# Natural Language processing (NLP) "Levels" of linguistic analysis

Thierry Poibeau (CNRS & PSL/ENS, Lattice)

thierry.poibeau@ens.psl.eu

DHAI Instensive Week, 29 March 2022

PR[AI]RIE
PaRis Artificial Intelligence Research InstitutE

ENS | PSL★

cnrs

# Intro (1): Natural Language Processing

- **Natural language processing** (**NLP**) is a subfield of linguistics, computer science, and artificial intelligence concerned with the interactions between computers and human language, in particular how to program computers to process and analyze large amounts of natural language data. The result is a computer capable of "understanding" the contents of documents, including the contextual nuances of the language within them. The technology can then accurately extract information and insights contained in the documents as well as categorize and organize the documents themselves.

https://en.wikipedia.org/wiki/Natural_language_processing

# Intro (2): Natural Language Processing

- Challenges in natural language processing frequently involve [speech recognition](#), [natural language understanding](#), and [natural-language generation](#)


- As a summary, NLP is a super wide area of research

- Here, we will only address (a very tiny part of) natural language understanding

# Levels of Linguistic Analysis

Why is NLP Hard?
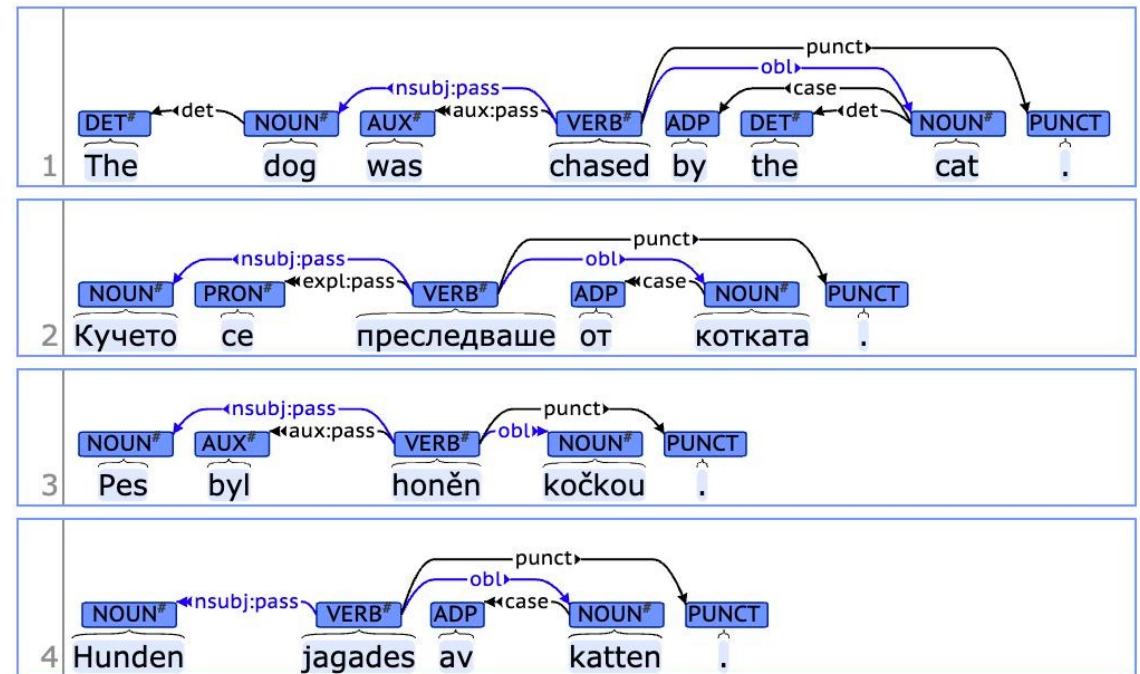Evaluation
Implementations
Conclusion

# Natural Language Pyramid



Natural Language Processing Pyramid

# Part of speech (pos) tagging

The_DT first_JJ time_NN he_PRP was_VBD shot_VBN in_IN the_DT
hand_NN as_IN he_PRP chased_VBD the_DT robbers_NNS outside_RB .\_.

