



Lumberjacking with GrETEL

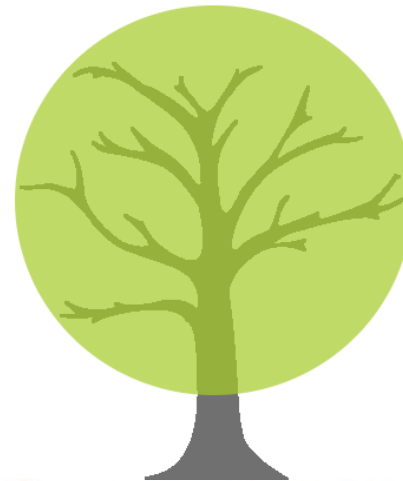
Liesbeth Augustinus
Vincent Vandeghinste
Ineke Schuurman
Frank Van Eynde

BKL taaldag - May 11, 2013



GrETEL

- **Greedy Extraction of Trees for Empirical Linguistics**
- Query engine for treebanks
- **Nederbooms project**
Exploitation of Dutch treebanks
for research in linguistics



CLARIN

Common Language Resources and Technology Infrastructure

KU LEUVEN

GrETEL

- **Greedy Extraction of Trees for Empirical Linguistics**
- Query engine for treebanks
- **Nederbooms project**
Exploitation of Dutch treebanks for research in linguistics
- **Goals**
 - User-friendly tools
 - Access to large data files
 - Fast and accurate



GrETEL

- **Greedy Extraction of Trees for Empirical Linguistics**
- Query engine for treebanks
- **Treebank** = syntactically annotated corpus
e.g. Penn Treebank (English), TüBa (German),
LASSY, CGN (Dutch)

TREEBANKS

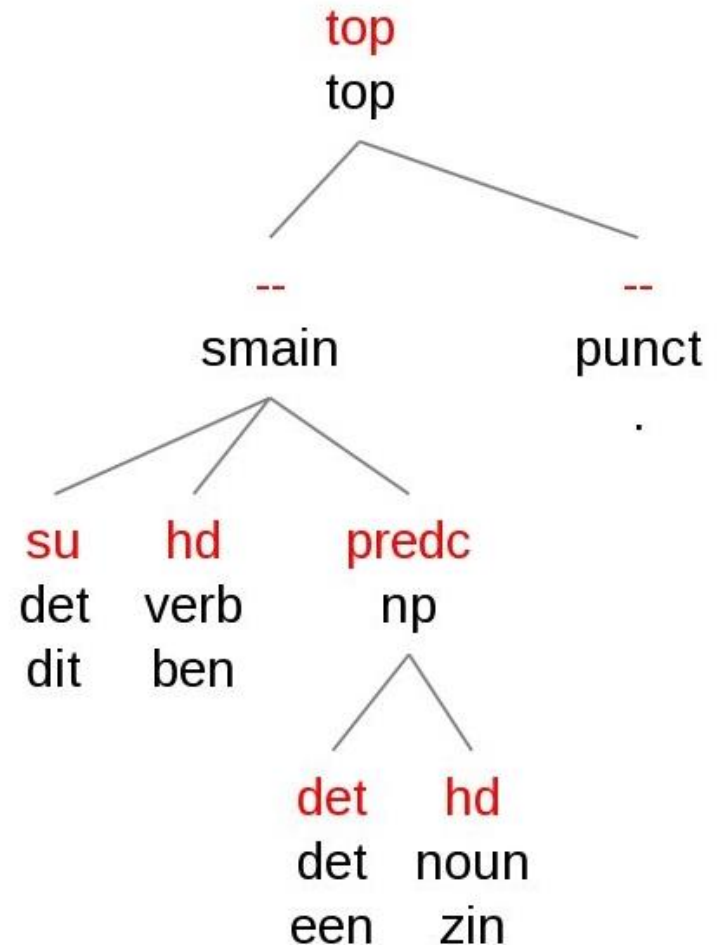
CGN core corpus	LASSY small
Spoken Dutch	Written Dutch
Stylistic & regional differences conversations vs read texts NL vs VL	Stylistic differences Wikipedia vs legal texts
± 1M tokens	± 1M tokens
130k sentences	65k sentences
Manually corrected	Manually corrected

