

Starten met NLP

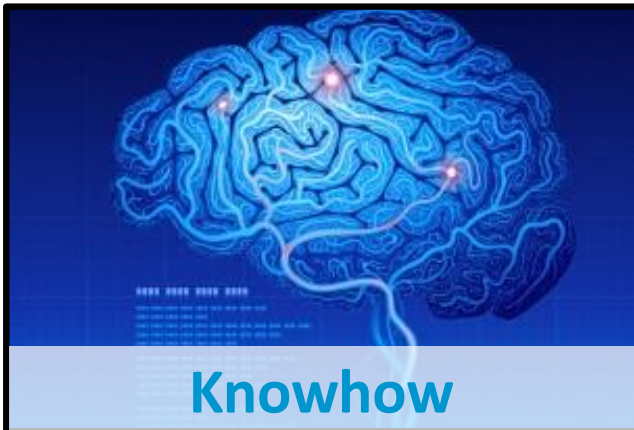
in het Nederlands

Joachim Ganseman
Smals Research

30/03/2021



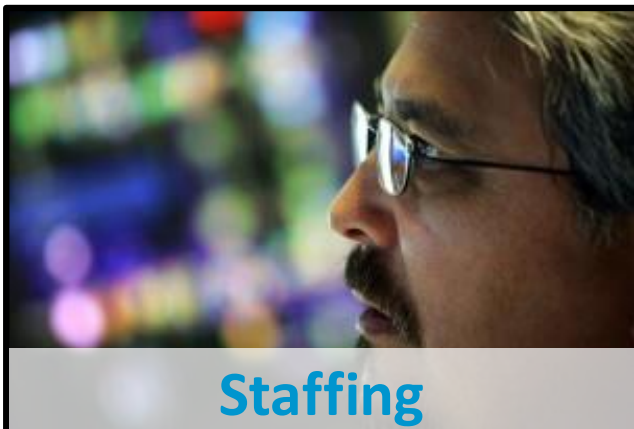
SUPPORT FOR E-GOVERNMENT



Knowhow



Development



Staffing



Infrastructure



WWW.SMALS.BE

Smals Research 2021



**Innovation with
new technologies**



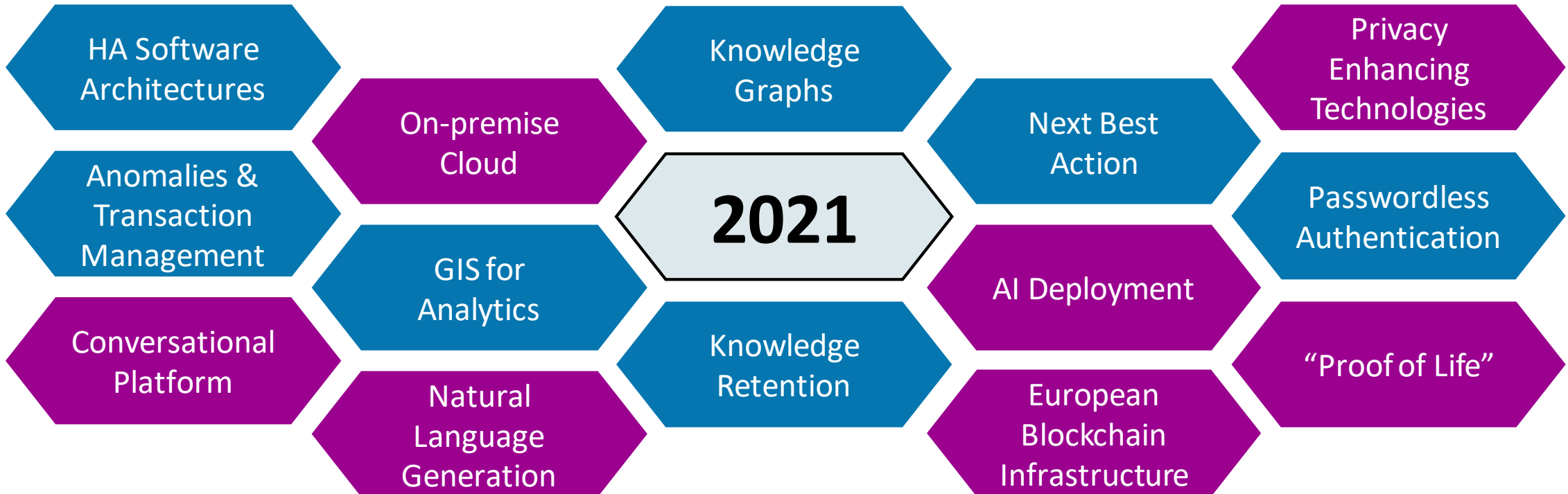
