

# **Cornetto and WordNet: network of relations based on the words of a language**

Isa Maks

e.maks@vu.nl

Computational Lexicology & Terminology Lab

Faculty of Arts

VU University Amsterdam

# Cornetto

- STEVIN 2005-2008

Combination of two existing lexical resources  
(Dutch WordNet and Reference Resource Dutch)

- CLARIN 2012

Standardized to LMF (lexicon markup framework)

ISO data categories

demo

# Cornetto

A lexical resource for Dutch which combines

- Syntagmatic information
- Lexical units with detailed information on specific word meanings ranging from form, morphology, syntax, case frames, combinatorics, semantics to pragmatics
- Paradigmatic information
- Synsets built from synonymous lexical units and related through semantic relations

about words and their senses

# Overview of the presentation

- Cornetto
  - Lexical Units
  - Synsets and wordnets
  - Links to other resources
- Applications
- Availability

# syntagmatic vs. paradigmatic

<b>John</b>	<b>bought</b> <b>purchased</b> <b>got</b> <b>acquired</b>	<b>a(n)</b>	<b>expensive</b> <b>costly</b> <b>pricey</b>	<b>old</b> <b>new</b>	<b>bicycle</b> <b>motorcycle</b> <b>car</b>
-------------	--	-------------	--	--------------------------	---

A syntagmatic relationship refers to the relationship a word has with other words that surround it: they are about the lexical company the word keeps.

A paradigmatic relationship refers to the relationship between words that are the same parts-of-speech and which can be substituted for each other in the same position within a given sentence: these words can be substituted.

Together they constitute the identity of an item

# Lexical Units

- Words are alphabetically ordered
- And divided in senses or lexical units

## opscheppen

opscheppen-v-1 toe | tō |

- Definition: **noun**  
1 any of the five digits at the end of the human foot: *he cut his big toe on a sharp stone.*

opscheppen-v-2

- Definition: **verb**  
2 the lower end, tip, or point of something, in particular:
  - the tip of the head of a golf club, furthest from the shaft.
  - the foot or base of a cliff, slope, or embankment.
  - a flattish portion at the foot of an otherwise steep curve on a graph.
  - a section of a rhizome or similar fleshy root from which a new plant may be propagated.

# Lexical Units

- Lemma and word forms
- Morphology and morpho-syntax
- Syntax, syntactic behaviour and syntactic subcategorisation
- Sense and semantics
- Combinatorics

# Lemma and wordforms

Plural forms : huis , huizen; (house , houses)

Comparative and superlative forms : groot, groter, grootst ; ( large, larger, largest)

Past tense and past participle : deed , gedaan  
(did done)

Related forms (spelling and form variants):  
bureau vs. buro organisation vs. organization



# Morphology and morpho- syntax

Separability (verbs) : hij **wast** de borden **af** ;  
he cleans the dishes

Gender (nouns): m, f, n

Adverbial usage (adjectives): slecht (bad,  
badly)

# Syntax and syntactic behaviour

## Complementation patterns

ofclause de vraag of ... (the question whether)

prep ... (opscheppen over , boast about)

## syntactic subcategorization frame

hij scheidt op over haar

PP (constituent) as prepositional object (function)

hij scheidt (het eten) op

NP as direct object

# Sense and semantics

Short definition or resume

Semantic type

Artefact een fabriek bouwen (to build a factory)

Place achter de fabriek wonen (to live behind the factory)

Institution de fabriek is gesloten (the factory is closed)

# Syntax and syntactic behaviour

## Complementation patterns

ofclause de vraag of ... (the question whether)

prep ... (opscheppen over , boast about)

## syntactic subcategorization frame

hij scheidt op over haar

PP (constituent) as prepositional object (function)

hij scheidt (het eten) op

NP as direct object

# Combinatorics

Illustrative examples

*We saw five monkeys in the zoo*

Lexical collocations

*Wash the dishes*

Grammatical collocations

*In terms of*

Idioms

*Kick the bucket*

# DEMO

# Paradigmatic relations

WordNets and Synsets

# Relational model of meaning

animal

kitten

man

boy

cat

tail

dog

puppy

woman

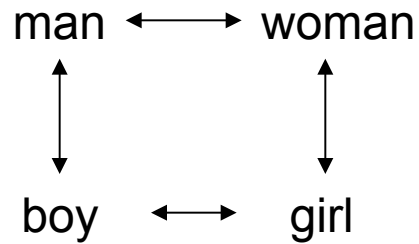




# Relational model of meaning

animal

kitten



cat

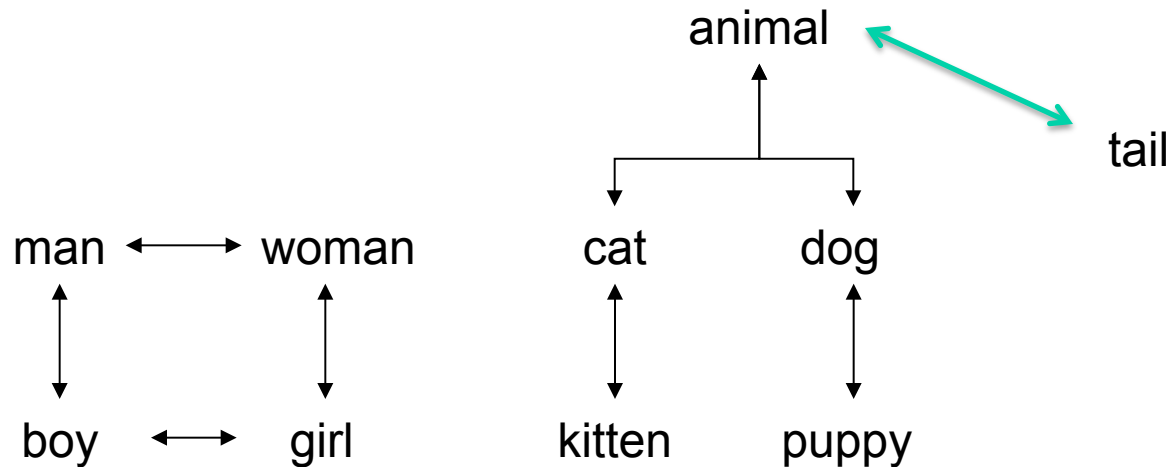
dog

puppy

tail



# Paradigmatic information: Relational model of meaning



# concepts

## Knowledge about concepts

Minimal knowledge about concepts:

- X is a kind of Y
- X has part Y
- an X Ys
- X is Y/has property Y

# Most important relation: synonymy

WordNet groups (roughly) synonymous, denotationally equivalent, words into unordered sets of synonyms (synsets)

{hit, beat}

{slaan, meppen}

{close, shut}

{sluiten, dichtdoen}

{big, large}

{groot, fors}

{queue, line}

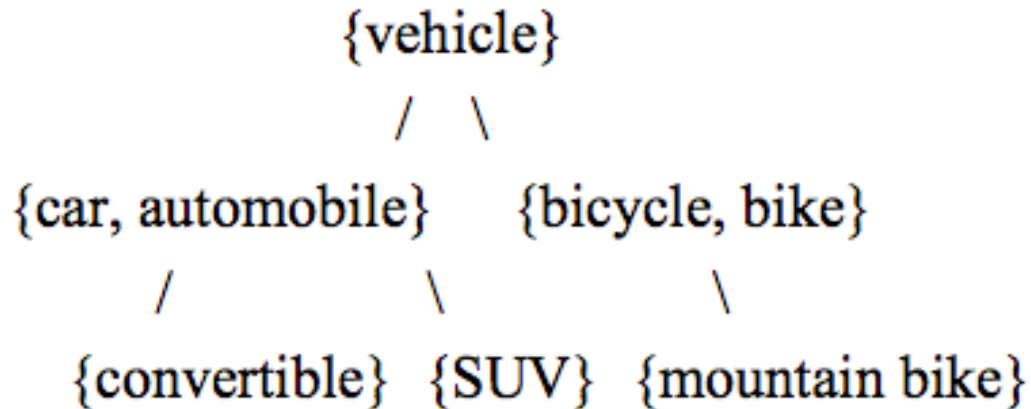
{rij, wachtrij}

Each synset expresses a distinct concept

Most synsets contain  $> 1$  member

# HYPONOMY/ HYPERONOMY

Denotes more/less general concepts:



“A car is is a kind of vehicle”

“The class of vehicles contains cars, SUVs, mountain bikes”

# Another relation: antonymy

“clang” association between pairs of antonymous words

{big-little}

{groot-klein}

{wet-dry}

{droog-nat, vochtig}

{big, large}

{groot, fors}

{sluiten-openen}

{dood-leven}

# Conceptual-semantic relations

Synsets, the nodes of the network are interrelated via conceptual-semantic relations

- Traditional lexicographer's relations
- Relations from psycholinguistic evidence (word association norms) : the mental lexicon

# Meronymy among noun synsets

Meronymy/holonymy (part/whole)

{car, automobile}

|

{engine}

/ \

{spark plug} {cylinder}

“An engine has spark plugs”

“Spark plus and cylinders are parts of an engine”



# Role relations

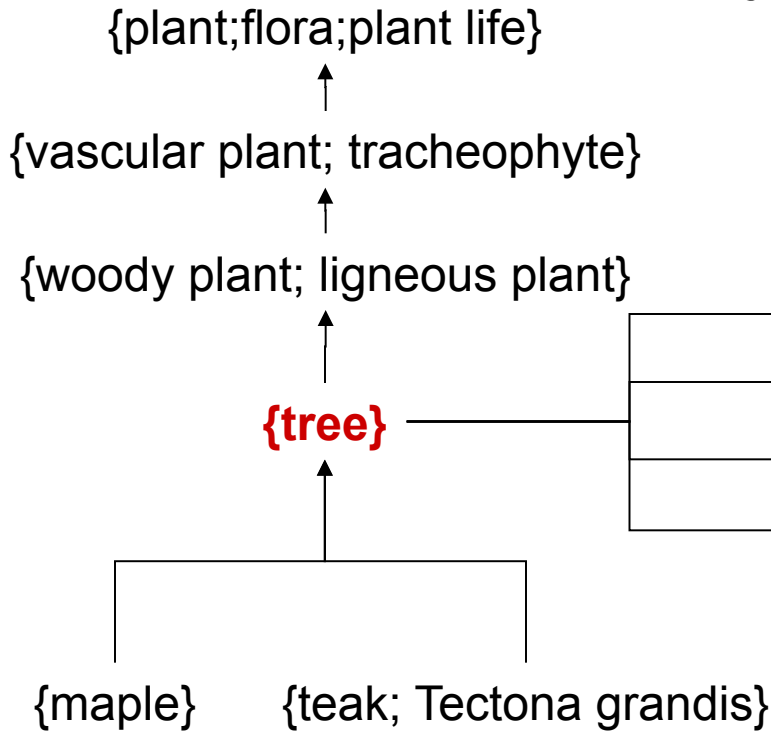
{knife}	ROLE_INSTRUMENT	{to cut}	
{to cut}	INVOLVED_INSTRUMENT	{knife}	<i>reversed</i>
{school}	ROLE_LOCATION	{to teach}	
{to teach}	INVOLVED_LOCATION	{school}	<i>reversed</i>

Useful when hyponymy does not clarify the position of the concept network, but the word is still closely related to another word.