GrETEL

- **Greedy Extraction of Trees for Empirical Linguistics**
- Query engine for treebanks
- **Treebank** = syntactically annotated corpus
e.g. Penn Treebank (English), TüBa (German),
LASSY, CGN (Dutch)
- **Parser**
e.g. Alpino (Van Noord 2006)

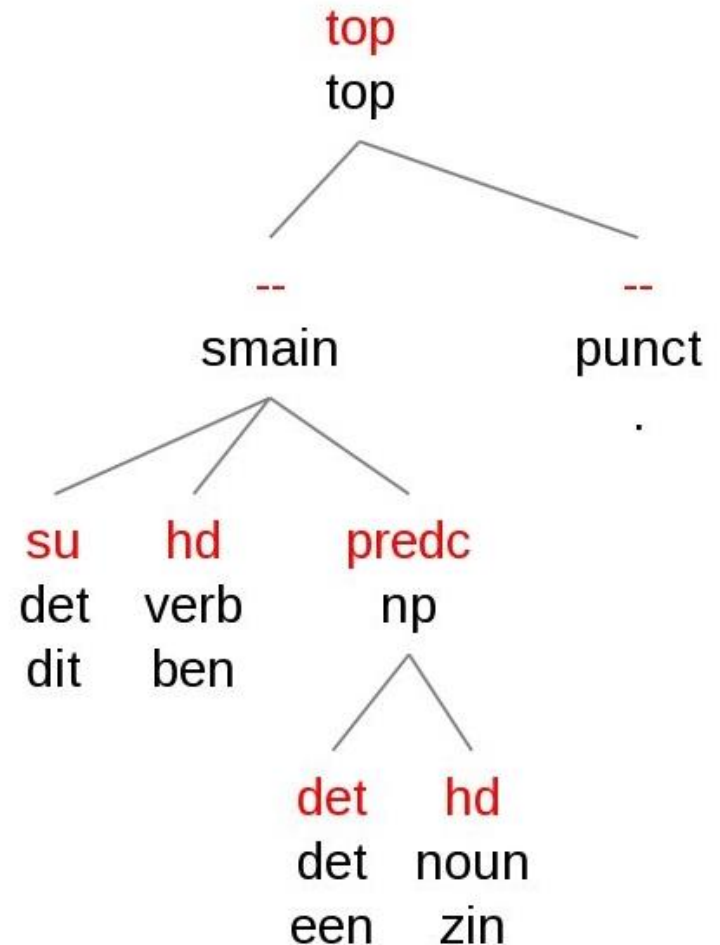
ALPINO PARSER

Dit is een zin. >> ALPINO parser >>
“This is a sentence.”



ALPINO PARSER

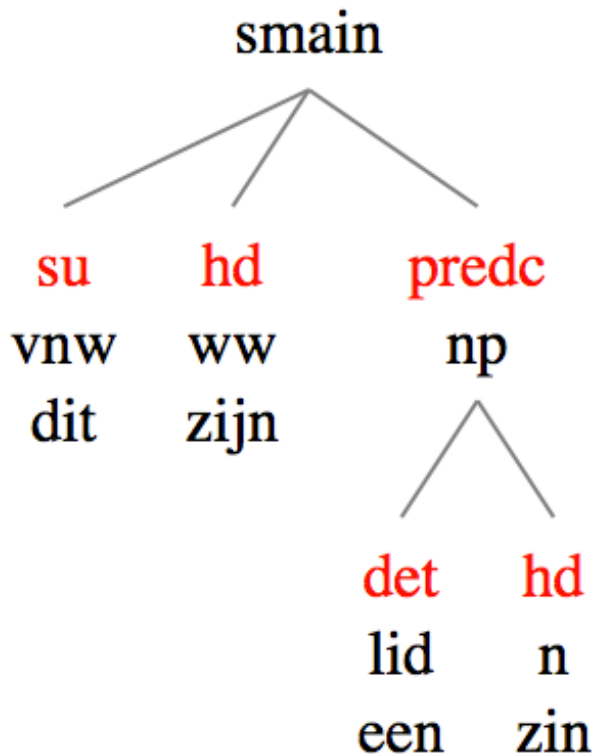
Dit is een zin. >> ALPINO parser >>
"This is a sentence."



