# Advanced Topics in NLP
## Introduction

**Dr. Paul Buitelaar & Dr. Omnia Zayed**

**Data Science Institute**

**University of Galway**

# Learning Objectives of this Course

- Gain insights into knowledge extraction from text, in particular around entities and relations

- Gain insights into opinion mining, in particular on emotion analysis, dynamic identification of aspect and analysis of figurative language

- Gain insights into language generation, in particular in the context of machine translation and chatbot development

OLLSCOIL NA GAILLIMHE
UNIVERSITY OF GALWAY

# Administrative Issues

# Lecturers



Dr. Paul Buitelaar & Dr. Omnia Zayed

# Lecture Plan

| Date | | Lecture |
|---|---|---|
| 9/1/23 | | Introduction |
| 16/1/23 | | Knowledge Extraction I: Entities |
| 23/1/23 | | Knowledge Extraction II: Relations |
| 30/1/23 | | Opinion Mining I: Emotions |
| 6/2/23 | | bank holiday |
| 13/2/23 | | Opinion Mining II: Aspect |
| 20/2/23 | | Opinion Mining III: Figurative Language |
| 27/2/23 | | Language Generation I: Machine Translation I |
| 6/3/23 | | Language Generation II: Machine Translation II + NLG |
| 13/3/23 | | Language Generation III: Dialog Systems |
| 20/3/23 | | Summary |
| 27/3/23 | | Industry Talk |

OLLSCOIL NA GAILLIMHE
UNIVERSITY OF GALWAY

# Labs

**Time & Venue:** Fridays 11am-1pm, IT102

Practical exercises covering the course content

Prerequisites are basic knowledge of Python, see:
https://docs.python.org/3/tutorial/

We will use Google colab Jupyter notebook, see:
https://colab.research.google.com


Dhairya Dalal


Ali Hatami

# Lab Exercises

| Date | | Lab |
|---|---|---|
| 13/1/23 | | intro to base NLP tools and methods |
| 20/1/23 | | Named Entity Recognition |
| 27/1/23 | | relation prediction |
| 3/2/23 | | emotion classification |
| 10/2/23 | | no lab - bank holiday |
| 17/2/23 | | aspect-based sentiment analysis |
| 24/2/23 | | metaphor classification |
| 3/3/23 | | build a machine translation model step by step |
| 10/3/23 | | machine translation evaluation |
| 17/3/23 | | no lab - bank holiday |
| 24/3/23 | | chatbot development with RASA |

OLLSCOIL NA GAILLIMHE
UNIVERSITY OF GALWAY

# Assignments

Two assignments

- Assignment 1: Released Jan 30th, Due Feb 20th

- Assignment 2: Released Mar 6th, Due Apr 3rd

Assignments count for 50% of final grade

# Recommended Reading

Lectures

- Jurafsky and Martin, *SPEECH and LANGUAGE PROCESSING: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*, 3rd edition: https://web.stanford.edu/~jurafsky/slp3/ed3book.pdf

- Chris Manning and Hinrich Schütze, *Foundations of Statistical Natural Language Processing*, MIT Press. Cambridge, MA: May 1999: https://nlp.stanford.edu/fsnlp/

Labs

- *Dive into Deep Learning*: https://d2l.ai/

OLLSCOIL NA GAILLIMHE
UNIVERSITY OF GALWAY

Summary of Intro NLP
"What we learned so far"

University
*of*Galway.ie

What have we learned in the Intro NLP course?