**Consultancy
& expertise**



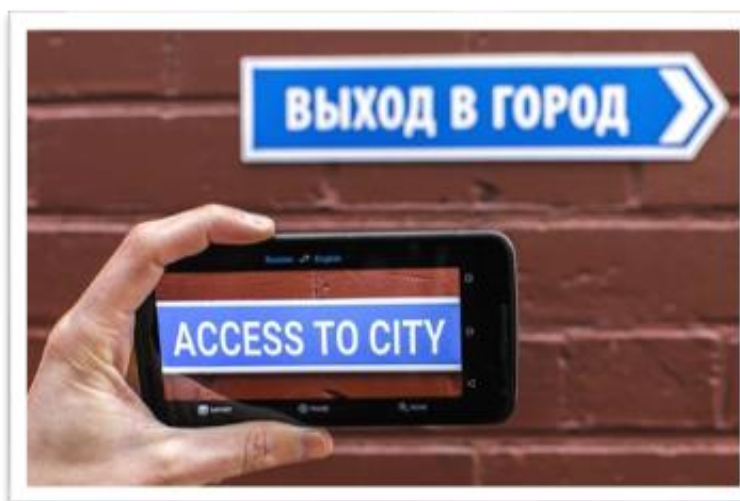
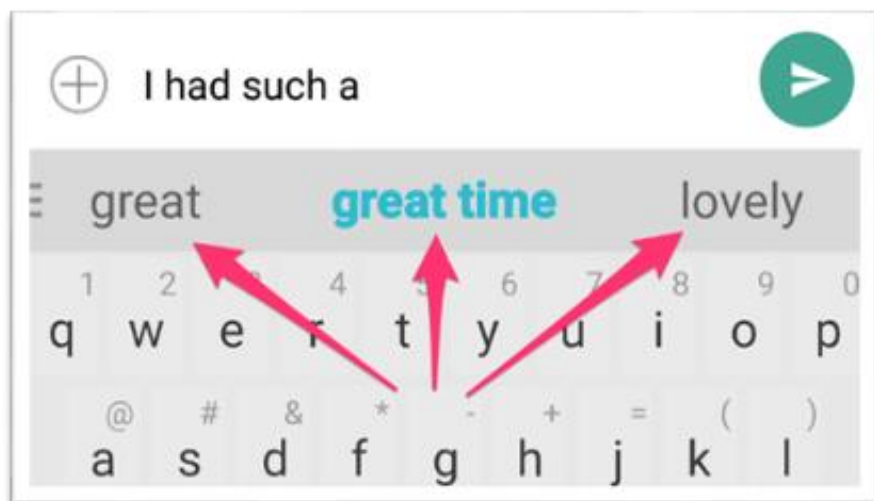
**Internal & external
knowledge transfer**



**Support for
going live**



NLP: de computationele technieken waarmee we geschreven / gesproken tekst kunnen **ontleden, analyseren en interpreteren**.



Statistisch lettercombinaties associëren ≠ taal begrijpen!