| first | time | shot | in | hand | as | chased | outside |
|-------|------|------|-----|------|-----|--------|---------|
| JJ | NN | NN | IN | NN | IN | JJ | IN |
| RB | VB | VBD | RB | VB | RB | VBD | JJ |
|    |    | VBN | RP |    |    | VBN | NN |
|    |    |     |    |    |    |     | RB |

- Lots of tools available (TreeTagger, Stanford Tagger, UD Tagger…) for numerous languages
- Accuracy (F-measure): often .9-.97 on standard text

# Parsing (automatic syntactic analysis)

- Large diversity of tools and resources
- [https://universaldependencies.org/](https://universaldependencies.org/) ~100 languages, ~200 treebanks with similar annotations
-  Accuracy (F-measure): hard to predict, generally .7-.9 on standard text

# Semantics

- A wide diversity of tasks
  - Word sense disambiguation (WSD)
  - Named entity recognition
  - Term recognition / Terminology
  - Text zoning
  - Event recognition
  - Sentiment / Opinion mining

When **Sebastian Thrun** `PERSON` started at **Google** `ORG` in **2007** `DATE` , few people outside of the company took him seriously. "I can tell you very senior CEOs of major **American** `NORP` car companies would shake my hand and turn away because I wasn't worth talking to," said **Thrun** `PERSON` , now the co-founder and CEO of online higher education startup Udacity, in an interview with **Recode** `ORG` **earlier this week** `DATE` .

A little **less than a decade later** `DATE` , dozens of self-driving startups have cropped up while automakers around the world clamor, wallet in hand, to secure their place in the fast-moving world of fully automated transportation.

From Simone Teufel (2006)

# Applications

- Spell and grammatical checker
- Search engine
- Information extraction
- Text summarization
- Machine translation
- Dialogue, conversational agents
- Opinion mining
- etc.

Levels of Linguistic Analysis
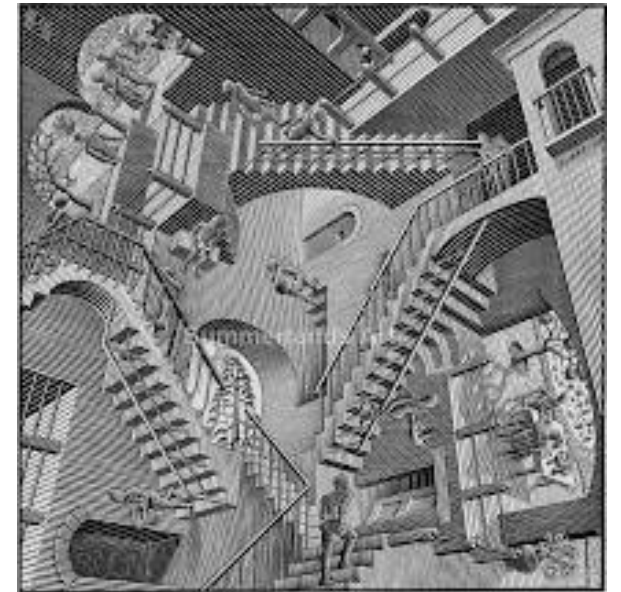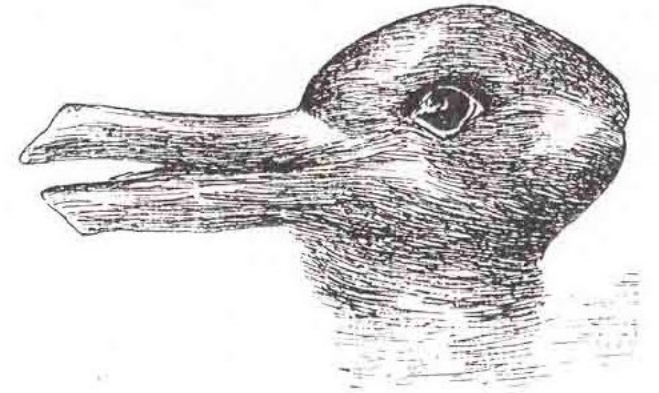Why is NLP Hard?
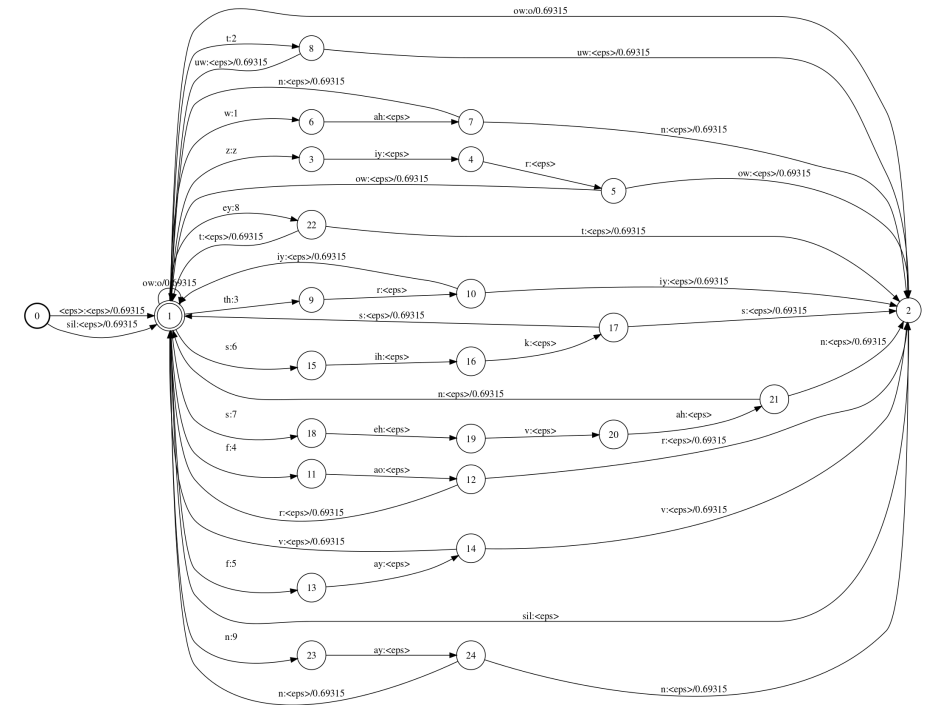Evaluation
Implementations
Conclusion