XML trees

Query language: **XPath**

XPATH



```
//node[@cat="smain" and  
node[@rel="su" and  
@pt="vnw" and @lemma="dit"]  
and node[@rel="hd" and  
@pt="ww" and @lemma="zijn"]  
and node[@rel="predc" and  
@cat="np" and  
node[@rel="det" and  
@pt="lid" and @lemma="een"]  
and node[@rel="hd" and  
@pt="n" and @lemma="zin"]]]
```

XPATH



```
//node[@cat="smain" and  
node[@rel="su" and  
@pt="vnw" and @lemma="dit"]  
and node[@rel="hd" and  
@pt="ww" and @lemma="zijn"]  
and node[@rel="predc" and  
@cat="np" and  
node[@rel="det" and  
@pt="lid" and @lemma="een"]  
and node[@rel="hd" and  
@pt="n" and @lemma="zin"]]]
```

XPATH



```
//node[@cat="smain" and  
node[@cat="su" and  
@pt="w" and @rel="dit"]  
and  
@pt="w" and @rel="zijn"]  
and node[@cat="dc" and  
@cat="n"  
node  
@pt="w" and @rel="seen"]  
and node[@rel="ho" and  
@pt="n" and @lemma="zin"]]]
```

XPATH



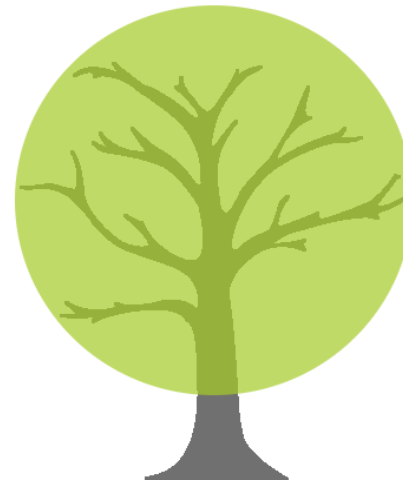
GrETEL

- **Greedy Extraction of Trees for Empirical Linguistics**
- **Query treebanks by example**



GrETEL

- **Greedy Extraction of Trees for Empirical Linguistics**
- **Query treebanks by example**
- First version
 - => only for LASSY treebank
- New release
 - => GrETEL for CGN treebank
 - => update based on user reviews





the user

- Example sentence

- Indicate relevant items of the sentence

- (Adapt XPath)
- Select treebank

- Inspect results



GrETEL

- Parser (Alpino)

- Automatically generate XPath expression

- Present results

OUTLINE

- GrETEL in a nutshell
- **GrETEL demo**
 - **Case study**
 - Search options
- Conclusions and future work

CASE STUDY

- Collective noun constructions
 - E.g. Een aantal bomen zijn omgevallen.
'A number of trees fell down.'
 - DET + NOUN + PLURAL NOUN

CASE STUDY

- Collective noun constructions
 - E.g. Een aantal bomen zijn omgevallen.
'A number of trees fell down.'
 - DET + NOUN + PLURAL NOUN
- Discontinuous constructions!
 - E.g. Een groot aantal oude bomen zijn omgevallen.
'A large number of old trees fell down.'

GrETEL ONLINE

Contact



Nederbooms

Home > Tools > GrETEL

GrETEL

Greedy Extraction of Trees for Empirical Linguistics

GrETEL is a query engine in which linguists can use a natural language example as a starting point for searching a treebank with limited knowledge about tree representations and formal query languages. By allowing linguists to search for constructions which are similar to the example they provide, we hope to bridge the gap between traditional and computational linguistics.