- Automatic speech recognition
- CCG
- Common sense
- Constituency parsing
- Coreference resolution
- Dependency parsing
- Dialogue
- Domain adaptation
- Entity linking
- Grammatical error correction
- Information extraction
- Language modeling
- Lexical normalization
- Machine translation
- Missing elements
- Multi-task learning
- Multi-modal
- Named entity recognition
- Natural language inference
- Part-of-speech tagging
- Question answering
- Relation prediction
- Relationship extraction
- Semantic textual similarity
- Semantic parsing
- Semantic role labeling
- Sentiment analysis
- Shallow syntax
- Simplification
- Intent Detection and Slot Filling
- Stance detection
- Summarization
- Taxonomy learning
- Temporal processing
- Text classification
- Word sense disambiguation

Current state-of-the-art: <https://nlpprogress.com/>

Named Entity Recognition

In 1917, Einstein applied the general theory of relativity to model the large-scale structure of the universe. He was visiting the United States when Adolf Hitler came to power in 1933 and did not go back to Germany, where he had been a professor at the Berlin Academy of Sciences. He settled in the U.S., becoming an American citizen in 1940. On the eve of World War II, he endorsed a letter to President Franklin D. Roosevelt alerting him to the potential development of "extremely powerful bombs of a new type" and recommending that the U.S. begin similar research. This eventually led to what would become the Manhattan Project. Einstein supported defending the Allied forces, but largely denounced using the new discovery of nuclear fission as a weapon. Later, with the British philosopher Bertrand Russell, Einstein signed the Russell-Einstein Manifesto, which highlighted the danger of nuclear weapons. Einstein was affiliated with the Institute for Advanced Study in Princeton, New Jersey, until his death in 1955.

Tag colours:

LOCATION TIME PERSON ORGANIZATION MONEY PERCENT DATE

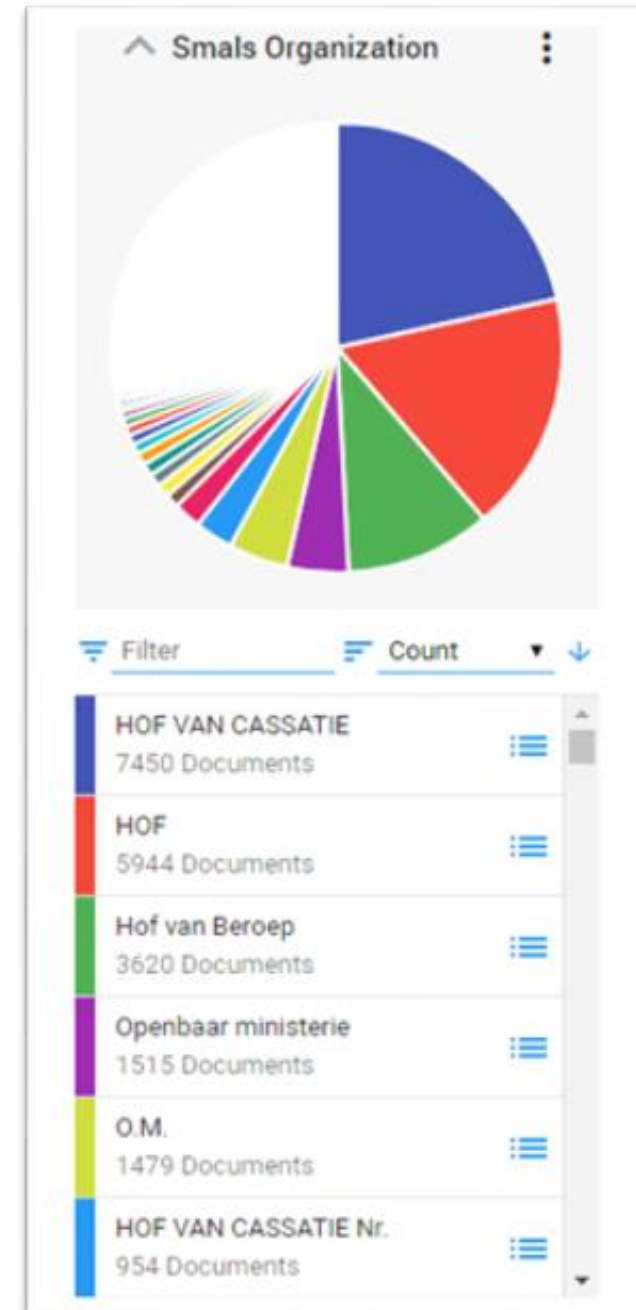
(voorbeeld © Europeana)