# Natural Language Processing

# Linguistics

*text text text ….*

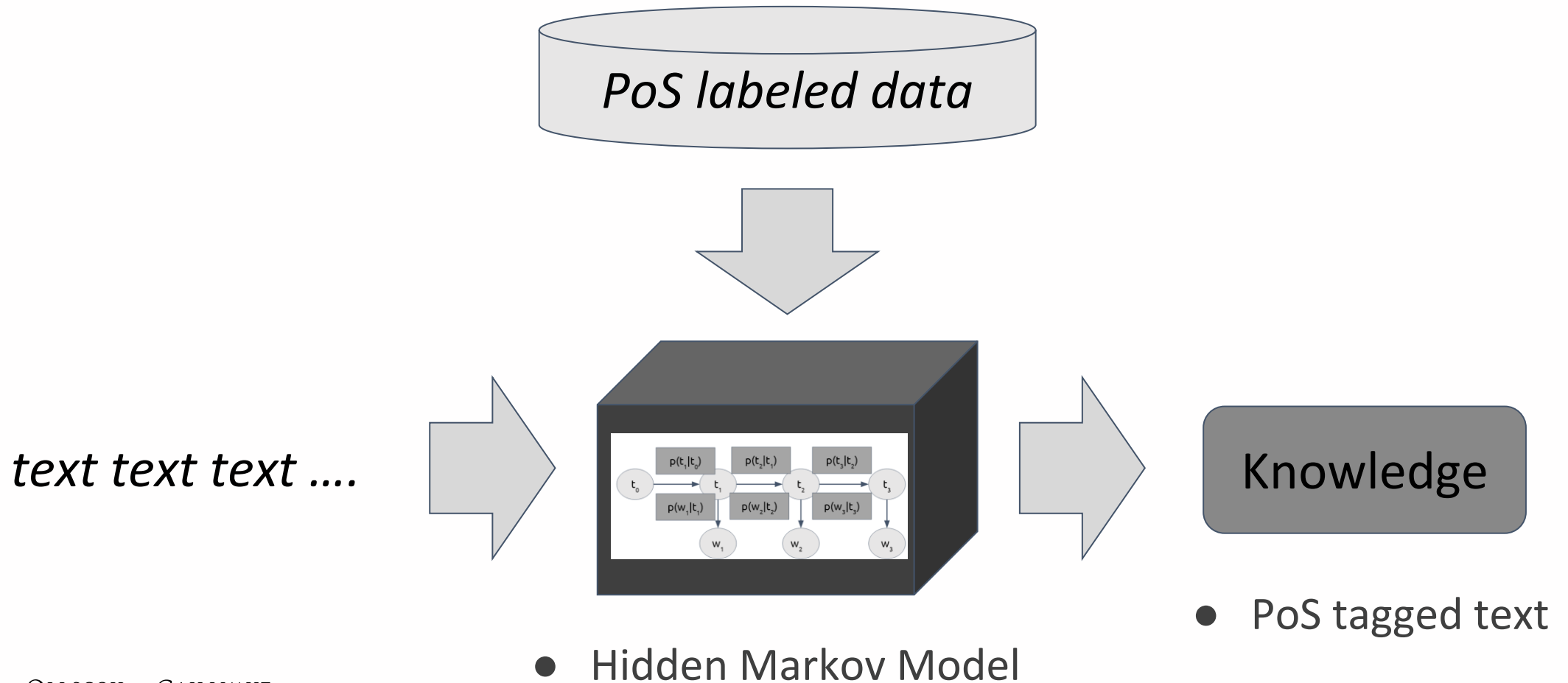## Linguistic Structure & Levels of Analysis

- **Morphology**: tokenization (MWEs), inflection, derivation, stemming, lemmatization

- **Syntax**: part-of-speech, grammar (constituency vs. dependency)

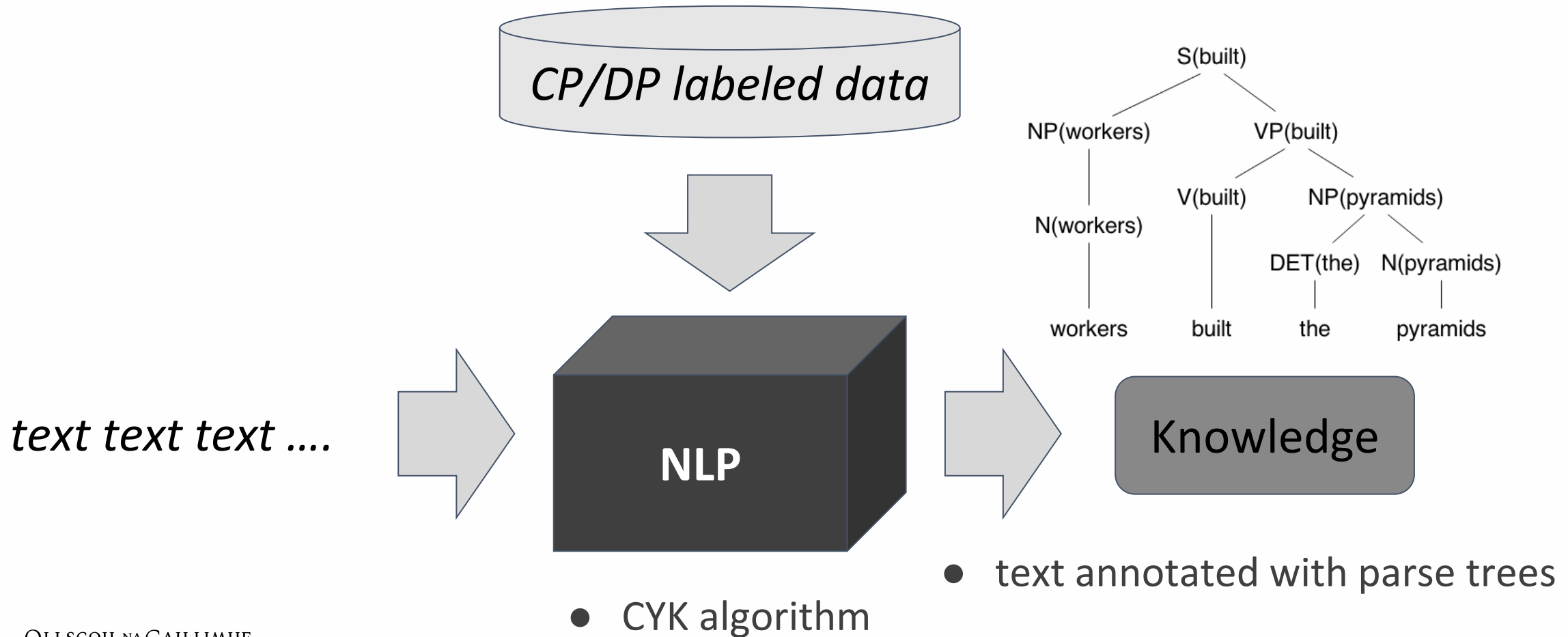- **Semantics**: word sense, semantic roles, coreference