KU LEUVEN

About

Tools

GrETEL

GrETEL online

GrETEL for LASSY

GrETEL for CGN

Documentation

INPUT



Nederbooms

[Home](#) > [Tools](#) > [GrETEL](#) > [GrETEL online](#) > [GrETEL for CGN](#)

◊ About

▼ Tools

▼ GrETEL

▼ GrETEL online

○ GrETEL for
LASSY

○ GrETEL for CGN

GrETEL for CGN (v1.1)

Please provide an **input example**

ANNOTATION MATRIX

GrETEL for CGN

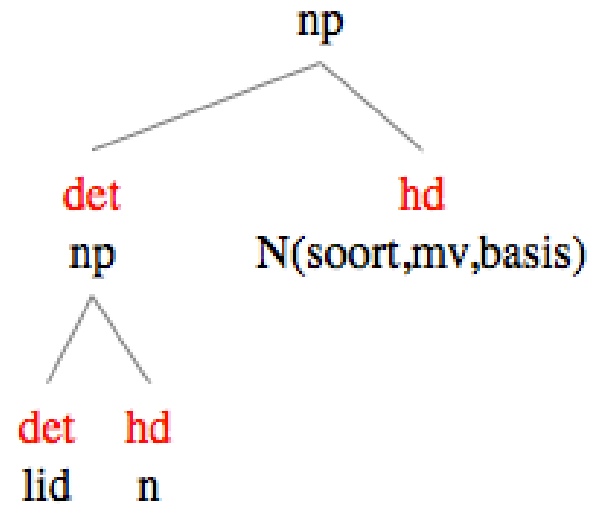
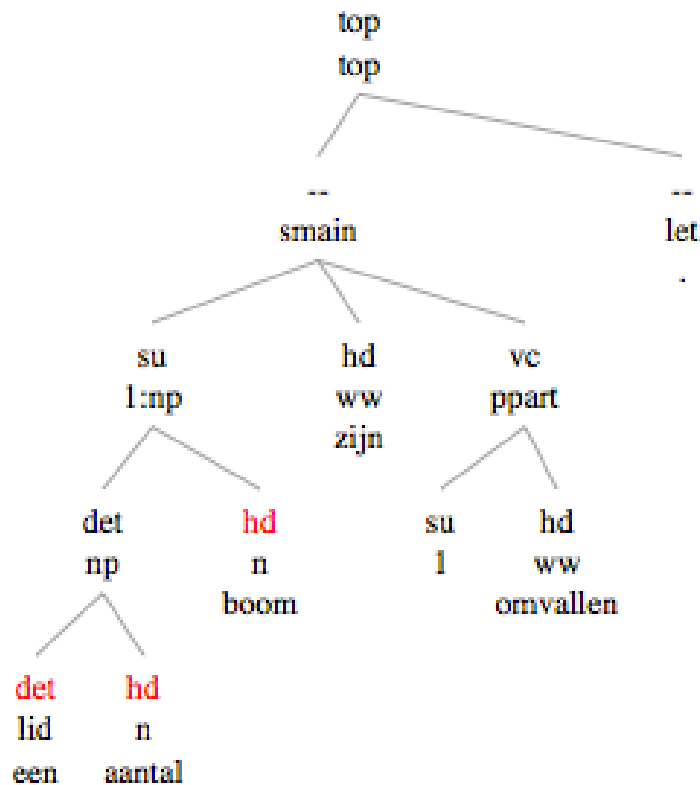
Please indicate the relevant parts of the sentence

sentence		Een	aantal	bomen	zijn	omgevallen	.
relevant nodes	pos	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
	extended pos	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
	lemma	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
	token	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
optional nodes		<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>

ANNOTATION GUIDELINES

GUIDELINES

- **pos:** Short part-of-speech tag. The different tags are: *n*, *ww*, *adj*, *lid*, *vnw*, *vg*, *bw*, *tw*, *vz*, *tsw*, *spec*, and *let*.
- **extended pos:** Long part-of-speech tag. For example:
N(soort,mv,basis), *WW(pv,tgw,ev)*,
VNW(pers,pron,nomin,vol,2v,ev). [list of all pos tags]
- **lemma:** Word form that generalizes over inflected forms. For example: *zin* is the lemma of *zin*, *zinnen*, and *zinnetje*; *gaan* is the lemma of *ga*, *gaat*, *gaan*, *ging*, *gingen*, and *gegaan*. Lemma is case insensitive (except for proper names).
- **token:** The exact word form. This is a case sensitive feature.