(De Tijd, 12/02/2021)

INTERVIEW

‘Alle Sky ECC-berichten lezen, zou ons 685 jaar kosten’

Het gerecht zou met de huidige middelen 685 jaar nodig hebben om alle berichten te lezen die onderschepd zijn bij de kraak van de misdaadtelefoons van Sky ECC. Federaal procureur Frédéric Van Leeuw



Enkele NLP libraries



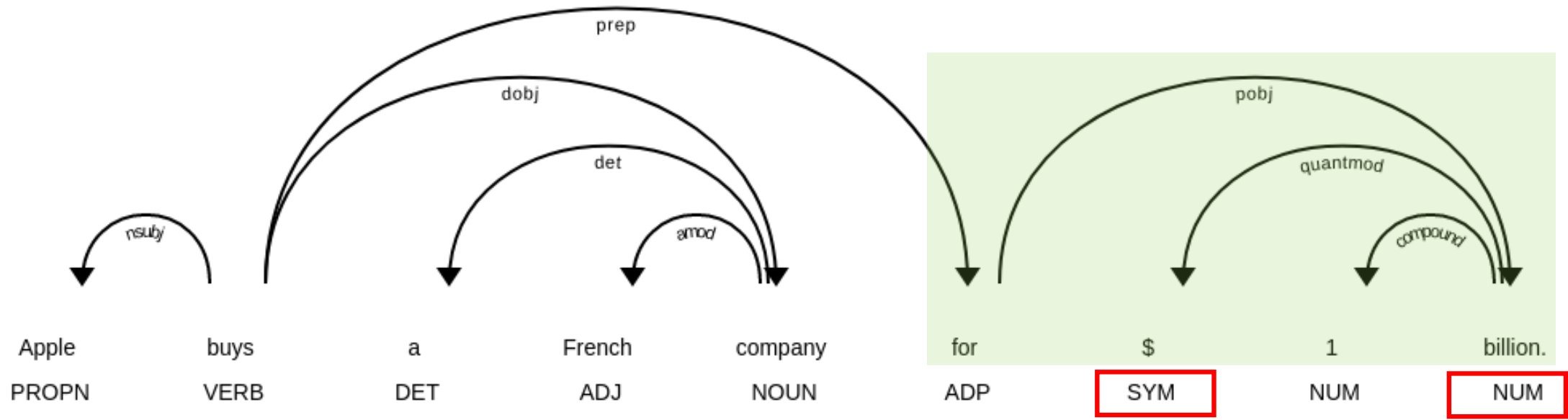
- Apple **ORG** buys a **French NORP** company for **\$1 billion MONEY** .

- Apple **PERSON** koopt een **Frans NORP** bedrijf voor \$ **1 CARDINAL** miljard.

