



NEDERBOOMS

Treebank Mining for Data-based Linguistics

Liesbeth Augustinus
Vincent Vandeghinste
Ineke Schuurman
Frank Van Eynde

NEDERBOOMS

- **Exploitation of Dutch treebanks for research in linguistics**
- CLARIN-VL project
- October, 2010 – February, 2012



NEDERBOOMS

- **Exploitation of Dutch treebanks for research in linguistics**
- CLARIN-VL project
- October, 2010 – February, 2012
- **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 treebank	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

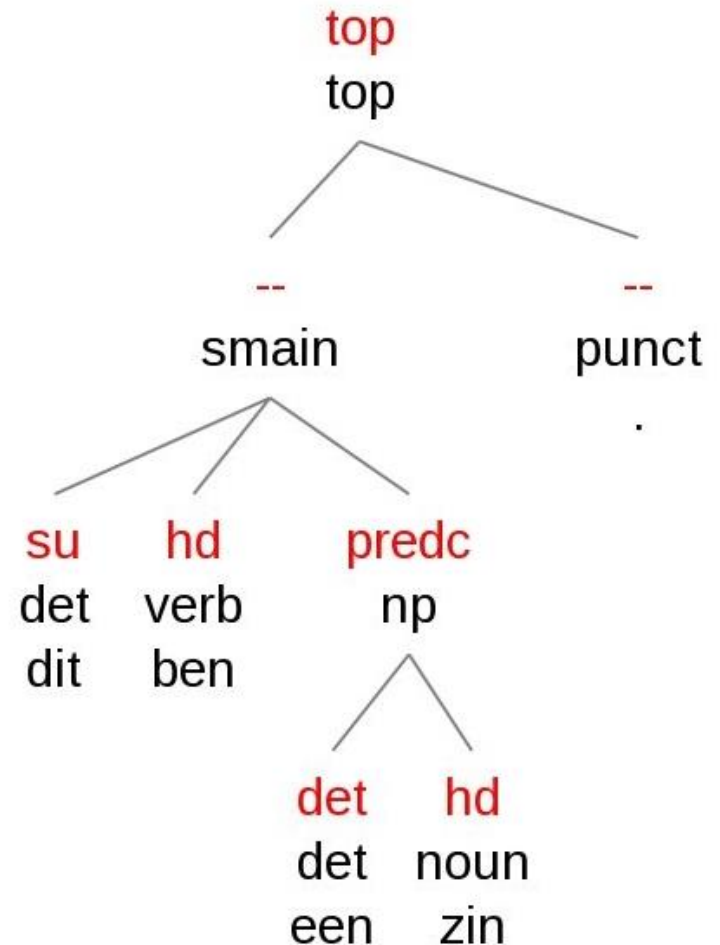
GrETEL

- **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

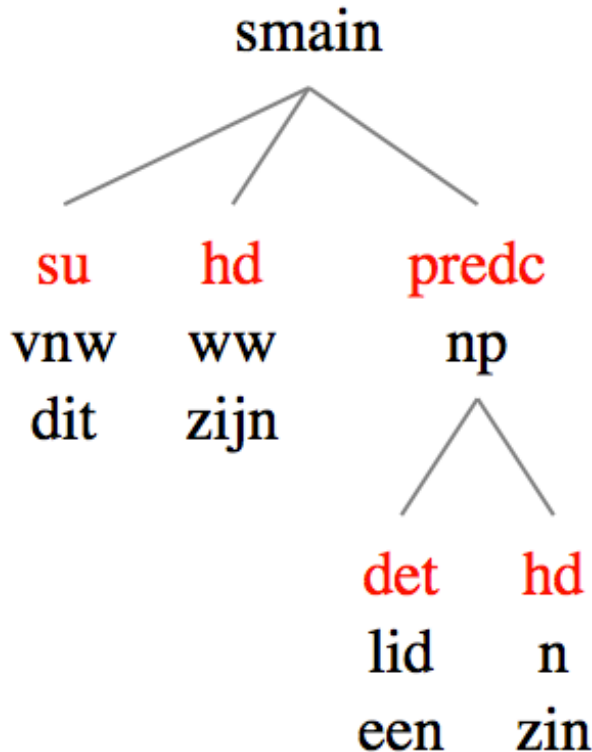
Dit is een zin. >> ALPINO parser >>
“This is a sentence.”