# Wordnet terminology

hypernym/hyperonym

↑  
H  
Y  
P  
O  
N  
Y  
M  
Y  
↓



holonym

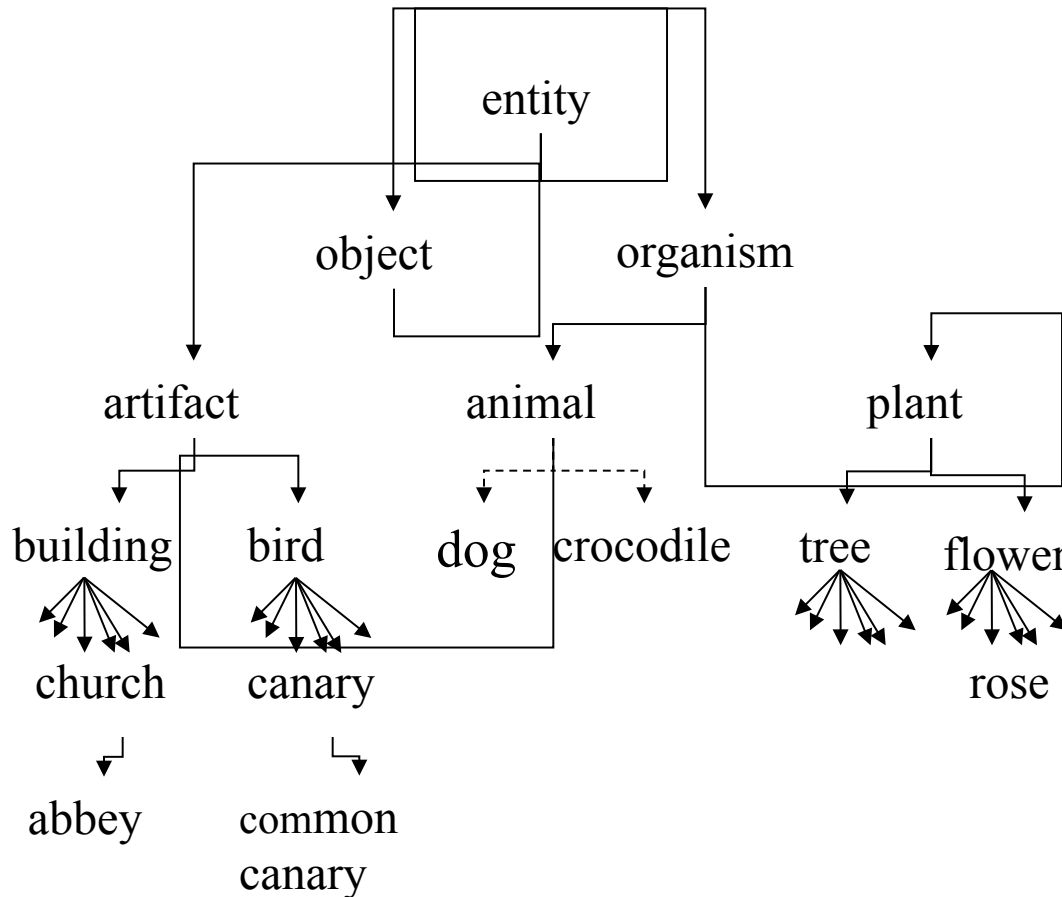
meronym

← MERONYMY →

hyponym



# Lexicalization patterns



← top-layer

← 25 unique beginners

← Basic Level Concepts



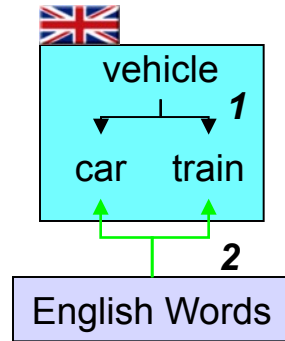
# DEMO

# What kind of resource is the WordNet?

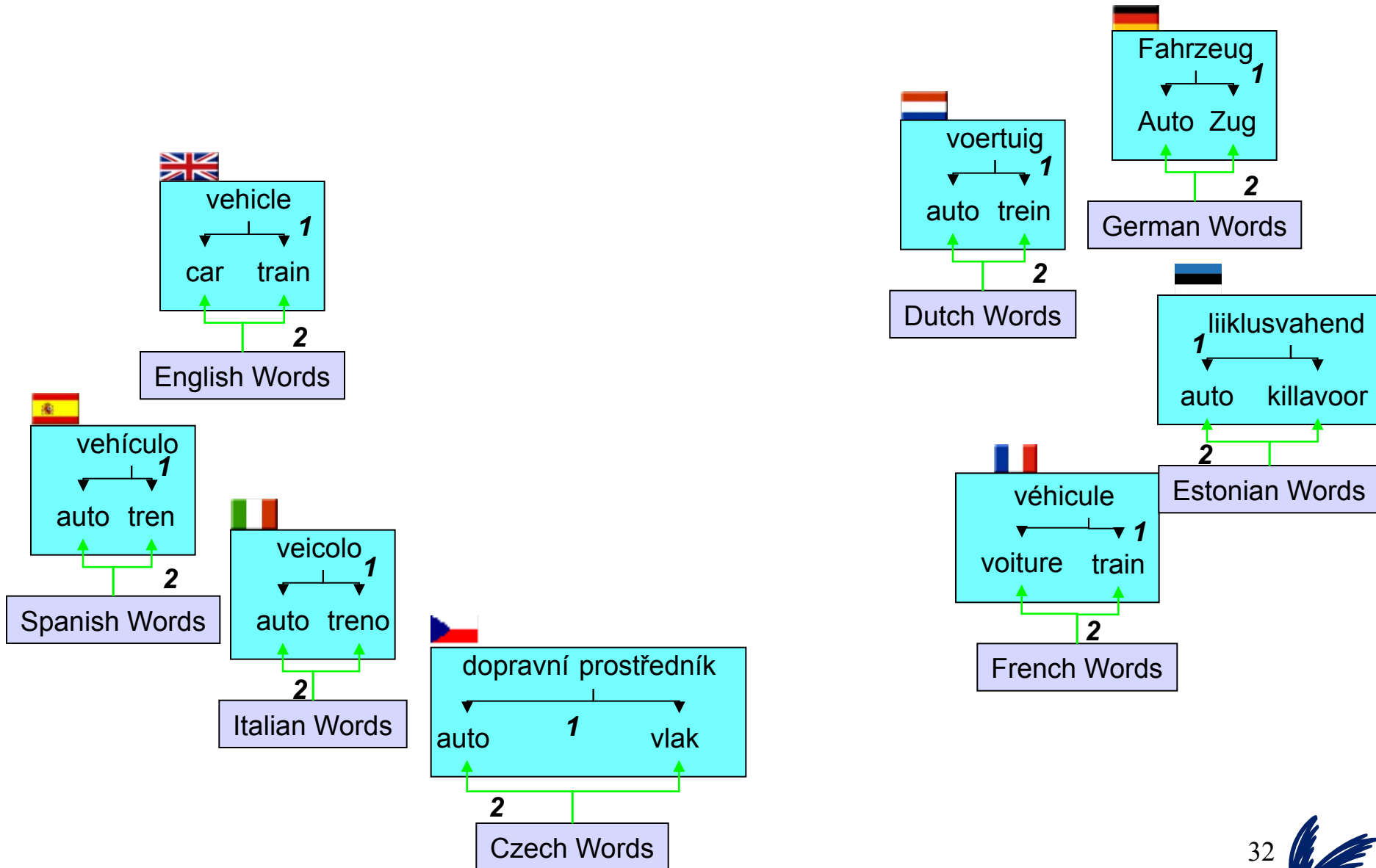
- <http://wordnet.princeton.edu/>
- Developed by George Miller and his team at Princeton University, as the implementation of a mental model of the lexicon
- Mostly used database in language technology
- Enormous impact in language technology development
- Large
- Free and downloadable
- English
- Links to many other knowledge bases

# Links to other repositories

# Euro WordNet Model

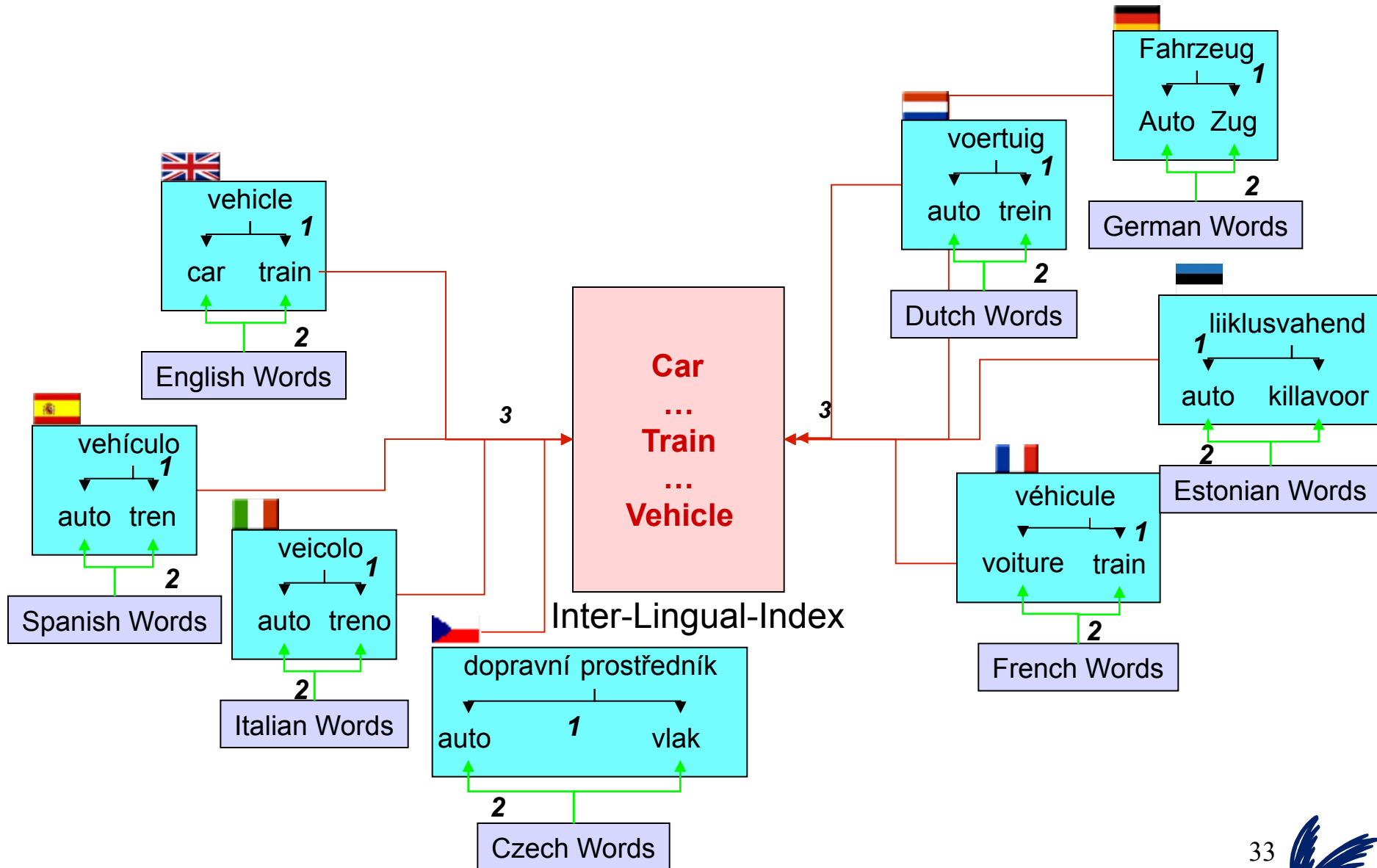


# EuroWordNet Model





# Euro WordNet Model



# Typical gaps in the English wordnet

- Dutch:

*doodschoppen* (to kick to death):

eq\_hyperonym {kill} V and to {kick} V

*aardig* (Adjective, to like)

eq\_near\_synonym {like} V

*cassière* (female cashier)

eq\_hyperonym {cashier} and to {woman}

*kunstproduct* (artifact substance)

eq\_hyperonym {artifact} and to {product}

- Spanish:

*alevín* (young fish):

eq\_hyperonym {fish} and eq\_be\_in\_state {young}

*cajera* (female cashier)

eq\_hyperonym {cashier} and to {woman}



# Differences among wordnets

## English Wordnet

## Dutch Wordnet

large number of synsets	62 synonyms
asshole, bastard, cocksucker, dickhead, shit, mother fucker, motherfucker, prick, whoreson, son of a bitch, SOB	naarling:1/r_n-24518, beroerling:1/d_n-26921, ellendeling:1/r_n-12324, etterbak:1/d_n-75936, etterbuil:2/d_n-75940, fielt:1/d_n-80137, fluum:2/d_n-81948, gemenerik:1/r_n-14607, hond:2/r_n-79023, hondenlul:1/r_n-17019, kankerlijer:1/d_n-130709, kelerelijder:1/d_n-540923, kelerelijer:1/d_n-147148, klerelijer:1/r_n-19790, kloot:1/r_n-19887, kloothommel:1/d_n-137246, klootspiraal:1/d_n-412711, klootzak:1/r_n-19888, kwal:2/r_n-21077, lamgat:1/d_n-152244, lammeling:1/r_n-21272, lamstraat:1/d_n-152396, lamzak:1/r_n-21286, lazersteen:1/d_n-413025, lazerstraat:1/d_n-154087, loeder:1/r_n-22410, lul:2/r_n-22757, lulhannes:1/d_n-161976, lulletje:1/d_n-541138, miesgasser:1/d_n-172163, mispunt:1/r_n-24006, onverlaat:1/r_n-26320, paardelul:1/d_n-228940, paardenlul:1/n_n-501022, patjakker:1/d_n-212558, pleurislijder:1/r_n-28842, ploert:1/r_n-28881, plurk:1/d_n-220067, etc. etc.
cad, bounder, blackguard, dog, hound, heel	
gasbag, windbag	
rotter, rat, skunk, stinker, bum, puke, crumb, lowlife, scum_bag, so-and-so	
pain, pain_in_the_neck, nuisance	
worm, louse, insect, dirt_ball	

*insulting terms for people who are stupid, ridiculous, irritating, lazy, slow, .....<sup>35</sup>*



# Wordnet versus ontology

- **Wordnet:**
  - autonomous **language-specific** and **complex** lexicalization pattern in a relational network.
  - **Usage:** to predict **substitution** in text for information retrieval, text understanding, text generation, machine translation.
- **Ontologies:**
  - data structure with formally defined concepts.
  - **Usage:** making semantic inferences.

# Ontology and lexicon

- Hierarchy of disjunct types:

Canine => PoodleDog; NewfoundlandDog;  
GermanShepherdDog; Husky

- Lexicon:

- NAMES for TYPES:

{**poodle**}EN, {**poedel**}NL, {**pudoru**}JP

=> ((instance x Poodle)

- LABELS for ROLES:

{**watchdog**}EN, {**waakhond**}NL, {**banken**}JP

=>((instance x Canine) and (role x GuardingProcess))



# Wordnet to ontology mappings

{teacher} Noun, English

-> sc\_domainOf **human**

-> sc\_playRole **done-by**

-> sc\_participantOf **teach**

{leraar} Noun, Dutch // lit. *male teacher*

-> sc\_domainOf **man**

-> sc\_playRole **done-by**

-> sc\_participantOf **teach**

{lerares} Noun, Dutch // lit. *female teacher*

-> sc\_domainOf **woman**

-> sc\_playRole **done-by**

-> sc\_participantOf **teach**

<http://www.ontologyportal.org>

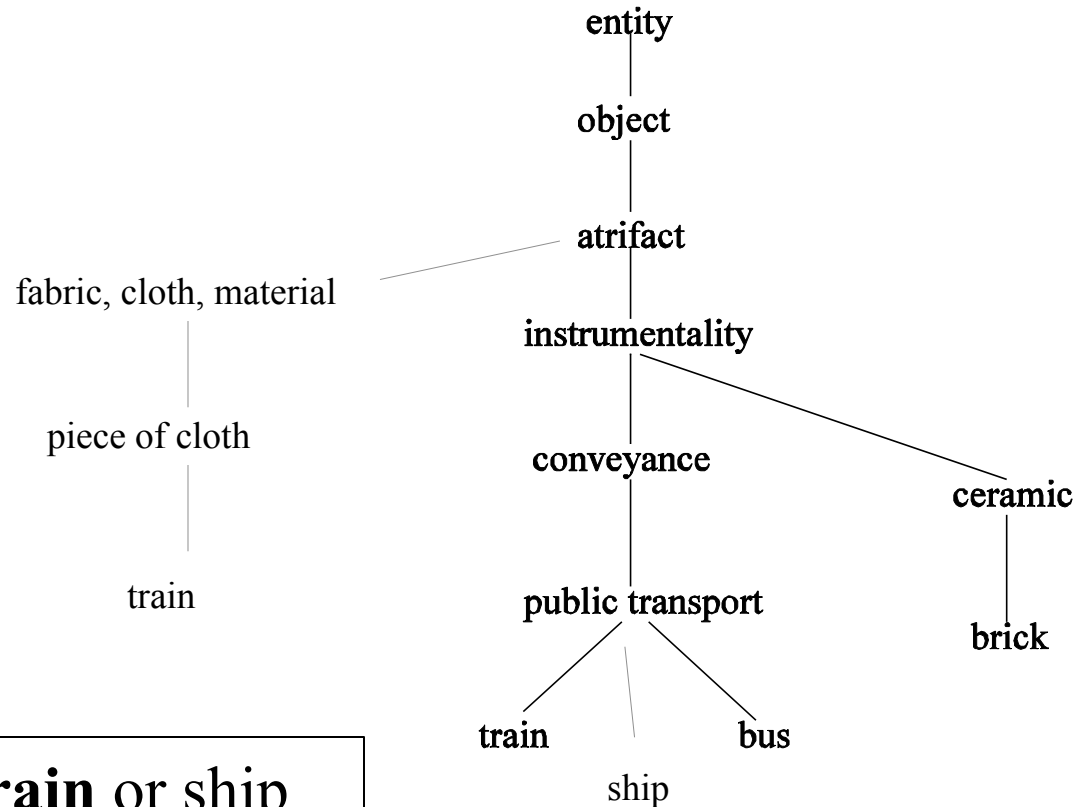


# Usage of Cornetto

Similarity

Sentiment Analysis

# Similarity/ Relatedness



It can be delivered by **train** or ship

## **train**

1. public transport provided by a line of railway cars coupled together and drawn by a locomotive
2. piece of cloth forming the long back section of a gown that is drawn along the floor





# Generating Sentiment Lexicons

Subjective vs. objective

a burglar alarm (objective)

He cried out in alarm (subjective)

He disapproved of drinking (subjective)