XPATH GENERATOR

XPath query generated from the input example. You can adapt it if necessary.

```
//node[@cat="np" and node[@rel="det" and @cat="np" and node[@rel="det" and @pt="lid"] and node[@rel="hd" and @pt="n"]] and node[@rel="hd" and @pt="n" and @postag="N(soort,mv,basis)"]]
```

TREEBANK SELECTION

CGN core corpus

Treebank		Contents	# Sentences	# Words	# Sentences	# Words	# Sentences	# Words
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>		NL		VL		TOTAL	
NL	VL							
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Spontaneous conversations ('face-to-face')	50,239	302,828	22,881	147,418	73,120	450,246
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Interviews with teachers of Dutch	2,484	25,724	4,289	34,158	6,773	59,882
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Telephone conversations (recorded via a switchboard)	11,649	70,084	3,142	19,984	14,791	90,068
<input type="checkbox"/>	<input checked="" type="checkbox"/>	Telephone conversations (recorded on MD)	0	0	929	6,309	929	6,309
<input checked="" type="checkbox"/>	<input type="checkbox"/>	Simulated business negotiations	3,123	25,524	0	0	3,123	25,524
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Interviews/discussions/debates (broadcast)	6,290	75,167	2,617	25,122	8,907	100,289
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	(Political) discussions/debates /meetings (non-broadcast)	1,166	25,125	543	9,009	1,709	34,134
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Lessons recorded in the classroom	3,064	26,004	1,395	10,116	4,459	36,120
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Live (sports) commentaries (broadcast)	2,251	25,002	1,026	10,147	3,277	35,149

RESULTS

Collective noun constructions

- E.g. *Een aantal bomen* *zijn omgevallen.*
'A number of trees fell down.'
- DET + NOUN + PLURAL NOUN

➔ 594 matches in 576 sentences

RESULTS: table

SEARCH RESULTS

treebank	hits	matching sentences	sentences	ratio (matching sentences/sentences)
NA	92	88	50239	0.18 %
VA	79	76	22881	0.33 %
NB	29	28	2484	1.13 %
VB	24	23	4289	0.54 %
NC	28	28	11649	0.24 %
VC	9	9	3142	0.29 %
VD	2	2	929	0.22 %
NE	5	5	3123	0.16 %
NF	60	56	6290	0.89 %
VF	27	27	2617	1.03 %
NG	27	26	1166	2.23 %
VG	11	11	543	2.03 %
NH	11	11	3064	0.36 %
VH	5	5	1395	0.36 %

RESULTS: data

Sentence ID	Matching sentences	
fvb400145__607	maar xxx na zo'n xxx bijvoorbeeld we heb*a we hebben dan vroeger soaps helemaal uitgeschreven in draaiboeken een paar scènes hè want anders is dat veel te lang en dan die in scène gezet en gefilmd en dan moet daarna de theorie komen rond soap van wat is een soap wie kijkt er naar soaps kenmerken van ...	[full screen] [XML]
fvb400165__202	alleen d'r zijn toch een een een pak woorden waar er geen regel voor bestaat .	[full screen] [XML]
fvb400165__357	en in die toneelvereniging we zaten daar met een heel uh bende spelers laat ons zeggen met vijftwintig spelers waarvan er zeker twintig uh toneel gevolgd hadden .	[full screen] [XML]
fvb400165__58	die werd wel gelezen doo*a door een bepaald aa*a door een een ja een groep studenten eigenlijk uh .	[full screen] [XML]
fvb400165__63	hij herhaalt zichzelf uh na een aantal bladzijden hè .	[full screen] [XML]