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**
 - ➔ No or limited knowledge of data structures and/or formal query languages needed





the user

1. Example sentence

2. Indicate relevant items
of the sentence

3. (Adapt XPath)
Select treebank

4. Inspect results



• 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

Number agreement in quantifying noun constructions

Een + noun + plural noun + singular/plural verb form

- *Een aantal treinen heeft vertraging.*
A number trains has-SG delay
- *Een aantal treinen hebben vertraging.*
A number trains have-PL delay

“A number of trains are running late.”

GrETEL ONLINE

Contact



Nederbooms

Home > Tools > GrETEL

GrETEL

What is *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 users to search for constructions which are similar to the example they provide, we hope to bridge the gap between traditional and computational linguistics.

- About
- Projects
- ▾ Tools
 - GrETEL
 - GrETEL for LASSY
 - GrETEL for CGN
 - Manual and docs
 - History

INPUT



Nederbooms

- About
- Projects
- ▽ Tools
 - ▽ GrE TEL
 - GrE TEL for LASSY
 - GrE TEL for CGN
 - Manual and docs
 - History

[Home](#) > [Tools](#) > [GrE TEL](#) > [GrE TEL for LASSY](#)

GrE TEL for LASSY (v1.2)

Please provide an **input example**

ANNOTATION MATRIX

Please indicate the relevant parts of the sentence. The syntactic properties of the relevant items are automatically included.

sentence		Een	aantal	treinen	heeft	vertraging	.
relevant nodes	pos	<input 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 checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>
	lemma	<input checked="" 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 type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>

