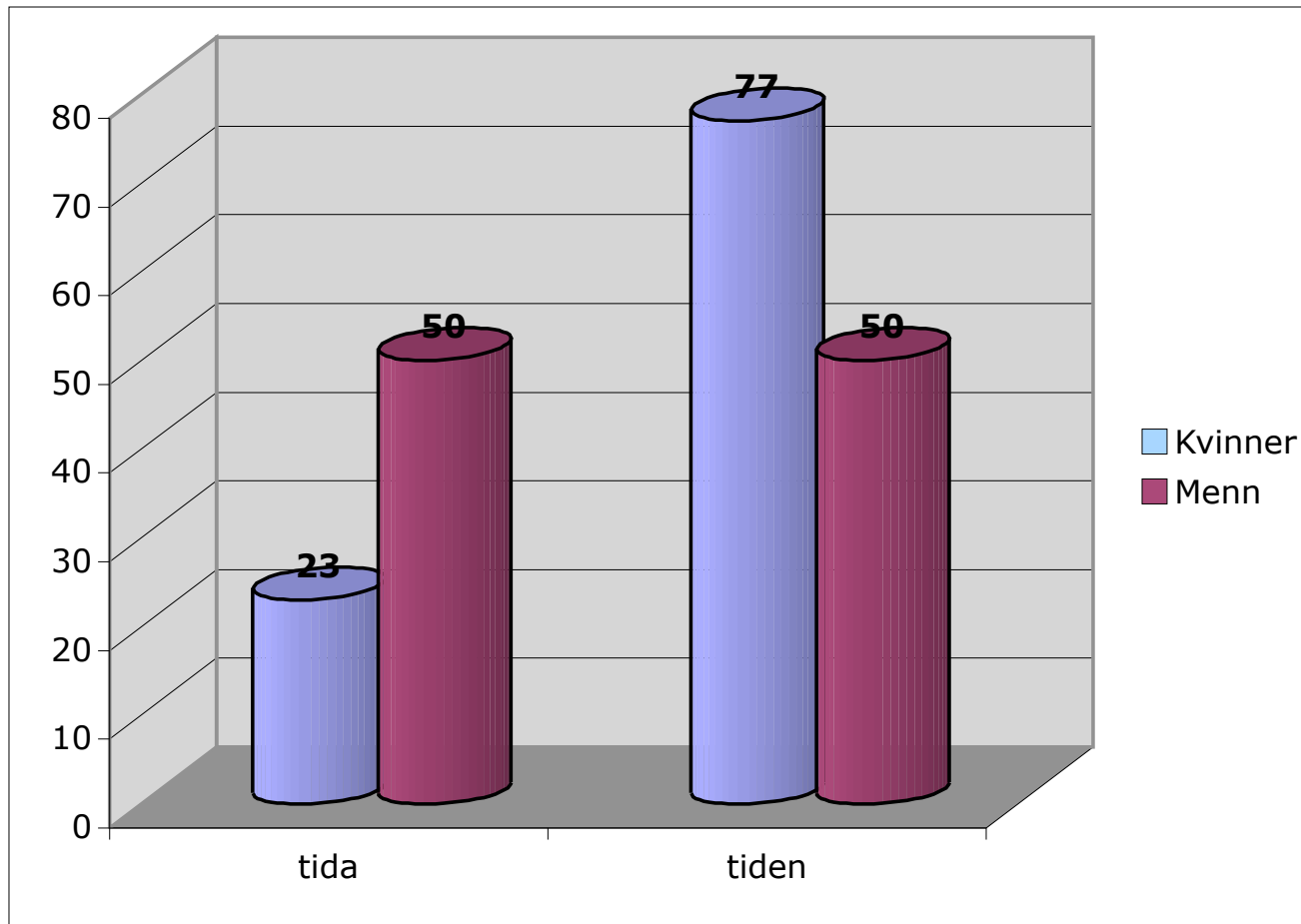


# Kjønn: *tida eller tiden?*

(Prosent, kjønnsvis)