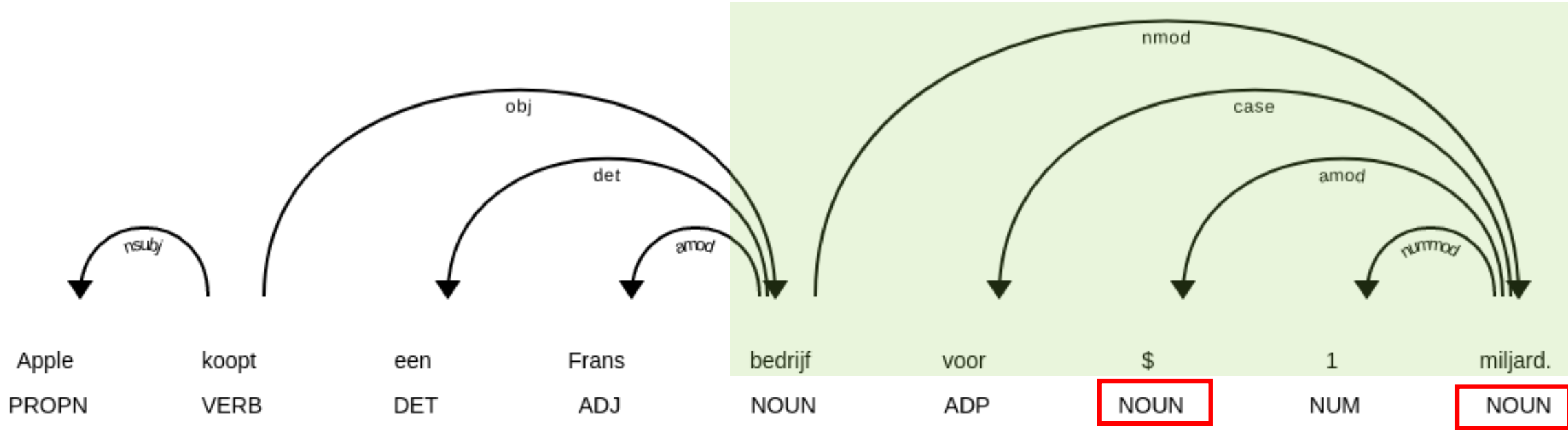


In de achtergrond

• EN:



• NL:



English

TAGGER



```
$, ', , , -LRB-, -RRB-, ., :, ADD, AFX,  
CC, CD, DT, EX, FW, HYPH, IN, JJ, JJR,  
JJS, LS, MD, NFP, NN, NNP, NNPS, NNS,  
PDT, POS, PRP, PRP$, RB, RBR, RBS, RP,  
SYM, TO, UH, VB, VBD, VBG, VBN, VBP,  
VBZ, WDT, WP, WP$, WRB, XX, ``
```

49 woordsoort-labels

"a.m.",
"Adm.",
"Bros.",
"co.",
"Co.",
"Corp.",
"D.C.",

Een 100-tal afkortingen
("tokenizer exceptions")

Nederlands

```
N|soort|ev|basis|zijd|stan__Gender=Com|Number=Sing,  
N|soort|ev|dim|onz|stan__Gender=Neut|Number=Sing,  
N|soort|mv|basis__Number=Plur, N|soort|mv|dim__Number=Plur,  
SPEC|afgebr, SPEC|afk__Abbr=Yes, SPEC|deeleigen, SPEC|enof,  
SPEC|meta, SPEC|symb, SPEC|vreemd__Foreign=Yes, TSW,  
TW|hoofd|nom|mv-n|basis, TW|hoofd|nom|mv-n|dim,  
TW|hoofd|nom|zonder-n|basis, TW|hoofd|nom|zonder-n|dim,  
TW|hoofd|prenom|stan, TW|hoofd|vrij, TW|rang|nom|mv-n,
```

... → 233 woordsoort-labels

"b.v.",
"b.ver.coll.gem.gem.comm.",
"b.verg.r.b.",
"b.versl.",
"b.vl.ev"

... → 1500+ afkortingen

English

OntoNotes 5.0

- 300000 zinnen
- 2.9mln woorden
- Nieuws, wiki, fora, conversaties, bijbel, ...