## Language Data

- **Lexicon**: WordNet, FrameNet

- **Corpora**: annotation (data labeling), multilingual, domain-specific
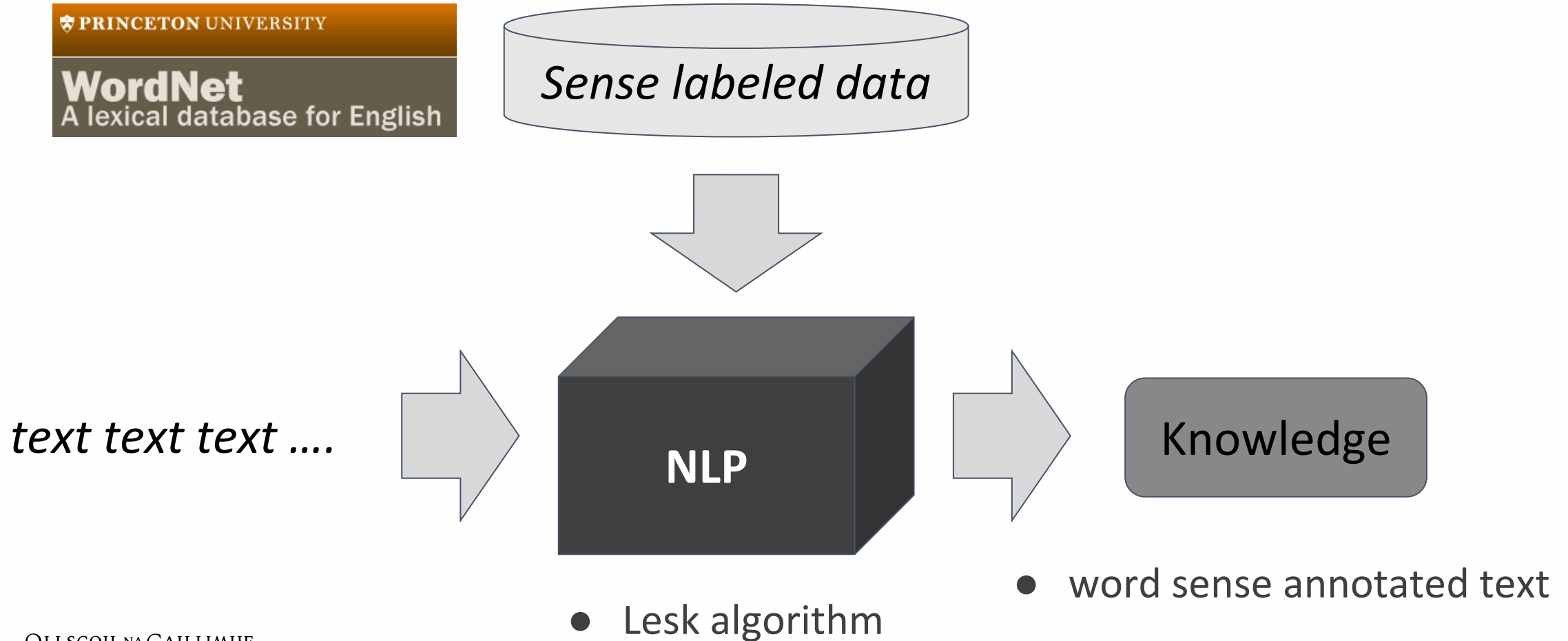
# Syntax: Part of Speech