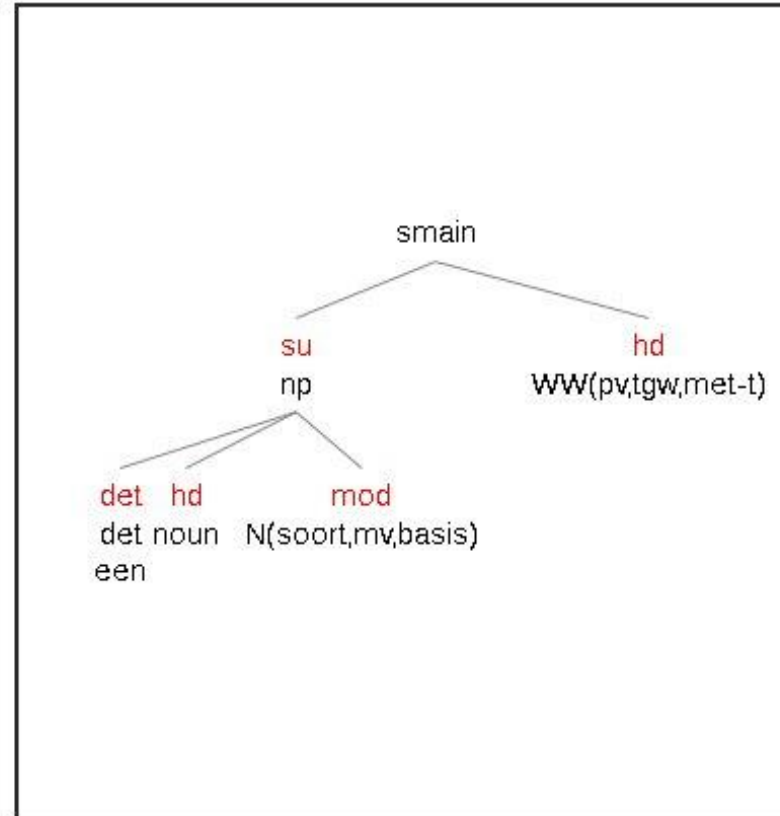
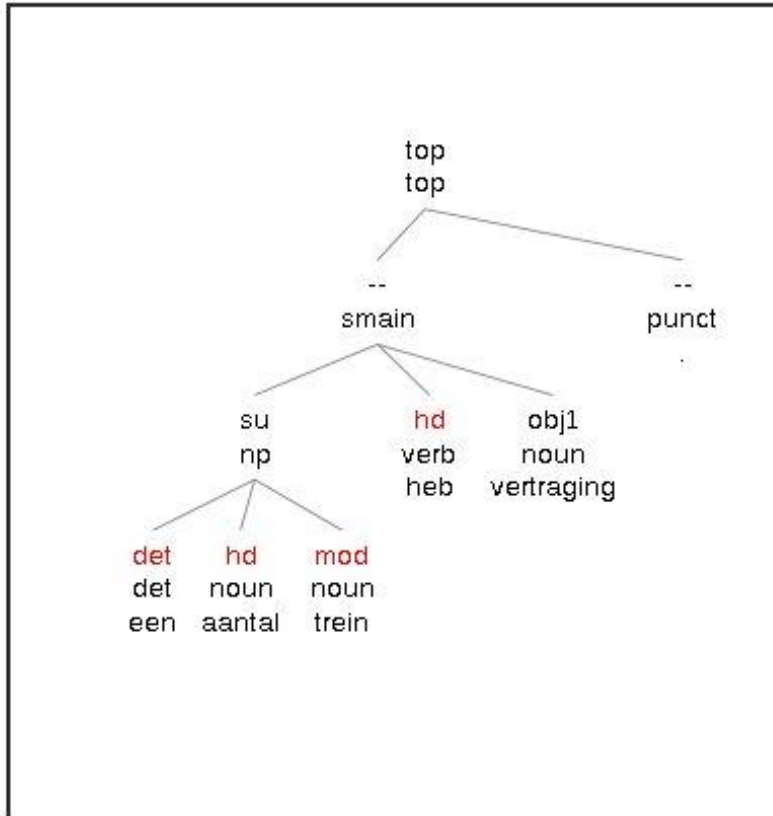
ANNOTATION GUIDELINES

GUIDELINES

- **pos**: Short part-of-speech tag (e.g. noun, verb, prep) [[list of all pos tags](#)]
- **extended pos**: Long part-of-speech tag. Besides the part-of-speech, the tags contain more detailed information on person, gender, number, tense, agreement etc. For example: N(soort,mv,basis), WWW(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). Note that the short part-of-speech tag will be taken into account as well.
- **token**: The exact word form. This is a case sensitive feature. Note that the short part-of-speech tag and the lemma will be taken into account as well.

Alpino parse of the input example [full screen]

Query tree [full screen]



XPATH GENERATOR

XPath query generated from the input example. You can adapt it if necessary. [\[download original XPath\]](#)

```
//node[@cat="smain" and node[@rel="su" and @cat="np" and node[@rel="det" and @root="een" and @pos="det"] and node[@rel="hd" and @pos="noun"] and node[@rel="mod" and @postag="N(soort,mv,basis)" and @pos="noun"]] and node[@rel="hd" and @postag="WW(pv,tgw,met-t)" and @pos="verb"]]
```

TREEBANK SELECTION

LASSY Small

<input checked="" type="checkbox"/>	Treebank	Contents	# Sentences	# Words
<input checked="" type="checkbox"/>	DPC	Dutch Parallel Corpus	11,716	193,029
<input checked="" type="checkbox"/>	Wikipedia	Dutch Wikipedia pages	7,341	83,360
<input checked="" type="checkbox"/>	WR-P-E	E-magazines, newsletters, teletext pages, web sites, Wikipedia	14,420	232,631
<input checked="" type="checkbox"/>	WR-P-P	Books, brochures, guides and manuals, legal texts, newspapers, periodicals and magazines, policy documents, proceedings, reports, surveys	17,691	281,424
<input checked="" type="checkbox"/>	WS-U	Auto cues, news scripts, text for the visually impaired	14,032	184,611
	LASSY Small	Complete treebank	65,200	975,055

RESULTS

Een + noun + plural noun + singular verb form

- *Een aantal treinen heeft vertraging.*
A number trains has-SG delay
→ **18 matches**

“A number of trains are running late.”