RESULTS: data

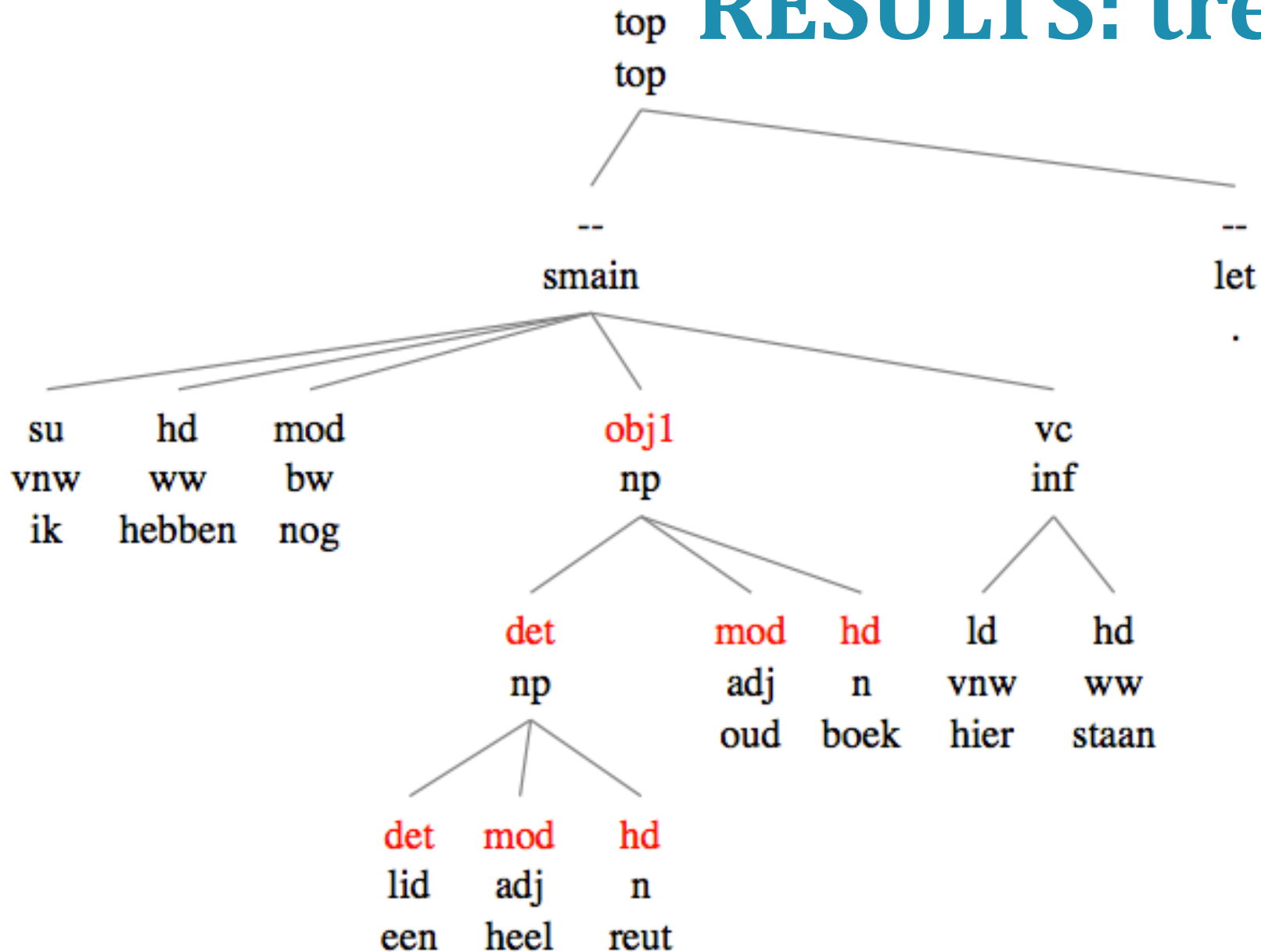
Sentence ID

Matching sentences

fnc008003__58	dan dan zijn jullie gewoon een stelletje mutsen die aan de telefoon hangen .	[full screen] [XML]
fnc008005__117	ik had vorige week had ik uh een paar batterijen gehaald .	[full screen] [XML]
fnc008007__272	ik bedoel uh de Queen Mum die dronk tenminste nog een paar whisky's elke avond om vervolgens in een diep coma weg te dommelen .	[full screen] [XML]
fnc008006__344	ik heb hier nog een hele reut ouwe boeken staan .	[full screen] [XML]
fnc008006__77	was een kerel en die had z'n die die timmerde z'n vrouw l*a al uh structureel in elkaar een aantal jaren .	[full screen] [XML]
fnc008007__71	ja d'r zijn toch een aantal kinderen die 't gezien hebben .	[full screen] [XML]
fnc008012__58	met een aantal auto's ?	[full screen] [XML]
fnc008013__100	want die eten in de zomer een beetje bessen of zo en en uh ete*a eten ze alleen maar vlees hè .	[full screen] [XML]

ik heb hier nog een hele reut ouwe boeken staan .