Positive vs. negative

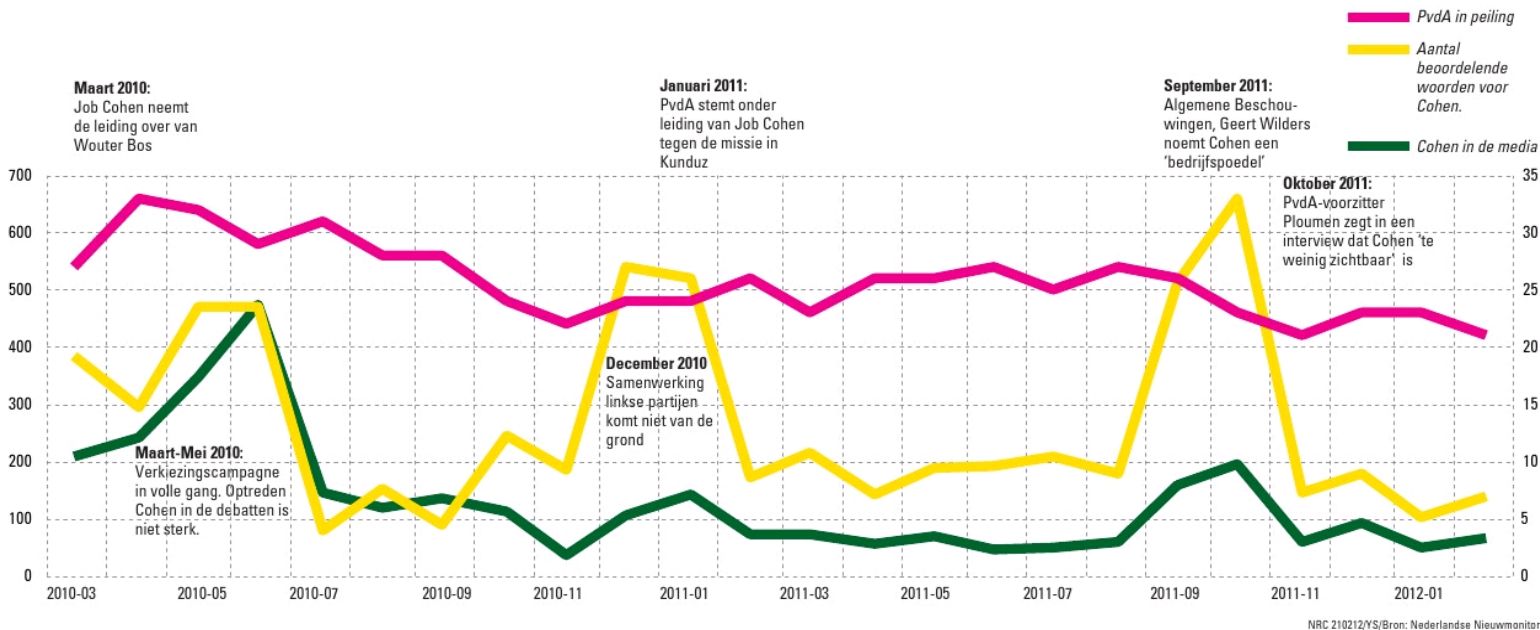
terrible (negative)

fantastic (positive)

review	ur	ann1	ann2
Hotel bezocht ivm jaarlijks weekendje weg met vrienden. <b>Gezellig</b> hotel. Te beginnen met de naam. Hotelbar aanwezig. <b>Persoonlijke</b> en <b>vriendelijke</b> benadering. Wij sliepen op de twee zolderkamers. <b>Niets mis mee</b> . Hotel ligt aan een <b>levendig</b> pleintje met kroegen en restaurants etc. <b>Goede</b> herinneringen. Ontbijt <b>uitstekend</b> .	3	4	4
Ga er zeker niet heen zolang het hotel in renovatie is. De receptie schijnt al klaar te zijn, het is net alsof je aan de incheckbalie bij Schiphol staat en bovendien chagrijnige receptionisten. We wilden een drankje drinken en geen bar te vinden behalve een afschuwelijke ongezellige ruimte. De ober zei dat we beter naar de overkant naar het Holiday Inn konden gaan om even gezellig te zitten. We hebben het ontbijt ook maar gelijk overgeslagen. Prijs/kwaliteit een 2	-4	-4	-5
Wat een <b>drama</b> hotel, echt vreselijk! <b>Kakkerlakken</b> in de kamer, niet een maar wel 6 en dat in december! 1 handdoekje (schuurpapier) voor 2 personen. Douche delen met 300 andere gasten. Lijkt wel een <b>gevangenis</b> . douche <b>niet schoon</b> , ook al vraag je er om. <b>wc overgelopen, maandverband overal. belachelijke</b> prijzen voor dit <b>crematorium</b> . <b>GEEN AANRADER!</b> P.s. niet via sun air en d reizen boeken. Van onze klachten en <b>beloofd</b> geld nooit iets gezien!	-5	-5	-5
Een Toilet in het midden van de kamer. Een douche zonder cabine. Een ontbijt vol vliegen. Weinig keuze. Lawaai buiten. Geen informatie en lange wachtrijen. Goed bed en netjes maar vanwege allerlei slordigheden, ongeduldig personeel, zeker niet voor herhaling vatbaar	-4	-5	-5
Zoals de andere bezoekers al zeggen. Het hotel is echt prima. Het eten is heerlijk nen op de sevice en de vriendelijkheid is ook niets aan te merken. Wacht alleen nog een paar jaar voor dit deel van Dubai. Het is een enorme bouwput en de sfeer van 1001-nacht is ver te zoeken.	4	-2	3
<b>Goed</b> Hotel, midden in Down Town. Echt een zaken hotel. voor vakantie zou ik iets anders uitzoeken.	4	-2	-2
Het receptie personeel moet cliënten met U aanspreken en niet met 'je'	4	-2	-2



## COHEN WEET STEEDS MINDER GOED DE MEDIA-AANDACHT OP TE EISEN



### JUIST OP DE BELANGRIJKE MEDIAMOMENTEN LIET COHEN HET AFWETEN

- Was het werkelijk de 'mediawerkelijkheid' die Job Cohen mede de das omdeed, zoals hij gisteren zelf verklaarde? Wetenschappelijk instituut voor de journalistiek De Nederlandse Nieuwsmonitor onderzocht dit samen met de afdeling Communicatiewetenschap van de Vrije Universiteit voor *nrc.next*.

- Cohen heeft moeite met de nieuwe realiteit, zo blijkt. Inmiddels wordt de manier waarop nieuws wordt gebracht allang niet meer door de inhoud, maar door de vorm bepaald. Met de televisie als meest bepalende medium. En precies daarmee had Cohen moeite. Het in één minuut duidelijk maken

van een standpunt tijdens een televisiedebat. Het meelachen met Rutger Castricum. Wilders aanpakken als het er echt op aankomt.

- De onderzoekers hebben een inhoudsanalyse uitgevoerd van alle artikelen in *de Volkskrant*, *Trouw*, *de Telegraaf*, *ADen NRC Handelsblad* tussen het aantreden van Cohen op 12 maart 2010 en gisteren. Er is gekeken naar de hoeveelheid media-aandacht voor Cohen, maar ook naar de aard van deze berichtgeving. Deze 'sentimentanalyse' kan op basis van trefwoorden bepalen of een artikel positief, neutraal of negatief van aard is. De betrouwbaarheid van deze analysevorm ligt volgens De

Nieuwsmonitor tussen de 60 en 70 procent.

- Als hij in maart 2010 leider van de PvdA wordt, spreken de krantenartikelen vol vertrouwen over Cohen, blijkt uit het onderzoek. Cohen is de gedoodverfde winnaar van de verkiezingen. Dit vertrouwen zakt snel weg. In het begin van de campagne heeft Cohen het zwaar in de debatten. Enkele maanden na zijn aantreden is de teneur van de artikelen al een stuk negatiever.