RESULTS: table

RESULTS: 18 matches in 18 sentences (out of 65,200 sentences)

[Download results\[`csv`\]](#)

[Show/hide detailed results](#)

TREEBANK	HITS	MATCHING SENTENCE	SENTENCES IN TREEBANK
DPC	1	1	11,716
Wikipedia	0	0	7,341
WR-P-E	5	5	14,420
WR-P-P	8	8	17,691
WS-U	4	4	14,032
TOTAL	18	18	65,200

RESULTS: data

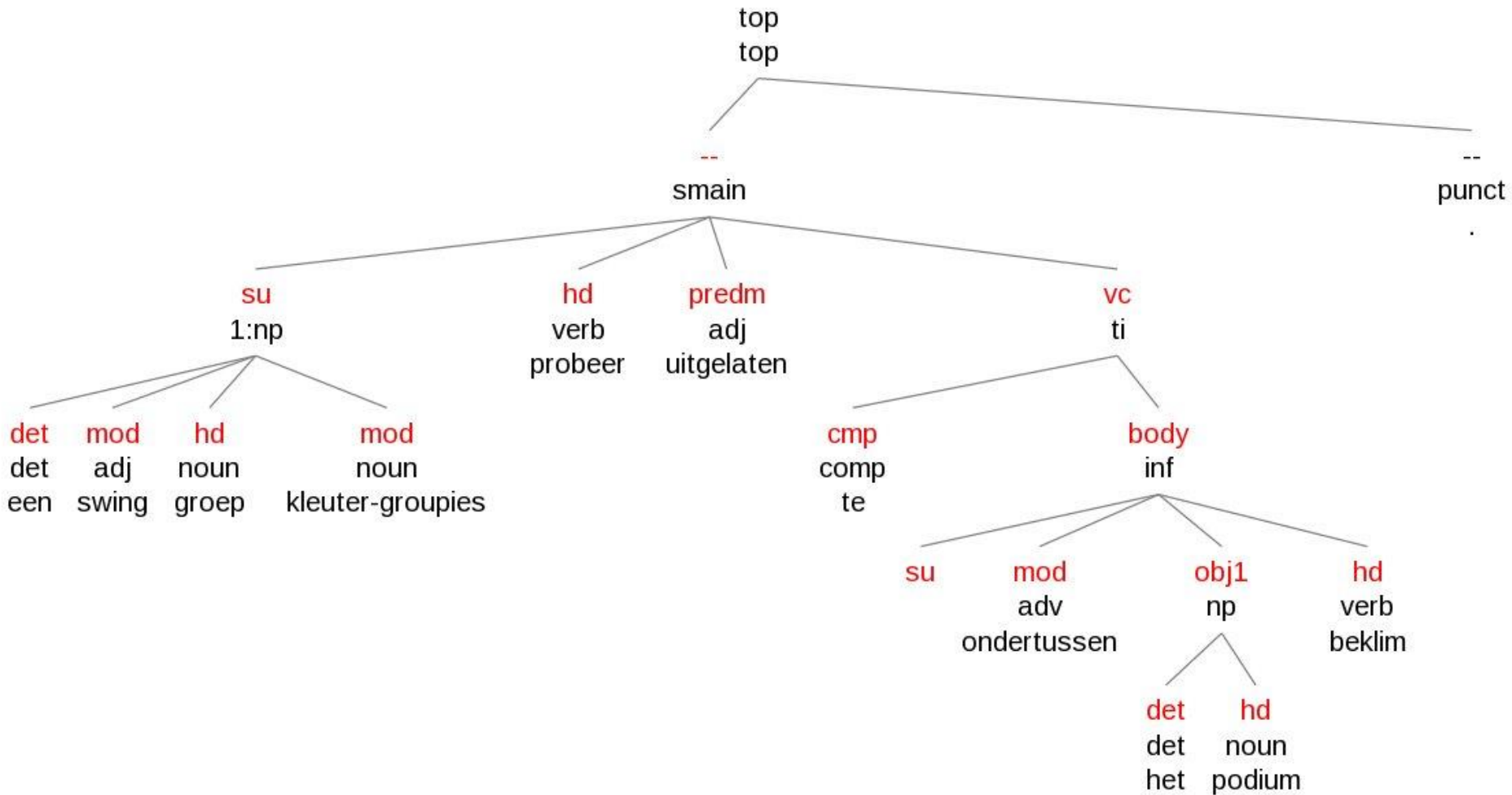
SENTENCE ID	MATCHING SENTENCE	HITS	DISPLAY OPTIONS
WR-P-E- I-0000020972.p.4.s.59.2	Een kolonie hommels sterft elk najaar , alleen de bevruchte jonge koninginnen blijven in leven en overwinteren alleen .	1	[full screen] [XML]
WR-P-P- H-0000000076.p.1.s.7	Een swingende groep kleuter-groupies probeert ondertussen uitgelaten het podium te beklimmen .	1	[full screen] [XML]
WR-P-P- I-0000000111.p.2.s.1	Een tiental commissies in het Amerikaanse Congres begint eveneens aan zijn onderzoek naar het failliet .	1	[full screen] [XML]
WS-U-E- A-0000000206.p.18.s.7	Eenmaal geland verslindt een wolk sprinkhanen in korte tijd de oogst van een heel seizoen .	1	[full screen] [XML]
WS-U-E- A-0000000011.p.24.s.1	Er dient zich een nieuwe generatie schaatsters aan , dat bewees Mark Tuitert afgelopen weekend in Heerenveen : want hij was na vier afstanden de sterkste : en is de nieuwe Europese Kampioen allround .	1	[full screen] [XML]
WR-P-E- I-0000039352.p.3.s.259	Gelukkig voor de talloze toeschouwers krijgt een aantal schutters weer snel een ritme te pakken .	1	[full screen] [XML]
WS-U-E- A-0000000215.p.25.s.2	In een voorstadje van Boedapest in Hongarije heeft zich een serie explosies voorgedaan in een vuurwerkfabriek .	1	[full screen] [XML]