*PoS labeled data*

*text text text ….*

$p(t_1|t_0)$  $p(t_2|t_1)$  $p(t_3|t_2)$

$t_0$  $t_1$  $t_2$  $t_3$

$p(w_1|t_1)$  $p(w_2|t_2)$  $p(w_3|t_3)$

$w_1$  $w_2$  $w_3$

Knowledge

- Hidden Markov Model

- PoS tagged text

# Syntax: Constituency / Dependency Parsing

CP/DP labeled data

```
                                    S(built)
                          NP(workers)        VP(built)
                          N(workers)      V(built)    NP(pyramids)
                                                   DET(the)  N(pyramids)
                          workers         built     the      pyramids
```

*text text text ....*

**NLP**

Knowledge

- CYK algorithm

- text annotated with parse trees

# Semantics: Word Sense Disambiguation



- Lesk algorithm
- word sense annotated text

16

# Probability



- probability model

$$P(B|A) = P(A|B)\frac{P(B)}{P(A)} \qquad TF\text{-}IDF = f_w \times \left( \log \left( \frac{N}{N_w} \right) + 1 \right)$$

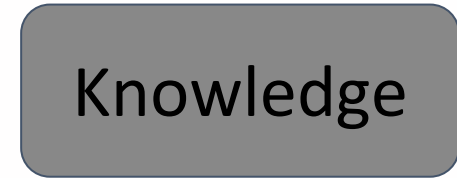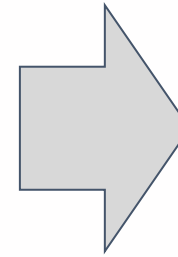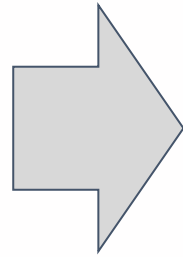- estimate word probability using Naive Bayes, TF-IDF

# Text Classification

**Set of Classes**

**Decisions**

- classifier

**text text text ….**

**NLP**

**Knowledge**

- probability model

$$P(B|A) = P(A|B)\frac{P(B)}{P(A)} \qquad TF\text{-}IDF = f_w \times \left( \log\left(\frac{N}{N_w}\right) + 1 \right)$$

OLLSCOIL NA GAILLIMHE
UNIVERSITY OF GALWAY

- estimate word probability using Naive Bayes, TF-IDF

# Language Modeling

text text text ....  →  **NLP**  →  Knowledge

- n-gram language model

$$p(w_1 w_2 \ldots w_n) = \prod_{k=1,\ldots,n} p(w_k | w_{k-m+1} \ldots w_{k-1})$$

- estimate the probability of a sentence using n-gram model

# Vector Space



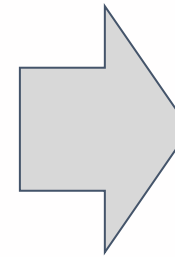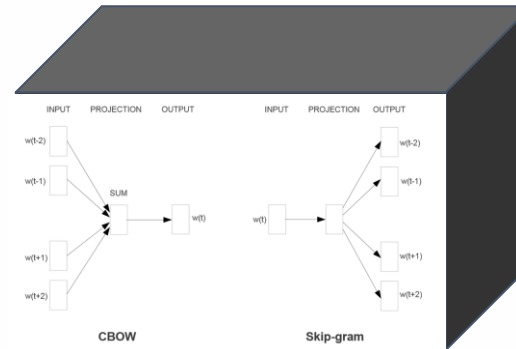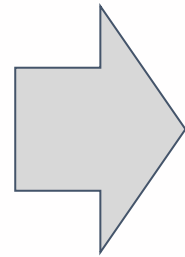text text text ....