Nederlands

UD LassySmall + Alpino (*)

- 20966 zinnen
- 306764 woorden
- wiki + nieuws uit begin jaren '00

(ref: Bijbelvertaling = ± 900000 woorden)

(* uitgez. word vectors: CommonCrawl+Wiki)

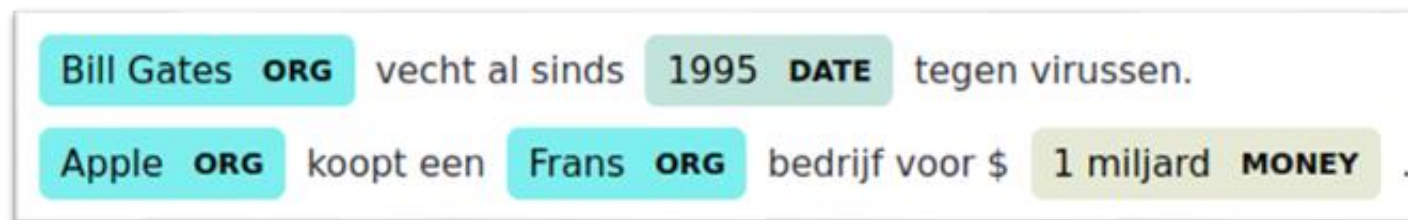
Model	Labeled dependencies	Part-Of-Speech (detailed)	NER F-score
Nl_core_news_sm (CNN)	81%	94%	72%
Nl_core_news_lg (CNN)	84%	95%	77%
En_core_web_sm (CNN)	90%	97%	84%
En_core_web_lg (CNN)	90%	97%	86%
En_core_web_trf (roBERTa)	94%	98%	90%

Situatie
25/03/2021

- Zelf gedefinieerde componenten invoegen

- Kloon een van de voorbeeldprojecten (<https://spacy.io/usage/projects>)
-
- Pas de “makefile” *project.yml* aan
- Voeg eigen trainingsdata en preprocessing scripts toe:

```
[ "OnePlus 9 Pro met nieuwe Sony-sensor verschijnt eind maart voor 899 euro.",  
  { "entities": [[0,7,"ORG"],[25,29,"ORG"],[64,72,"MONEY"]] },  
 [ "Gerucht: Discord voert gesprekken met Microsoft over mogelijke overname.",  
   { "entities": [[9,16,"ORG"],[38,47,"ORG"]] },  
 ...
```



Bill Gates **ORG** vecht al sinds **1995** **DATE** tegen virussen.
Apple **ORG** koopt een **Frans** **ORG** bedrijf voor \$ **1 miljard** **MONEY** .

- Risico: **catastrophic forgetting** → trainingsdata finetunen

- Vertrek ook hier van een bestaand project
 -
- Leg een eigen, volledige dataset aan
- Behoudt alleen labels die je nodig hebt

Back in 2000 , **People Magazine** **PUBLISHER** highlighted **Prince Williams'** **PERSON** style who at the time was a little more fashion-conscious , even making fashion statements at times .

Now-a-days the prince mainly wears **navy** **COLOR** **suits** **ITEM** (sometimes **double-breasted** **DESIGN**) , **light blue** **COLOR** **button-ups** **ITEM** with **classic** **LOOK** **pointed** **DESIGN** **collars** **PART** , and **burgundy** **COLOR** **ties** **ITEM** .