RESULTS: data

SENTENCE ID	MATCHING SENTENCE	HITS	DISPLAY OPTIONS
WR-P-E- I-0000020972.p.4.s.59.2	Een kolonie hommels sterft elk najaar , alleen de bevruchte jonge koninginnen blijven in leven en overwinteren alleen .	1	[full screen] [XML]
WR-P-P- H-0000000076.p.1.s.7	Een swingende groep kleuter-groupies probeert ondertussen uitgelaten het podium te beklimmen .	1	[full screen] [XML]
WR-P-P- I-0000000111.p.2.s.1	Een tiental commissies in het Amerikaanse Congres begint eveneens aan zijn onderzoek naar het failliet .	1	[full screen] [XML]
WS-U-E- A-0000000206.p.18.s.7	Eenmaal geland verslindt een wolk sprinkhanen in korte tijd de oogst van een heel seizoen .	1	[full screen] [XML]
WS-U-E- A-0000000011.p.24.s.1	Er dient zich een nieuwe generatie schaatsters aan , dat bewees Mark Tuitert afgelopen weekend in Heerenveen : want hij was na vier afstanden de sterkste : en is de nieuwe Europese Kampioen allround .	1	[full screen] [XML]
WR-P-E- I-0000039352.p.3.s.259	Gelukkig voor de talloze toeschouwers krijgt een aantal schutters weer snel een ritme te pakken .	1	[full screen] [XML]
WS-U-E- A-0000000215.p.25.s.2	In een voorstadje van Boedapest in Hongarije heeft zich een serie explosies voorgedaan in een vuurwerkfabriek .	1	[full screen] [XML]

“greedy” search

Een swingende groep kleuter-groupies probeert ondertussen uitgelaten het podium te beklimmen .



