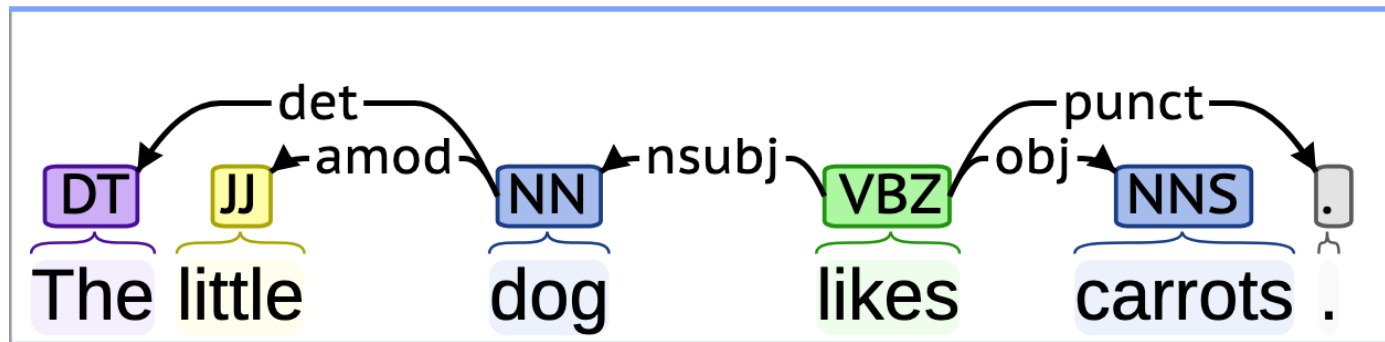# PARSING

# Les difficultés des parsers

- Les ambiguïtés

- Les structures complexes

- Les unités phraséologiques (MWE : multi-word expressions, MWU : multi-word units). Le projet PARSEME
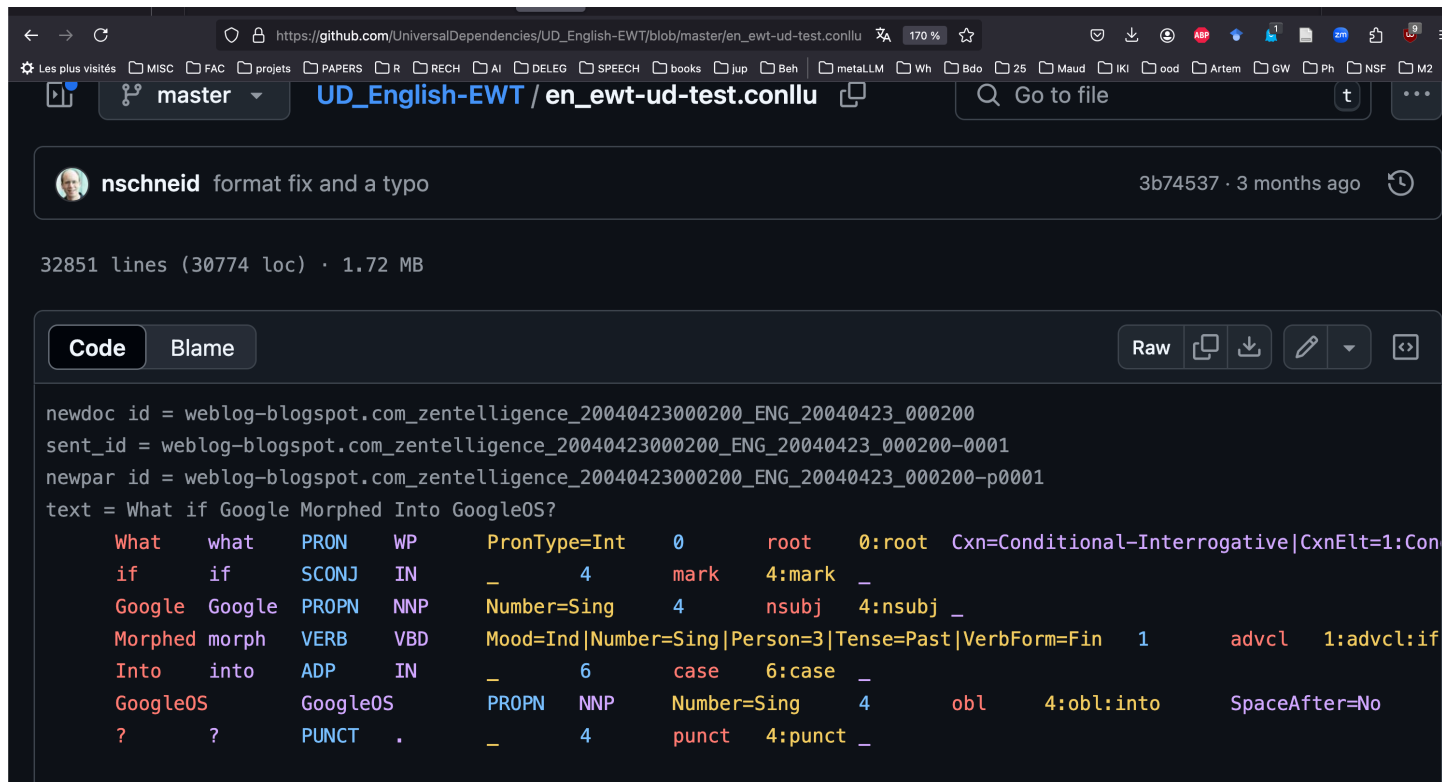
# Les PARSERS (analyseurs syntaxiques)