**NLP**

- co-occurrence matrix
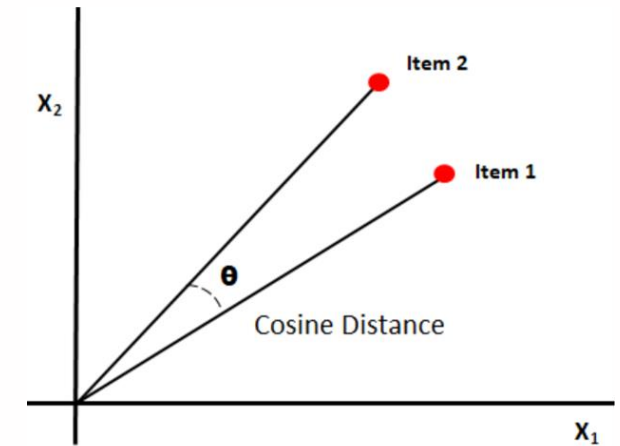
- Cosine similarity

Knowledge

- distributional model

OLLSCOIL NA GAILLIMHE
UNIVERSITY OF GALWAY

# Word Embeddings



*text text text ....*

- neural network architecture

- Cosine similarity

**Knowledge**

- word embeddings

*Image from https://www.oreilly.com/library/view/statistics-for-machine*

# RNNs: Text Classification



*Set of Classes*

*text text text ....*

Decisions

- RNN based classifier

Knowledge

- word embeddings

- neural network architecture

# Deep Learning: Transformers

*text text text ....*



- Transformer with self-attention

Knowledge

- Pretrained model

# Transformers: Text Classification



**Set of Classes**

**Decisions**

- Classifier with fine-tuning for transfer learning

**text text text ....**

**Knowledge**

- Pretrained model

- Transformer with self-attention

# Applications: Sentiment Analysis



- Sentiment analysis classifier

- Knowledge

- Pretrained model

- Transformer with self-attention

# Ethics

text text text ….

**NLP**

Decisions

- Explainable AI

Knowledge

- Data privacy (GDPR)
- Data Protection Impact Assessment
- Data Statement

- Ethical, bias-aware NLP

- Trustworthy AI

OLLSCOIL NA GAILLIMHE
UNIVERSITY OF GALWAY

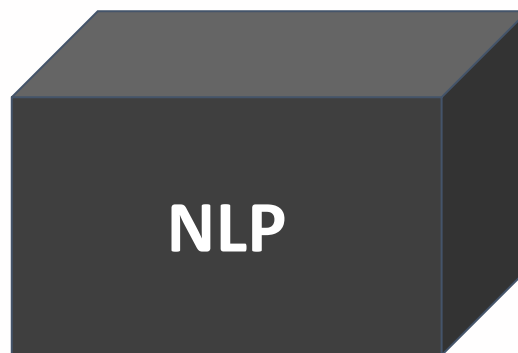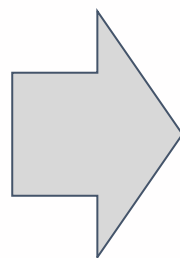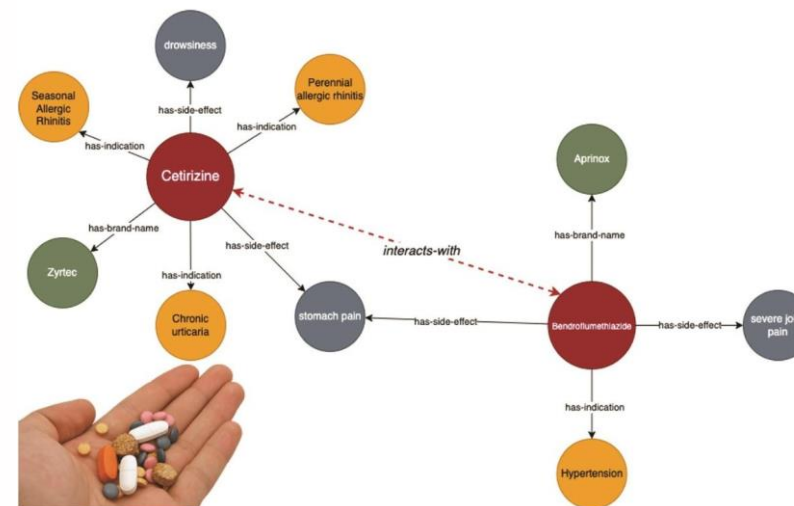What is missing?