RESULTS: trees

RESULTS

Een + noun + singular/plural verb form

- *Een aantal treinen heeft vertraging.*
A number trains has-SG delay
→ **18 matches**
- *Een aantal treinen hebben vertraging.*
A number trains have-PL delay
→ **30 matches**

“A number of trains are running late.”

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

PP-over-V

- V + PP
- ... *dat hij opstond met een kater*.
... that he woke-up with a hangover

- PP + V
- ... *dat hij met een kater opstond*.
... that he with a hangover woke-up

“... that he woke up with a hangover.”

SEARCH OPTIONS

PP-over-V

- V + PP
- ... *dat hij opstond met een kater*.
... that he woke-up with a hangover

WR-P-E-I-0000039352.p.3.s.39	De verscheidenheid qua bordjes is groot , zo loopt bijvoorbeeld Schutterij 't Zandakker Gilde Sint Jan Venray met een prachtig roodkoperen geslagen draagbord , terwijl andere verenigingen met een houten draagbord lopen .	1	[full screen] [XML]
dpc-vla-001171-nl-sen.p.29.s.6	Het was in het Kaaitheater dat Anne Teresa De Keersmaeker een podium vond in het begin van haar carrière .	1	[full screen] [XML]
dpc-eli-000943-nl-sen.p.24.s.1	Het Comité voor geneesmiddelen voor menselijk gebruik (CHMP) heeft geconcludeerd dat Ariclaim een bescheiden effect vertoonde in studies maar dat dit effect een voordeel kan zijn voor vrouwen met matige tot ernstige stress-urine-incontinentie , daar er momenteel geen alternatieve farmacologische behandeling voor deze aandoening bestaat .	1	[full screen] [XML]
dpc-vla-001161-nl-sen.p.54.s.1	Een goed jaar na de start van het nieuwe financieringsinstrument kunnen we besluiten dat ARKimedes na een zwakke start nu toch wel goed op gang komt .	1	[full screen] [XML]

2,895 matches in 2,769 sentences

But: results include PP + V as well!

SEARCH OPTIONS

PP-over-V

- V + PP + **word order option**
- ... *dat hij opstond met een kater*
... that he woke-up with a hangover

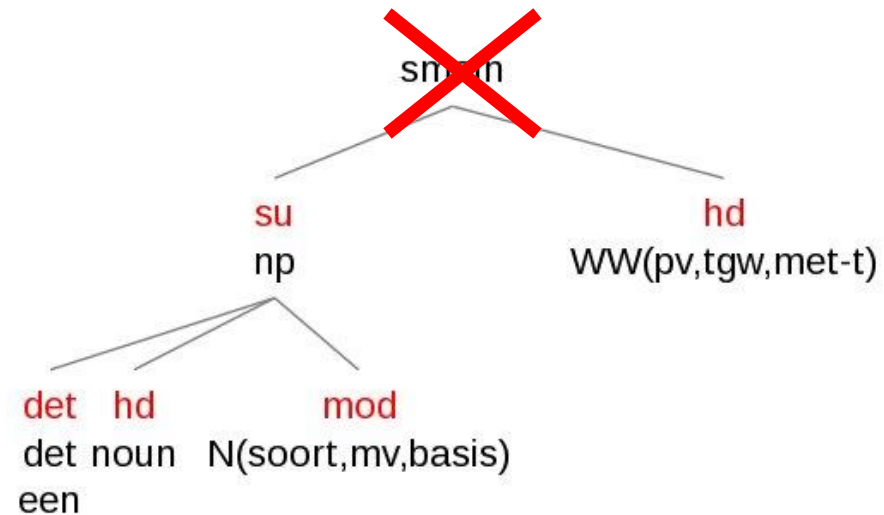
WS-U-E-A-0000000218.p.21.s.3	Het ging destijds eigenlijk al meteen mis : het voorprogramma met ondermeer André van Duin en de Fouryo's was veel te rustig voor een publiek dat kwam voor rok-en-rol .	1	[full screen] [XML]
dpc-med-000686-nl-sen.p.26.s.3	Het belangrijkste voordeel van het product is de associatie van een lage dosis ethyniloestradiol met een progestativum van de tweede generatie (cfr de polemie van het einde van vorige eeuw over het thrombo-embolisch risico dat lager is met tweede generatie progestativum in vergelijking met deze van de derde generatie .)	1	[full screen] [XML]
WR-P-E-I-0000024561.p.4.s.68.3	Brels eerste album op het Barclay-label bestaat voor een groot deel uit studio-versies van materiaal dat luisteraars al kenden van het Olympia-concert .	1	[full screen] [XML]
dpc-rou-000983-nl-sen.p.8.s.11	Het verhaal van Rodenbach , dat net als wijn rijpt op eikenhouten vaten , past daar perfect bij , maar was nog niet verteld . "	1	[full screen] [XML]