- Daarna zal het niet meer goed komen met de PvdA-leider. Cohen laat het juist op de belangrijke mediamomenten afweten, zo is te zien in

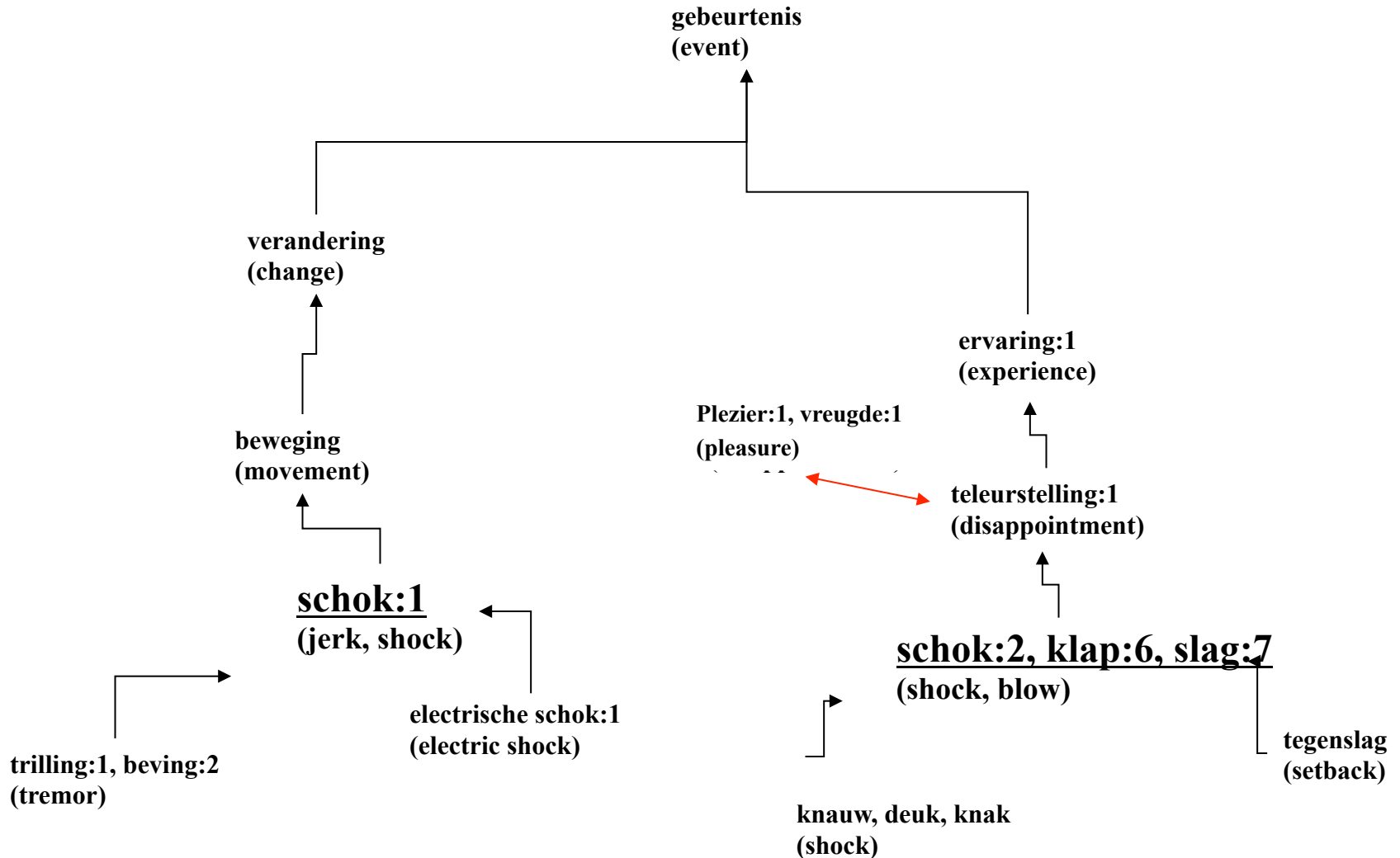
de grafiek hierboven. De pieken in de beoordeelingslijn laten zien wanneer het er bij Cohen echt om ging: de momenten dat hij het meest werd 'beoordeeld' in de media. En voor Cohen waren dat nu juist elke keer momenten waarop hij de meeste kritiek kreeg: de debatten, Wilders die hem een bedrijfspoezel noemde en PvdA-voorzitter Ploumen die Cohen een gebrek aan zichtbaarheid verweet.