# Why is NLP hard?

- Words are polysemous, language is ambiguous
  - Il a free, il a tout compris.
  - I buried $100 in the bank.
  - Flying planes are dangerous.
  - We saw her duck (R. Nordquist)
- Not only a linguist' / an artificial problem. Cf "*Good*"
  - "useful" or "functional" (*That's a good hammer*)
  - "exemplary" (*She's a good student*),
  - "pleasing" (*This is good soup*),
  - "moral" (*a good person*)
  - "righteous (*I have a good daughter*)

# How does it work? Rule-based systems

- 1950-1990: Rule-based systems
  - Dictionary + grammar
  - Finite state transducers
- Benefits
  - Easy to read and develop
  - "Naturalness" of the approach
- Limitations
  - Poor coverage
  - Hard to maintain
  - Unsuitable for some task (WSD)

# How does it work? Statistical systems

- ## 1990-2014: Statistical systems
  - Learn a model from representative data
  - Apply it to new data
- ## Benefits
  - Good coverage
  - Takes into account the statistical nature of language
- ## Limitations
  - Needs annotated data
  - Takes into account only the local context

# How does it work? Deep learning



- 2014-: Deep learning systems
  - A continuation of the previous approach
- Benefit
  - Better coverage, better generalizations
  - Takes into account larger contexts (transformers)
  - Universal representation (through vectors)
- Limitations
  - Needs even more annotated data
  - Exact nature of the generalizations unknown
  - Not fully reliable

# Beware of Language Diversity!



New Guinea

- 7000 languages in the world
- Very few of them have accurate NLP tools

Levels of Linguistic Analysis
Why is NLP Hard?
Evaluation
Implementations
Conclusion

# Overview on Evaluation

- Evaluation play a central role in NLP
  - Compare different approaches
  - Monitor progression of the field
  - Human evaluation is costly and often not so reliable
- A large number of metrics have been proposed (a research area in itself)
  - Bleu for Machine translation, Rouge for summarization, etc.
- However, most NLP tasks can be evaluated using precision and recall

# Evaluation: Precision and Recall

- For any element in the dataset
  - Have all the relevant elements been tagged?
  - Are all the tagged elements relevant?