789 matches in 777 sentences
Results only include V + PP

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. [[download original XPath](#)]

```
//node[@cat="sm n"] and node[@rel="su" and @cat="np" and node[@rel="det" and @root="een" and @pos="det"] and node[@rel="hd" and @pos="noun"] and node[@rel="mod" and @postag="N(soort,mv,basis)" and @pos="noun"]] and node[@rel="hd" and @postag="WW(pv,tgw,met-t)" and @pos="verb"]
```

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. [[download original XPath](#)]

```
//node[@cat="smain" and node[@rel="su" and @cat="np" and node[@rel="det" and @rest="con" and @pos="det"] and node[@rel="hd" and @pos="noun"] and node[@rel="mod" and @postag="N(soort,mv,basis)" and @pos="noun"]] and node[@rel="hd" and @postag="WW(pv,tgw,met-t)" and @pos="verb"]]
```

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. [[download original XPath](#)]

```
//node[@cat="smain" and node[@rel="su" and @cat="np" and node[@rel="det" and @root="een" and @pos="det"] and node[@rel="hd" and @pos="noun"] and node[@rel="mod" and @postag="N(soort,mv,basis)" and @pos="noun"]] and node[@rel="hd" and @postag="WW(pv,tgw,met-t)" and @pos="verb"]]
```

XPath query generated from the input example. You can adapt it if necessary. [[download original XPath](#)]

```
//node[@cat="smain" and node[@rel="su" and @cat="np" and node[@rel="det" and @root="een" and @pos="det"] and node[@rel="hd" and @pos="noun"] and node[@rel="mod" and @getal="mv" and @graad="basis" and @pos="noun" and @ntype="soort"]] and node[@rel="hd" and @pvtijd="tgw" and @wvform="pv" and @pos="verb" and @pvagr="met-t"]]
```

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. [[download original XPath](#)]

```
//node[@cat="smain" and node[@rel="su" and @cat="np" and node[@rel="det" and @root="een" and @pos="det"] and node[@rel="hd" and @pos="noun"] and node[@rel="mod" and @postag="N(soort,mv,basis)" and @pos="noun"]] and node[@rel="hd" and @postag="WW(pv,tgw,met-t)" and @pos="verb"]]
```

XPath query generated from the input example. You can adapt it if necessary. [[download original XPath](#)]

```
//node[@cat="smain" and node[@rel="su" and @cat="np" and node[@rel="det" and @root="een" and @pos="det"] and node[@rel="hd" and @pos="noun"]] and node[@rel="mod" and @getal="mv" and @graad="basis" and @pos="noun" and @ntype="soort"]] and node[@rel="hd" and @pvtijd="tgw" and @wvform="pv" and @pos="verb" and @pvagr="met-t"]]
```