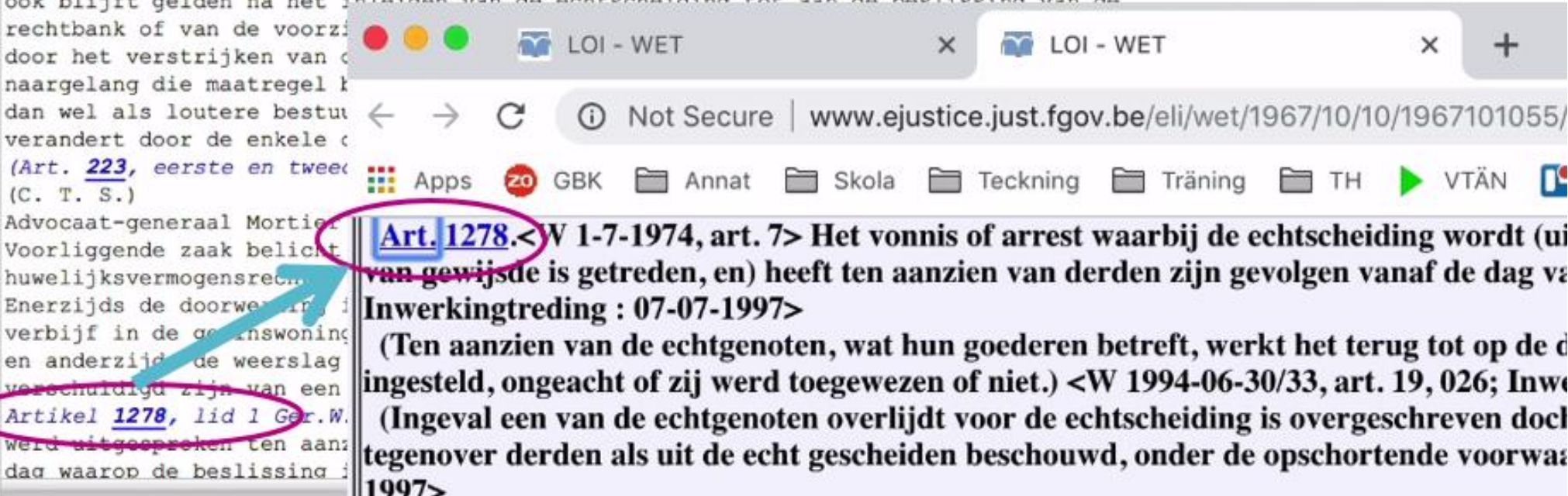
But who knows what the future holds ...

Duchess Kate **PERSON** did wear an **Alexander McQueen** **BRAND** **dress** **ITEM** to the **wedding** **OCCASION** in the **fall of 2017** **SEASON** .

- Dataset: gepubliceerde arresten

- Bvb. herformatteer herkende entiteiten als URLs

proof-of-concept uitgewerkt door TheMatchbox op NLP4Gov Hackathon, Informatie Vlaanderen, 2018:



The screenshot shows a web browser window with two tabs labeled 'LOI - WET'. The address bar displays 'www.ejustice.just.fgov.be/eli/wet/1967/10/10/1967101055/'. The browser's taskbar includes icons for 'Apps', 'GBK', 'Annat', 'Skola', 'Teckning', 'Träning', 'TH', and 'VTÄN'. The main content area displays a legal document with several annotations:

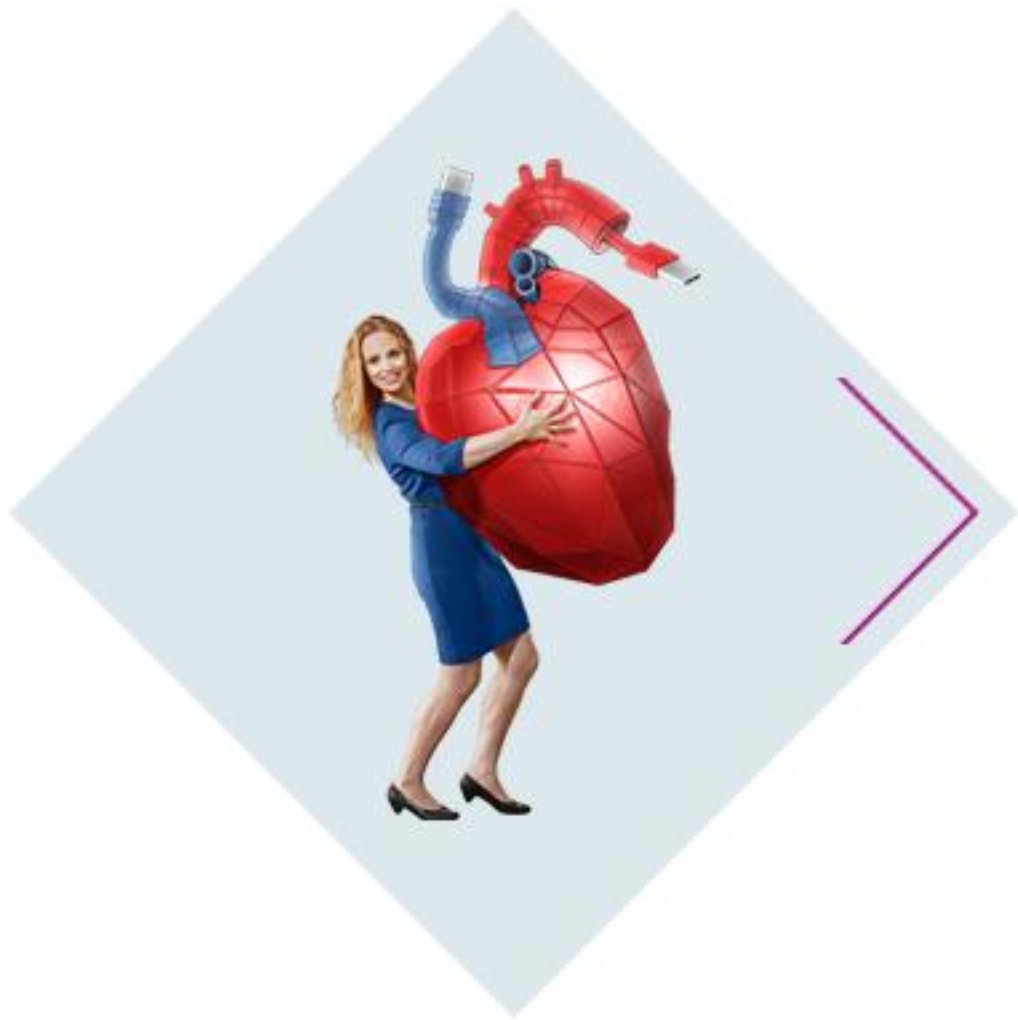
- A red circle highlights the text '(Art. 223, eerste en tweede (C. T. S.))' in the left margin.
- A red circle highlights the text 'Art. 1278.' in the main text.
- A red circle highlights the text 'Artikel 1278, lid 1 Ger.W.' in the left margin.
- A blue arrow points from the 'Art. 1278.' annotation to the text 'Inwerkingtreding : 07-07-1997>'.

The main text of the document includes: 'Advocaat-generaal Mortier Voorliggende zaak belicht huwelijksvermogensrecht. Enerzijds de doorwerking verbijf in de gemeenschap en anderzijds de weerslag verschuldigd zijn van een Inwerkingtreding : 07-07-1997> (Ten aanzien van de echtgenoten, wat hun goederen betreft, werkt het terug tot op de ingesteld, ongeacht of zij werd toegewezen of niet.) <W 1994-06-30/33, art. 19, 026; Inwe (Ingeval een van de echtgenoten overlijdt voor de echtscheiding is overgeschreven doch tegenover derden als uit de echt gescheiden beschouwd, onder de opschortende voorwa 1997>'.

- Met behulp van Knowledge Bases

- Zie Sofie Van Landeghem, "Training a custom Entity Linking model with spaCy", [YouTube](#)

- Resterende problemen
 - Structureel gebrek aan grote Nederlandstalige geannoteerde datasets
 - Nederlandse taalmodellen lopen achter op Engelstalige
- Het goede nieuws
 - Zelf een taalmodel tweaken



Met dank aan

Katy Fokou
SpaCy devs & contributors
Sofie Van Landeghem
Yves Peirsman (NLP Town)
TheMatchbox
RU Groningen

...

www.smalsresearch.be

www.smals.be/jobs

Joachim Ganseman
joachim.ganseman@smals.be