- Deux générations d'analyse : analyse en arbres, puis analyse en relations de dépendances
- Des relations de dépendances (deprel) entre une tête (HEAD, GOV) et son dépendant

# Parsers en dépendance (suite)

- Des catégories posées comme universelles (upos)
- Des catégories issues de jeux d'étiquettes (tagsets : xpos)
- Un format tabulaire vertical :
- Des banques d'exemples (Treebanks) pour le train, dev, test
- Des outils d'apprentissage : UDPipe
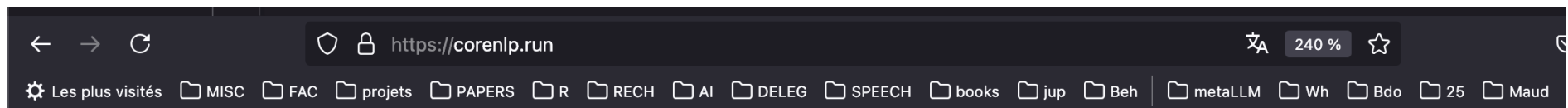
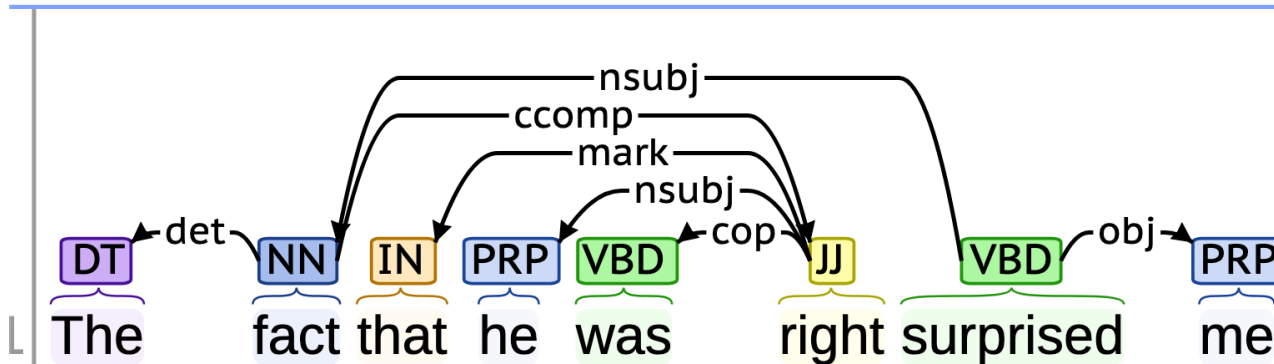# Dependency parsing : format CONNL-U