- Evaluation indicators
  - Precision: # relevant tagged  element /# tagged elements
  - Recall: # relevant tagged  element /# elements to be tagged
  - F-measure = 2 * (P * R) / P + R

Source : Wikipedia, https://en.wikipedia.org/wiki/Precision_and_recall

Levels of Linguistic Analysis

Why is NLP Hard?

Evaluation

**Implementations**

Conclusion

# Lots of resources available!

- Lots of data, corpora and open source code available online (remember UD)

- Most recent models are open source and easy to integrate (including models from companies like Facebook, cf. FastText, or Google, cf. Bert)

- Deep learning means GPU are requited (Google Colab can help)

- But also note: there are lots of languages with few or nearly no resources!   (~7000 languages in the world)

# Hugging Face

- Recent developments (Transformers, cf. Bert in different languages etc. https://huggingface.co/

# Spacy

- General NLP: Spacy  ([https://spacy.io/](https://spacy.io/))
- Integrates Huggingface work on recent large scale NLP models
- Easier to use than directly manipulating Hugging Face code

# Spacy

- Models for different languages

# Github / Jupyter Notebooks

- Jupyter notebooks, for example Sentiment analysis in French on a corpus made of Allocine reviews

# Named entity recognition using Spacy



Jupyter notebook,
adapted from C. Brando,
Master HN PSL

Levels of Linguistic Analysis

Why is NLP Hard?

Evaluation

Implementations

**Conclusion**

# Conclusion

- A field that evolve quickly
- Lots of progress, but lots of challenges remaining
- Bigger may not always be better

- Keep in mind language diversity, and under resource languages
- Keep in mind ethical issues (e.g. gender bias in language models)

# To go further...

- Check the Bible: https://web.stanford.edu/~jurafsky/slp3/ (free!)

## Speech and Language Processing (3rd ed. draft)

**Dan Jurafsky** and **James H. Martin**

**Here's our December 30, 2020 draft! Includes:**

- new version of Chapter 8 (bringing together POS and NER in one chapter),
- new version of Chapter 9 (with Transformers)
- Chapter 11 (MT)
- neural span parsing and CCG parsing moved into Chapter 13 (Constituency Parsing) and Statistical Constituency Parsing moved to Appendix C
- new version of Chapter 23 (QA modernized)
- Chapter 26 (ASR + TTS)
- Plus a modernizing pass (and typo fixing, thanks to all of you!!!) on all the other chapters.

We are really grateful to all of you for finding bugs and offering great suggestions!

Individual chapters are below; here is a single pdf of all the chapters in the December 30, 2020 draft of the book-so-far

As always, typos and comments very welcome (just email slp3edbugs@gmail.com and let us know the date on the draft)!
(Due to reorganizing, still expect some missing latex cross-references throughout the pdfs, don't bother reporting those missing ref/typos.)

Feel free to use the draft slides in your classes.
We are in the process of updating the slides now; so far the slides for Chapters 2, 3, 4, 5, 6, 20, and 24 have been updated.

When will the whole book be finished?
We're shooting for well before the end of 2021 for the 3 remaining chapters (Intro, Contextual Embeddings, Semantic Parsing) + random missing sections.

And if you need last year's draft chapters, they are here.

| | Chapter | Slides | Relation to 2nd ed. |
|---|---|---|---|
| 1: | Introduction | | [Ch. 1 in 2nd ed.] |
| 2: | Regular Expressions, Text Normalization, Edit Distance | 2: Text Processing [pptx] [pdf] <br> 2: Edit Distance [pptx] [pdf] | [Ch. 2 and parts of Ch. 3 in 2nd ed.] |
| 3: | N-gram Language Models | 3: N-grams [pptx] [pdf] | [Ch. 4 in 2nd ed.] |
| 4: | Naive Bayes and Sentiment Classification | 4: Naive Bayes + Sentiment [pptx] [pdf] [new in this edition] |  |