UNIVERSITETET I OSLO  
DET HUMANISTISKE FAKULTET

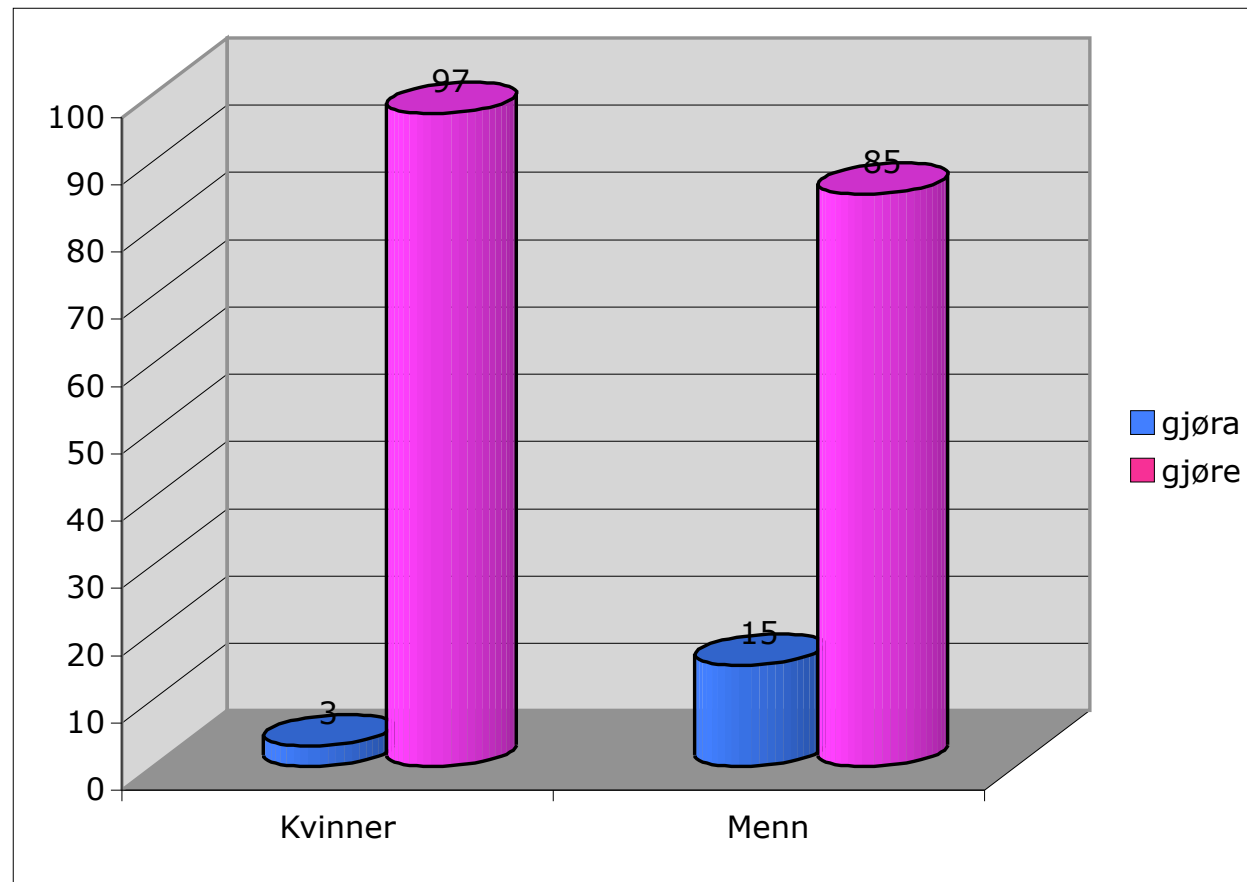


# Kjønn:

## Kløyvd infinitiv (*gjøra* eller *gjøre*)



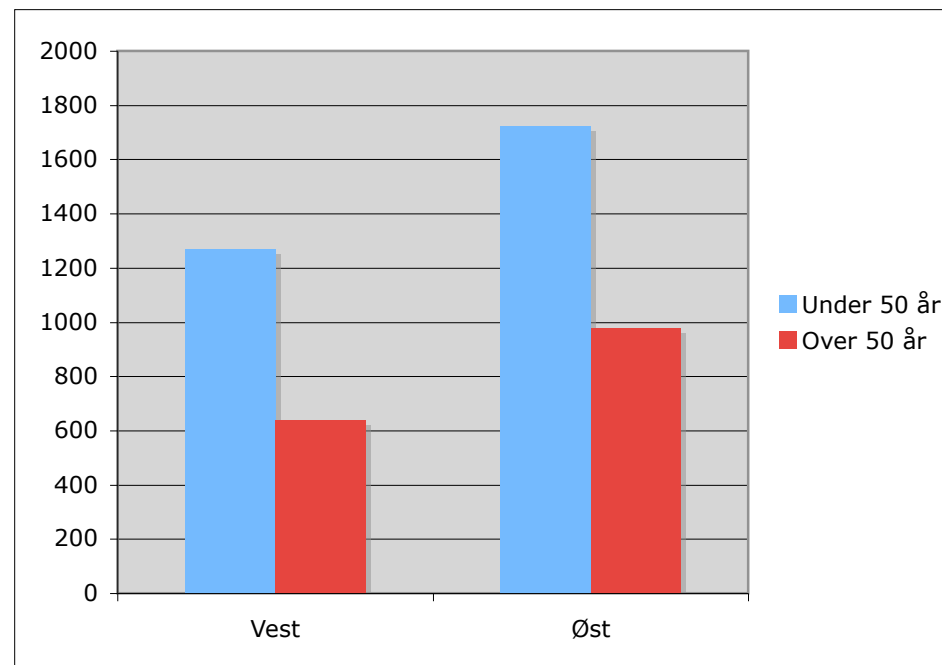
UNIVERSITETET I OSLO  
DET HUMANISTISKE FAKULTET