# CoNLL-U "Conference on Natural Language Learning - Universal Dependencies"

- a standardized format for annotating linguistic data, particularly for tasks in computational linguistics and natural language processing (NLP)

- . The CoNLL-U format is an extension of the original CoNLL format, specifically designed to support the Universal Dependencies project

- .**Key Features of CoNLL-U**

1. **Structure**: CoNLL-U files are plain text files (UTF-8 encoded) that contain three types of lines[1]

2. : Word lines with 10 tab-separated fields for token annotation

3. Blank lines marking sentence boundaries

4. Comment lines starting with '#' for sentence-level metadata

5. **Token Annotation**: Each word or token is described using 10 fields, including[1]

6. : ID: Token identifier

7. FORM: The word form as it appears in the text

8. LEMMA: The base or dictionary form of the word

9. UPOS: Universal part-of-speech tag

10. HEAD: Syntactic head in dependency parsing

11. DEPREL: Dependency relation to the HEAD

12. **Extensibility**: The CoNLL-U format can be extended to CoNLL-U Plus for additional annotation layers while maintaining compatibilit

13. **Standardization**: CoNLL-U has become a standard format in NLP due to its simplicity and effectiveness in handling annotated linguistic data

- The CoNLL-U format facilitates various NLP tasks, including part-of-speech tagging, syntactic parsing, and named entity recognition, by providing a structured and standardized way to represent linguistic annotation

# CoreNLP (Stanford)

# Des Treebanks disponibles sur Github

- https://github.com/UniversalDependencies

- Documentation :
- https://universaldependencies.org/

# En anglais



https://universaldependencies.org/

# En français

**French treebanks**

| | | | | | | |
|---|---|---|---|---|---|---|
| ▸ | **GSD** | 400K | ⓛⒻ | 📅📖👍W | CC BY SA | ★★★★☆ |
| ▸ | **Sequoia** | 70K | ⓛⒻ | ✏️📖ⓘW | GPL | ★★★★☆ |
| ▸ | **ParTUT** | 28K | ⓛⒻ | 🔨📖W | CC BY NC SA | ★★★★☆ |
| ▸ | **ParisStories** | 42K | ⓛⒻ | 💬 | CC BY SA | ★★★★☆ |
| ▸ | **Rhapsodie** | 44K | ⓛⒻ | 💬 | CC BY SA | ★★★★☆ |
| ▸ | **PUD** | 24K | ⓛⒻ | 📖W | CC BY SA | ★★★☆☆ |
| ▸ | **FQB** | 23K | ⓛⒻ | 📖ⓘ | GPL | ★★☆☆☆ |

See here for comparative statistics of French treebanks.

https://universaldependencies.org/

# Stanza

# Spacy

```
EXAMPLE

import spacy_udpipe


spacy_udpipe.download("en") # download English model


text = "Wikipedia is a free online encyclopedia, created and e
nlp = spacy_udpipe.load("en")


doc = nlp(text)
for token in doc:
    print(token.text, token.lemma_, token.pos_, token.dep_)
```

- https://spacy.io/universe/project/spacy-udpipe

# Le réentraînement de Treebanks spécialisées

- Stanza uses multiple English treebanks for different purposes:

- For general English, Stanza uses the Universal Dependencies (UD) treebanks, specifically the English Web Treebank (EWT)
- . For biomedical text, Stanza provides two separate syntactic analysis pipelines: a. One trained on the CRAFT (Colorado Richly Annotated Full Text) treebank
- .
- b. Another trained on the GENIA treebank
- . For clinical text, Stanza uses a combination of: a. The English Web Treebank (EWT).
- b. A silver-standard MIMIC treebank created from clinical notes in the MIMIC-III database
- .

- These treebanks are combined in various ways to improve robustness and performance across different domains. For instance, the CRAFT pipeline combines the English Web Treebank with the CRAFT treebank, while the clinical pipeline combines the EWT with the silver-standard MIMIC treebank.