SEARCH OPTIONS

LASSY Small

<input checked="" type="checkbox"/>	Treebank	Contents	# Sentences	# Words
<input checked="" type="checkbox"/>	DPC	Dutch Parallel Corpus	11,716	193,029
<input checked="" type="checkbox"/>	Wikipedia	Dutch Wikipedia pages	7,341	83,360
<input checked="" type="checkbox"/>	WR-P-E	E-magazines, newsletters, teletext pages, web sites, Wikipedia	14,420	232,631
<input checked="" type="checkbox"/>	WR-P-P	Books, brochures, guides and manuals, legal texts, newspapers, periodicals and magazines, policy documents, proceedings, reports, surveys	17,691	281,424
<input checked="" type="checkbox"/>	WS-U	Auto cues, news scripts, text for the visually impaired	14,032	184,611
	LASSY Small	Complete treebank	65,200	975,055

OPTION

Include context (one sentence before and after the matching sentence)

SEARCH OPTIONS

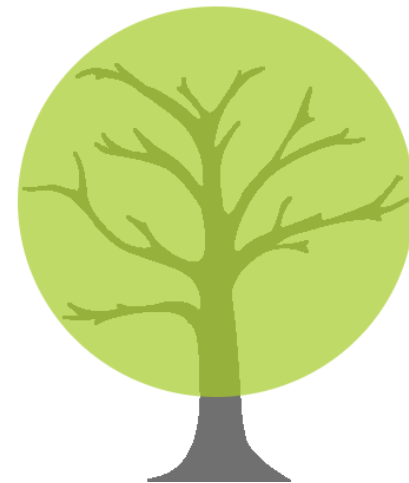
SENTENCE ID	MATCHING SENTENCE	HITS	DISPLAY OPTIONS
WR-P-E- I-0000020972.p.4.s.59.2	Levenscyclus Een kolonie hommels sterft elk najaar , alleen de bevruchte jonge koninginnen blijven in leven en overwinteren <i>alleen</i> .	1	[full screen] [XML]
WR-P-P- H-0000000076.p.1.s.7	Tien meter lang . ' Een swingende groep kleuter-groupies probeert ondertussen uitgelaten het podium te beklimmen . Het gaat hier niet om lieve luistermuziek maar om stevige kleuterpop .	1	[full screen] [XML]
WR-P-P- I-0000000111.p.2.s.1	Een tiental commissies in het Amerikaanse Congres begint eveneens aan zijn onderzoek naar het failliet . Een van de commissies wil informatie over het energiebeleid krijgen van de Amerikaanse vice-president Dick Cheney .	1	[full screen] [XML]
WS-U-E- A-0000000206.p.18.s.7	Een sprinkhaan kan per dag zijn eigen gewicht aan voedsel opeten . Eenmaal geland verslindt een wolk sprinkhanen in korte tijd de oogst van een heel seizoen . Alleen in Mali werden de afgelopen dagen al 42 springhanenzwermen waargenomen .	1	[full screen] [XML]
WS-U-E- A-0000000011.p.24.s.1	Er dient zich een nieuwe generatie schaatsters aan , dat bewees Mark Tuitert afgelopen weekend in Heerenveen : want hij was na vier afstanden de sterkste : en is de nieuwe Europese Kampioen allround .	1	[full screen] [XML]
WR-P-E- I-0000039352.p.3.s.259	De eerste ronde valt het zwaarste , getuige het feit dat meestal meer dan de helft van de deelnemers dan reeds ' sneuvelt ' . Gelukkig voor de talloze toeschouwers krijgt een aantal schutters weer snel een ritme te pakken .	1	[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**
 - Query very large treebanks
 - Improve user interface
 - Include SoNaR corpus (ca 500M tokens, 42M sentences)

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



Try it yourself at

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

Questions? Comments? Suggestions?

Contact us at gretel@ccl.kuleuven.be

Thanks for your attention!



KU LEUVEN