# Alder: Substantiver på -a (*kona, gata*)



UNIVERSITETET I OSLO  
DET HUMANISTISKE FAKULTET

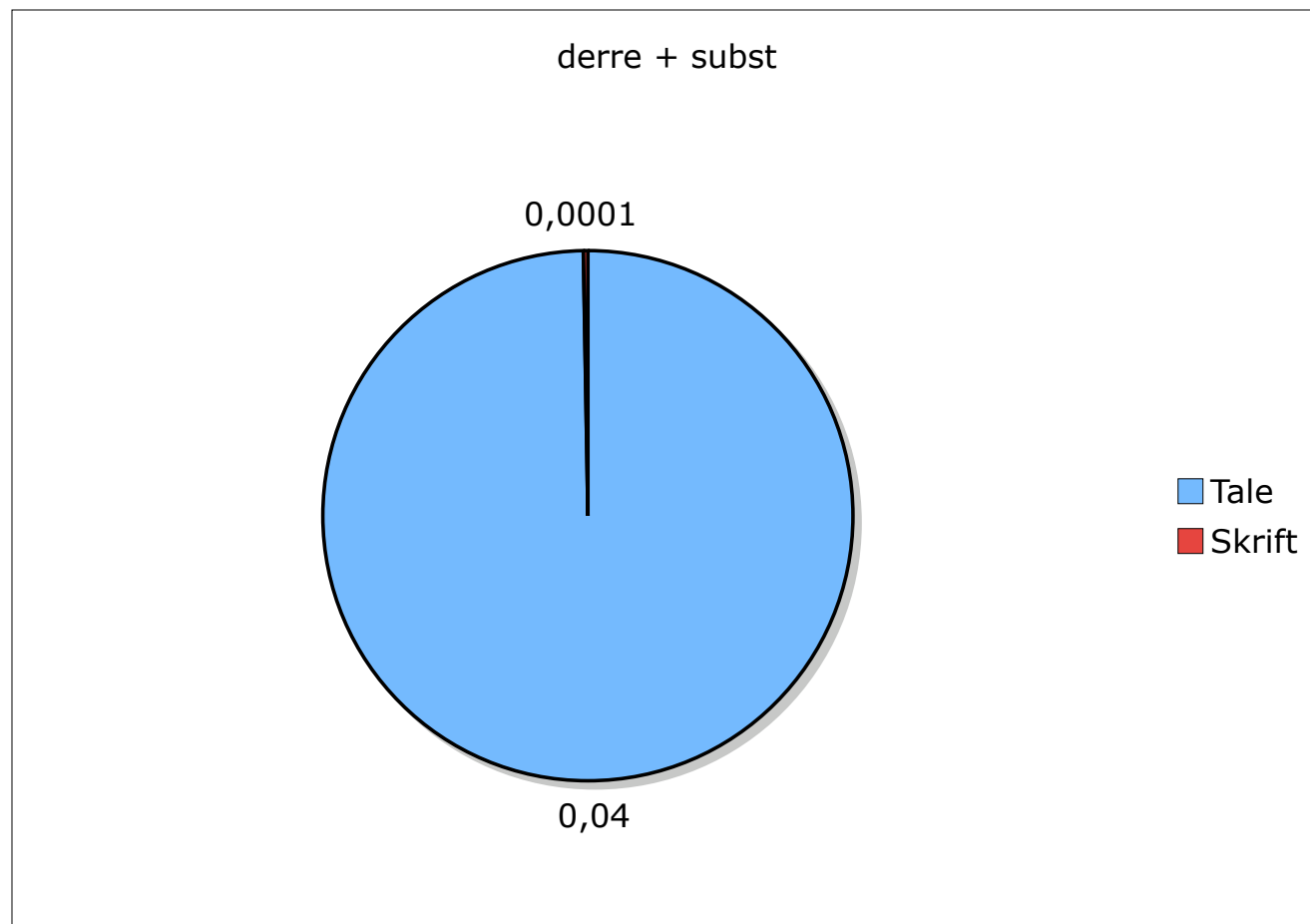


# Skrift og tale: *derre* + substantiv

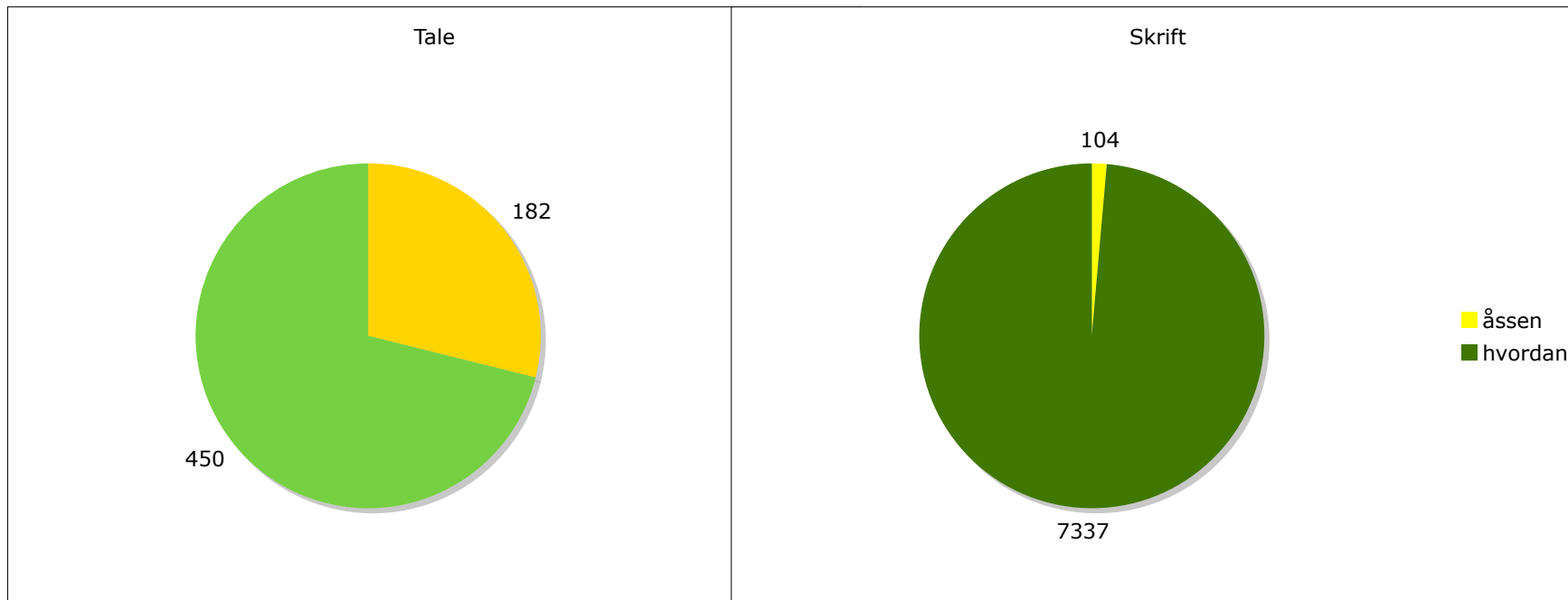
(NoTa-korpuset og Oslo-korpuset av taggede, norske tekster)



UNIVERSITETET I OSLO  
DET HUMANISTISKE FAKULTET



# Tale og skrift: *åssen* og *hvordan*





# And from now on: Spoken corpora of some of the Ethiopian languages!