RESULTS: trees



OUTLINE

- GrETEL in a nutshell
- **GrETEL demo**
 - Case study
 - **Search options**
- Conclusions and future work

SEARCH OPTIONS

→ Below annotation matrix

OPTIONS

- Respect word order
- Ignore properties of the dominating node
- Split extended pos tags

SEARCH OPTIONS

Green versus red word order in Dutch

- green: past participle – auxiliary

*De NAVO stelt dat ze er alles aan **gedaan heeft***

- red: auxiliary – past participle

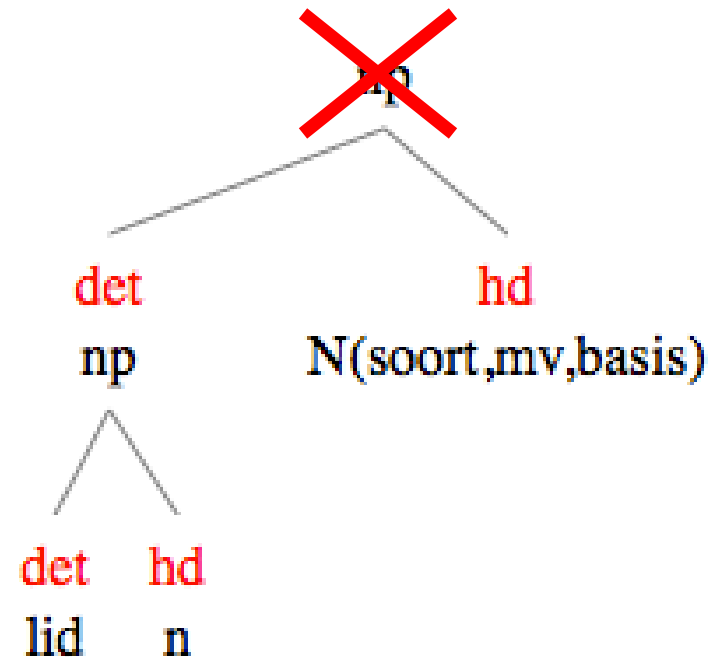
*De NAVO stelt dat ze er alles aan **heeft gedaan***

“The NATO claim that they have done everything in their power”
(deredactie.be)

SEARCH OPTIONS

OPTIONS

- Respect word order
- Ignore properties of the dominating node
- Split extended pos tags



XPath query generated from the input example. You can adapt it if necessary.

```
/*node[@cat="np"] and node[@rel="det" and @cat="np" and node[@rel="det" and @pt="lid"] and node[@rel="hd" and @pt="n"]] and node[@rel="hd" and @pt="n" and @postag="N(soort,mv,basis)"]]
```

SEARCH OPTIONS

OPTIONS

- Respect word order
- Ignore properties of the dominating node
- Split extended pos tags

XPath query generated from the input example. You can adapt it if necessary.

```
//node[@cat="np" and node[@rel="det" and @cat="np" and node[@rel="det" and @pt="lid"] and  
node[@rel="hd" and @pt="n"]] and node[@rel="hd" and @pt="n" and @postag="N(soort,mv,basis)"]]
```

XPath query generated from the input example. You can adapt it if necessary.

```
//node[@cat="np" and node[@rel="det" and @cat="np" and node[@rel="det" and @pt="lid"] and  
node[@rel="hd" and @pt="n"]] and node[@rel="hd" and @pt="n" and @getal="mv" and  
@graad="basis" and @ntype="soort"]]
```

SEARCH OPTIONS

OPTIONS

- Respect word order
- Ignore properties of the dominating node
- Split extended pos tags

XPath query generated from the input example. You can adapt it if necessary.

```
//node[@cat="np" and node[@rel="det" and @cat="np" and node[@rel="det" and @pt="lid"] and node[@rel="hd" and @pt="n"]] and node[@rel="hd" and @pt="n" and @postag="N(soort,mv,basis)"]]
```