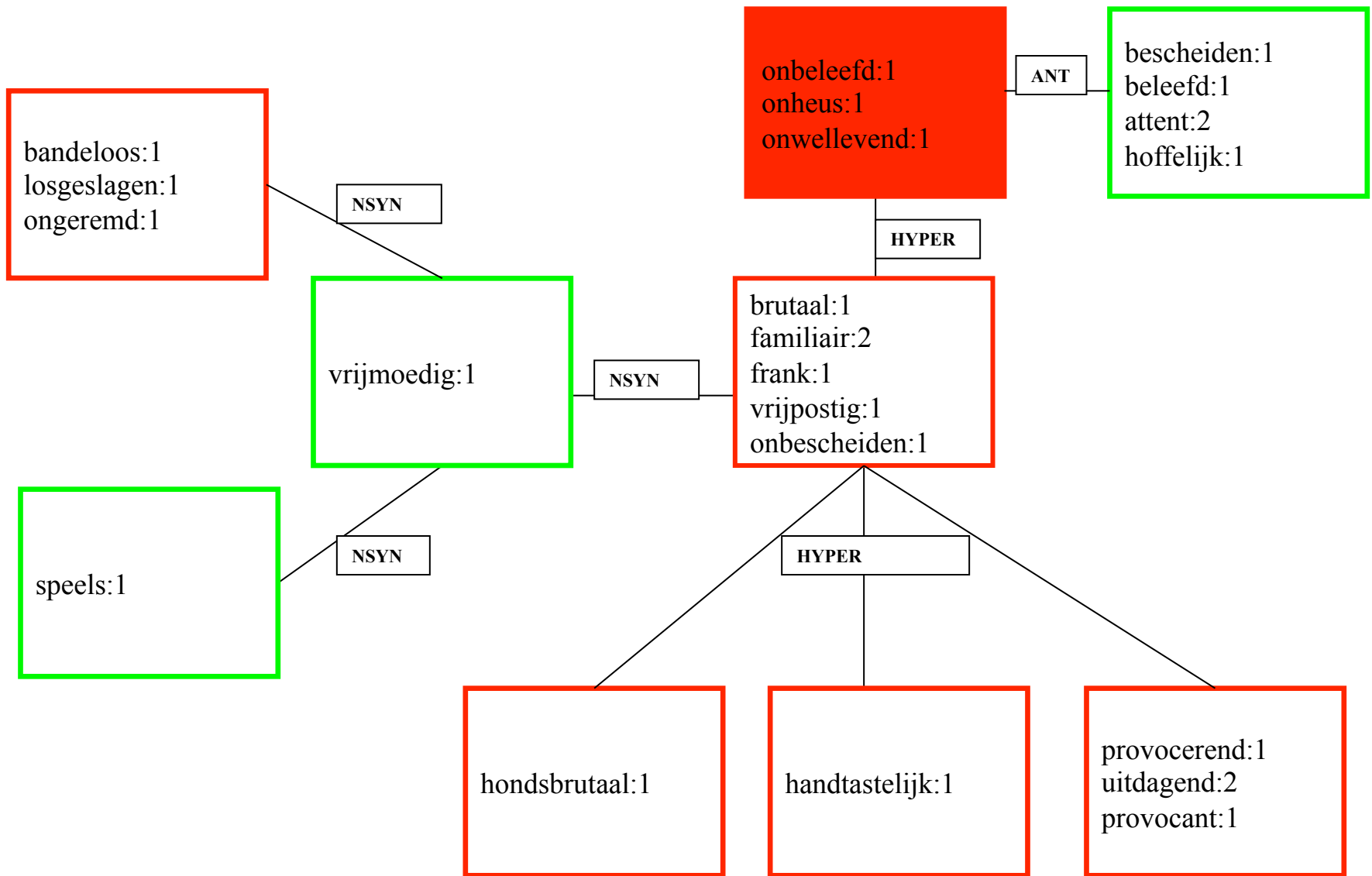
STIJN BRONZWAER  
Redacteur Media

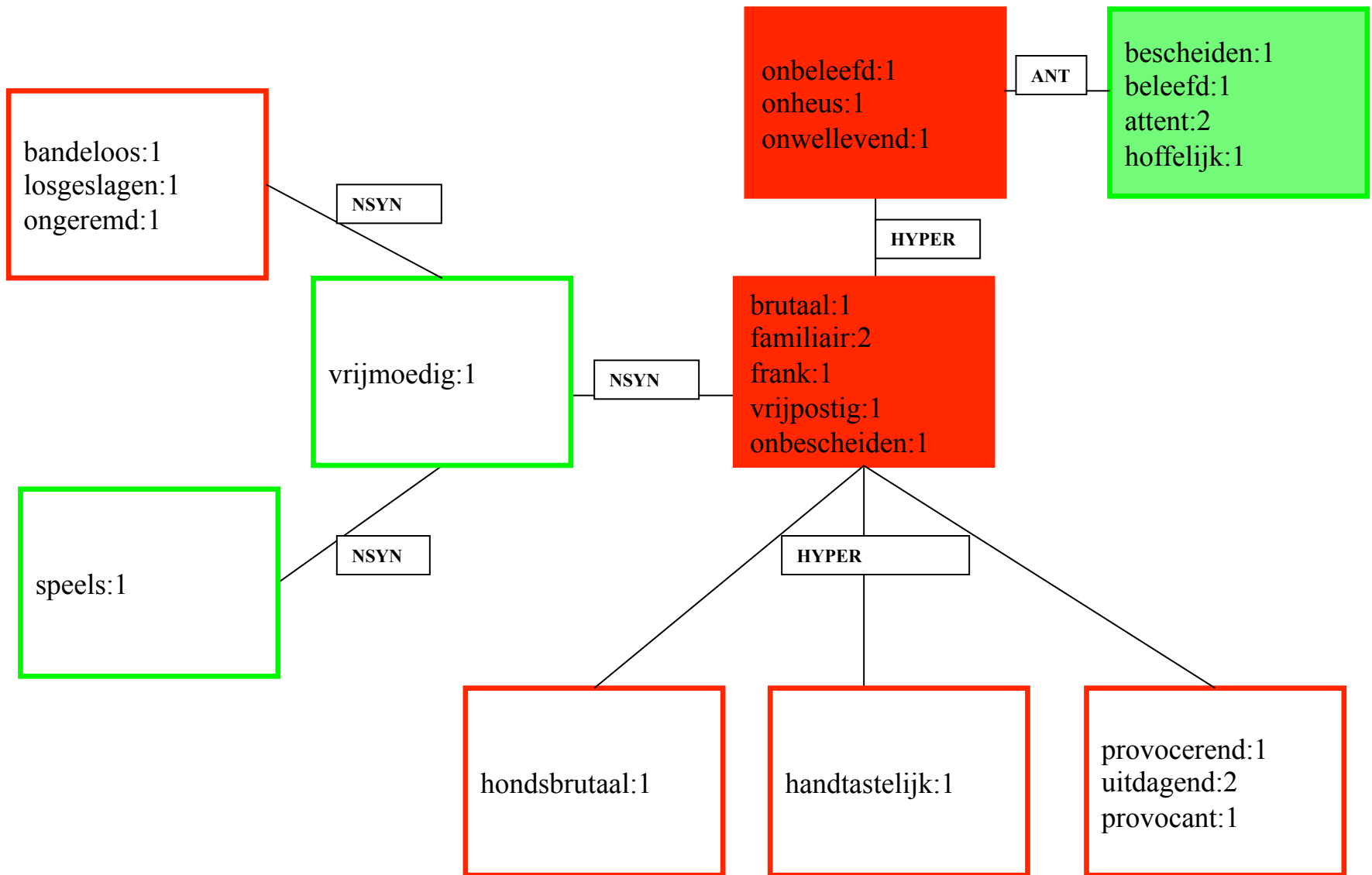
▶ Lees het hele onderzoek op [www.nieuwsmonitor.net](http://www.nieuwsmonitor.net)

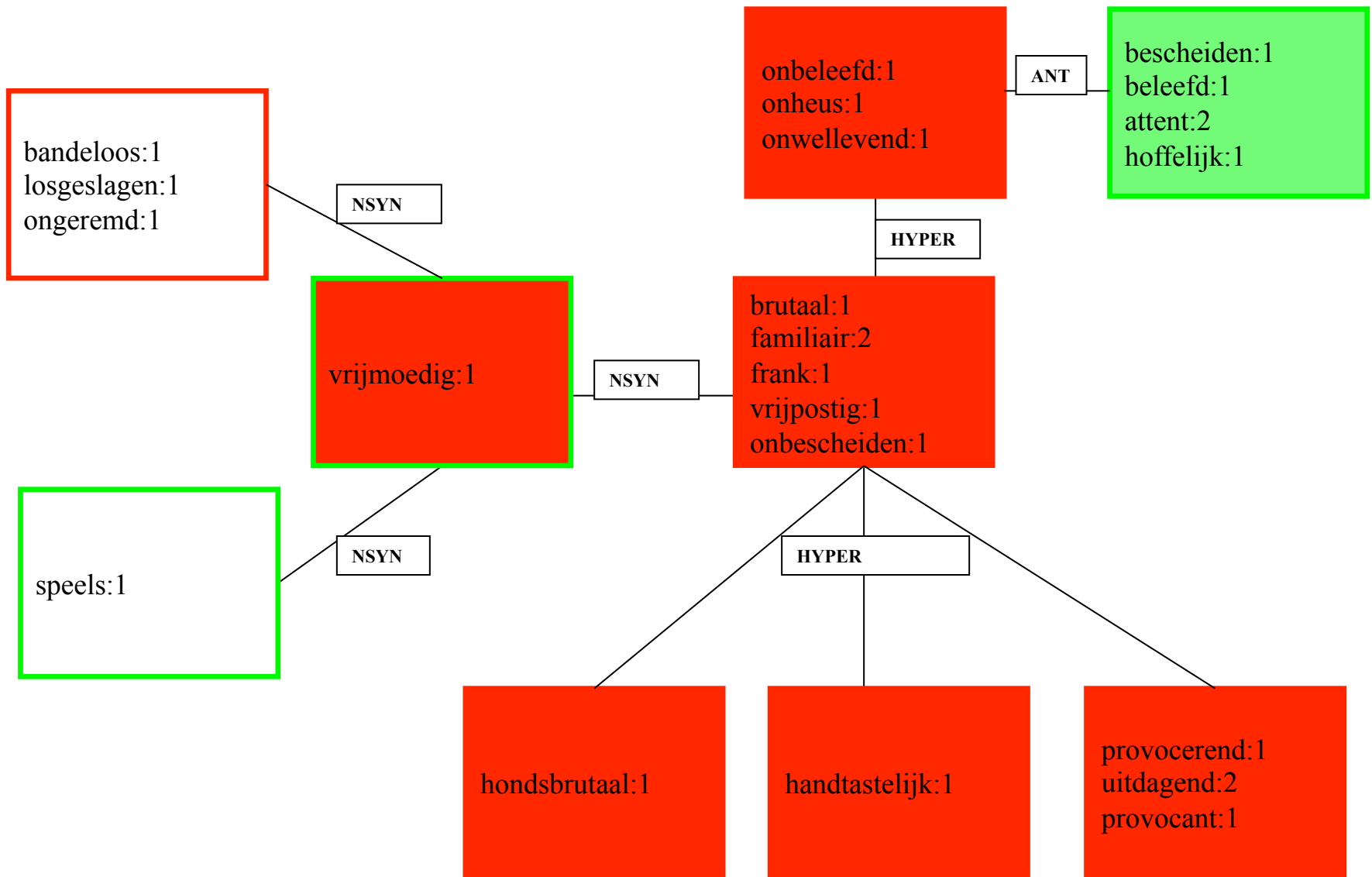


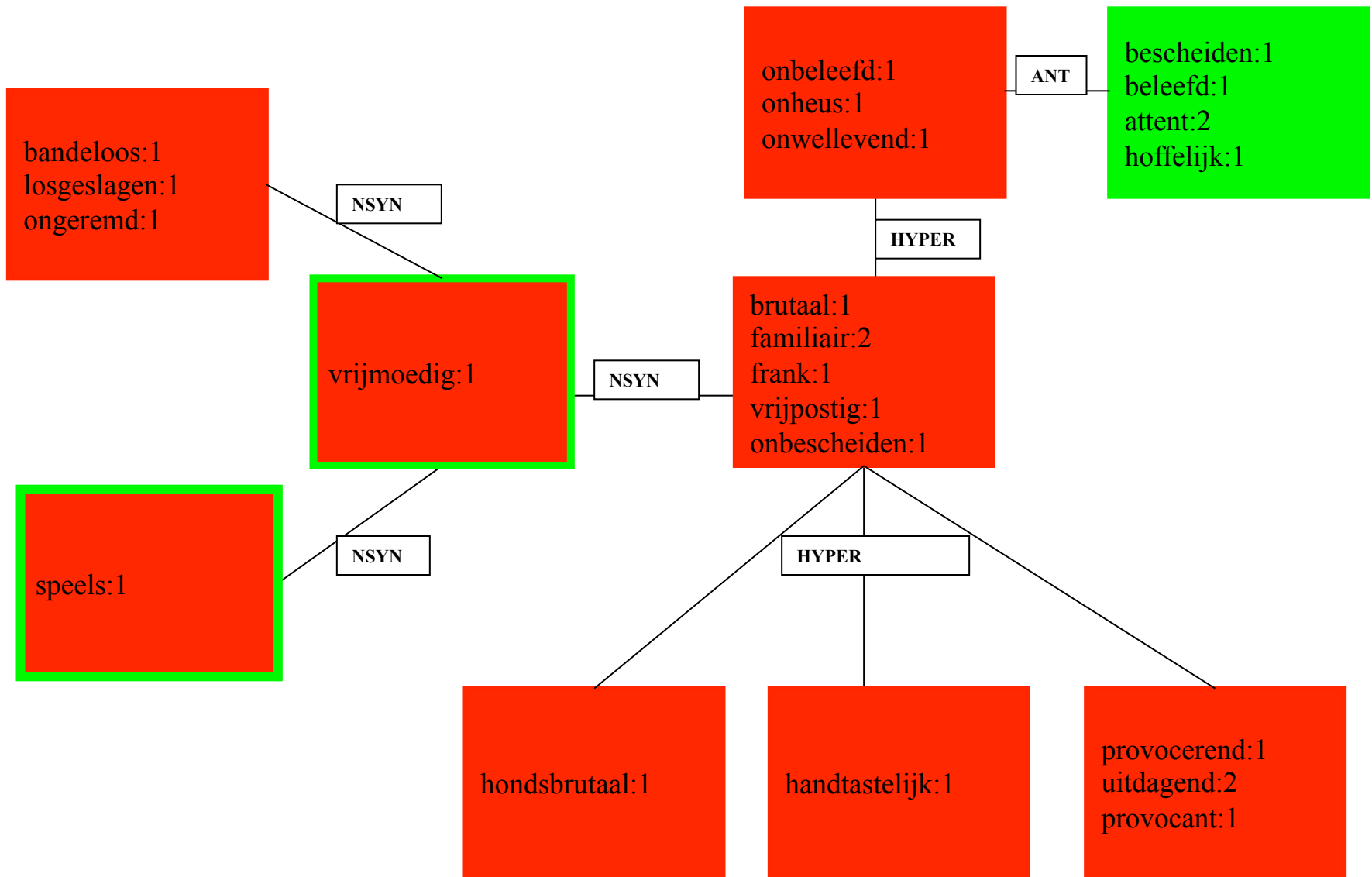
# Method I : Wordnet propagation Dutch Wordnet Cornetto













# Method II :

## machine learning with Cornetto features

### Lexical unit features

#### Gradability

*awful, more awful, most awful*

*Left hand, \*lifter hand, \*leftest hand*

#### Attributive vs. predicative use

*The chemical lab vs. \*the lab is chemical*

*Complementation patterns*

#### Syntactic Complementation (lack of), valency

*more syntax = less subjectivity and zero syntax = more subjectivity.*

*Hij zeurt. (he nags) vs. \*hij zegt. (he says)*

#### Semantic type

*Emotion, cognitive\_verb, (subjective words) Action, artifact (objective)*

#### Selectional restrictions (nouns and verbs)

*He annoyed me (subject:person/ object:person) vs. he pushed the car (object :non-human)*

#### Definition overlap

*Gloss overlap (Lesk): definitions of subjective words have more in common with each other than definitions of objective words*



# Method II : machine learning with Cornetto features

## Synset features

WN similarity: subjective concepts are more similar and/or more related to each other than to objective concepts

Synset size: synsets with subjective concepts have more synonyms on average than other synsets.

WordNet domain features

SUMO features

*Method II uses information from both the lexical unit and the synset part of Cornetto to generate a sentiment lexicon*

# Results

## Propagation (positive vs. negative)

goed/a/positive (good)

zeuren/v/negative (nag)

## Machine learning (objective vs. subjective)

goed/a/subjective

zeuren/v/subjective

voetballen/v/objective (play football)

		ML	Prop	ML&prop
ML	n	0.77	0.57	
ML	a	0.76	0.77	
ML	v	0.69	0.65	

## 2.2 Statistics

The current version of Cornetto (version 2.1) includes the following items:

Lexical entries	127,334
Noun Lexical entries	85,368
Verb Lexical entries	16,502
Adjective Lexical entries	15,458
Multiword Lexical Entries	9,397
Lexical Entries with polarity labels	22,202
Sense group relations	11,812
Sense examples	80,512
Lexical collocations	19,166
Grammatical collocations	10,373
Synsets	70,497
Synset relations (internal relations between Dutch synsets)	91,734
Synset equivalent relations (between Dutch and English Princeton WordNet synsets)	84,031
References to SUMO (from synsets to SUMO)	69,610
References to WordNet domains (from synsets to WordNet domains)	93,165

Table 1: statistics of Lexical Resource Cornetto (version 2.1, May 2013)

# Availability Cornetto

Lexical Resource CORNETTO		
demo	<a href="http://cornetto.inl.nl/">http://cornetto.inl.nl/</a>	
demo editor	<a href="https://debvisdic.let.vu.nl:9002/">https://debvisdic.let.vu.nl:9002/</a>	username: gast password: gast
download	<a href="http://tst-centrale.org/nl/producten/lexica/cornetto/">http://tst-centrale.org/nl/producten/lexica/cornetto/</a> Xml-LMF or RDF	Free for research (with a license)