XPath query generated from the input example. You can adapt it if necessary.

```
//node[@cat="np" and node[@rel="det" and @cat="np" and node[@rel="det" and @pt="lid"] and node[@rel="hd" and @pt="n"]] and node[@rel="hd" and @pt="n" and @getal="mv" and @graad="basis" and @ntype="soort"]]
```

SEARCH OPTIONS

CGN core corpus

Treebank		Contents	# Sentences	# Words	# Sentences	# Words	# Sentences	# Words
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>		NL		VL		TOTAL	
NL	VL							
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Spontaneous conversations ('face-to-face')	50,239	302,828	22,881	147,418	73,120	450,246
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Interviews with teachers of Dutch	2,484	25,724	4,289	34,158	6,773	59,882
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Telephone conversations (recorded via a switchboard)	11,649	70,084	3,142	19,984	14,791	90,068
<input type="checkbox"/>	<input checked="" type="checkbox"/>	Telephone conversations (recorded on MD)	0	0	929	6,309	929	6,309
<input checked="" type="checkbox"/>	<input type="checkbox"/>	Simulated business negotiations	3,123	25,524	0	0	3,123	25,524
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Interviews/discussions/debates (broadcast)	6,290	75,167	2,617	25,122	8,907	100,289
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	(Political) discussions/debates /meetings (non-broadcast)	1,166	25,125	543	9,009	1,709	34,134
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Lessons recorded in the classroom	3,064	26,004	1,395	10,116	4,459	36,120
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Live (sports) commentaries (broadcast)	2,251	25,002	1,026	10,147	3,277	35,149

Include context

SEARCH OPTIONS

Sentence ID

Matching sentences

fnc008003__58	weet ik eigenlijk niet wie dan 't meest kletst . <i>dan dan zijn jullie gewoon een stelletje mutsen die aan de telefoon hangen</i> . nou dat is helemaal niet waar .	[full screen] [XML]
fnc008005__117	wat moest jij nog in 't dorp vanmiddag ? <i>ik had vorige week had ik uh een paar batterijen gehaald</i> . want ik had batterijtjes nog 'ns een keer van die ouwe dingge uit één of andere koffiemolen of weet ik waar of ze in zaten maar die uh oplaadbare .	[full screen] [XML]
fnc008007__272	uh d'r zijn mensen die die die roken als als als ... <i>ik bedoel uh de Queen Mum die dronk tenminste nog een paar whisky's elke avond om vervolgens in een diep coma weg te dommelen</i> . ja ?	[full screen] [XML]
fnc008006__344	ik ga morgenvroeg wel effe naar 't post*a ... <i>ik heb hier nog een hele reut ouwe boeken staan</i> . of ouwe boeken .	[full screen] [XML]
fnc008006__77	ja dat was vandaag in dat waargebeurde maandagavonddrama ook . <i>was een kerel en die had z'n die die timmerde z'n vrouw l*a al uh structureel in elkaar een aantal jaren</i> . en toen op een gegeven moment ging ze bij d'r xxx weg en nieuw appartementje hele reutemeteut bijgestaan door d'r zuster .	[full screen] [XML]
fnc008007__71	hebben hebben andere kinderen 't gezien ? <i>ja d'r zijn toch een aantal kinderen die 't gezien hebben</i> . en dat is ook de reden waarom de leraren ge*a geattendeerd werden en uh ggg andere uh de de leraar die kwam ook aan die die zei nou ik ik rook 't meteen .	[full screen] [XML]

OUTLINE

- GrETEL in a nutshell
- GrETEL demo
 - Case study
 - Search options
- **Conclusions and future work**

CONCLUSIONS

- **GrETEL**: search engine for Dutch treebanks
- Input = natural language example
- Output = sample of similar sentences
- Syntactic concordancer
- Available online (via *Mozilla Firefox*)
- No installation required



FUTURE WORK

- **GrETEL 2.0**
 - Include SoNaR corpus (ca 500M tokens)
 - More generic

- **AfriBooms**
 - GrETEL for Afrikaans
 - Include other treebank formats



Try it yourself at

<http://nederbooms.ccl.kuleuven.be/eng/gretel>

Thanks for your attention!



KU LEUVEN