# Knowledge Extraction Tasks

## Symbolic Knowledge

- Named Entity Recognition & Entity Linking

- Relation Extraction



*text text text ….*

**NLP**

Knowledge

- Entities, Relations

OLLSCOIL NA GAILLIMHE
UNIVERSITY OF GALWAY

# Knowledge Extraction - Definition

Creation of knowledge from **unstructured (textual, language) data**

Extracted knowledge must **facilitate inferencing**

Requires **reuse and/or generation of formal knowledge**

# Knowledge Extraction - Example

*"Naomi Carey, who is the director of Hutchinson Care Homes said she is currently only able to operate at 85% capacity."*

*"Galway GP Martin Daley, former president of the Irish Medical Organisation, gave a statement today to this effect."*

**Entities**

| | | |
|---|---|---|
| PERSON: | *Martin Daley, Naomi Carey* | |
| ORGANISATION: | *Irish Medical Organisation, Hutchinson Care Homes* | |
| CITY: | *Galway* | |

**Relations**

| | | |
|---|---|---|
| *Martin Daley* → | at-organisation | ← *Irish Medical Organisation* |
| *Naomi Carey* → | at-organisation | ← *Hutchinson Care Homes* |
| *Martin Daley* → | has-occupation | ← *GP* |
| *Naomi Carey* → | has-occupation | ← *former president* |

OLLSCOIL NA GAILLIMHE
UNIVERSITY OF GALWAY

# Opinion Mining Tasks

Emotion Analysis

Aspect Mining

Figurative Language Processing

# Emotion Analysis



happy  sad  angry  disappointed

surprised  proud  in love  scared

# Emotion Analysis

# Emotion Analysis



Deep neural network

Input layer | Multiple hidden layers | Output layer
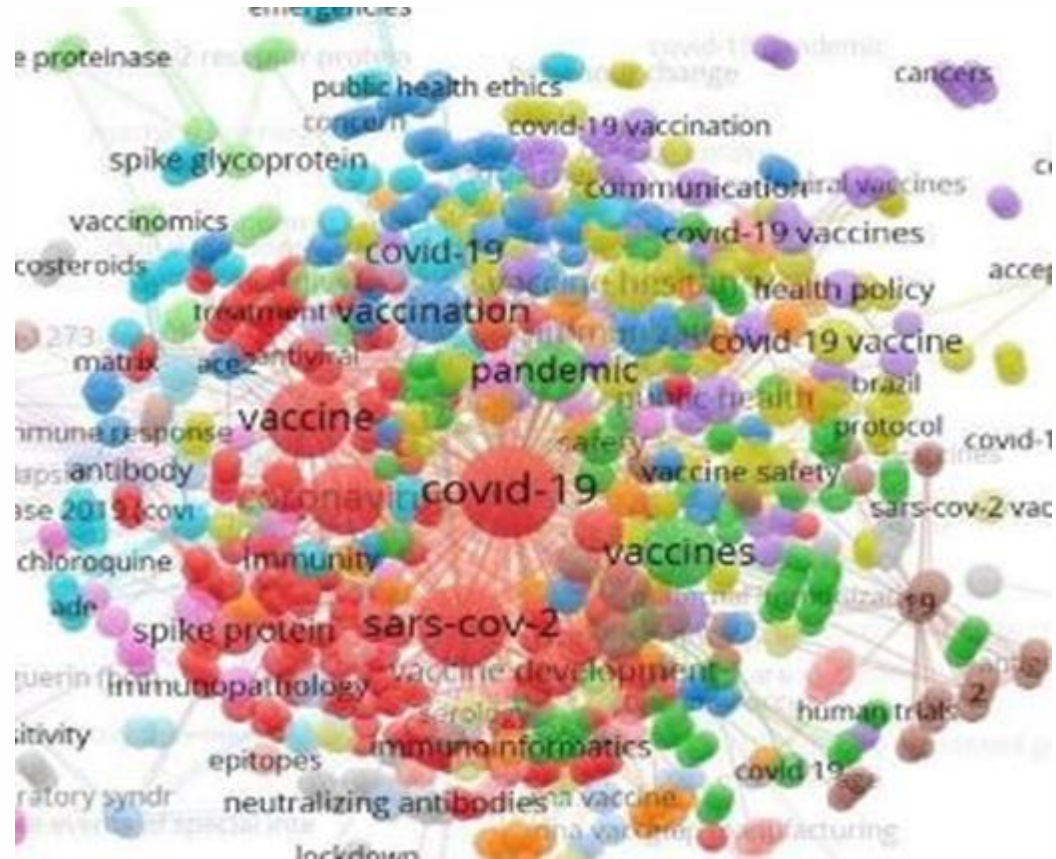
# Aspect Mining

"The camera's *focus* was **bad**, but has a **great** *size* and is **easy-to-use**."

# Aspect Mining
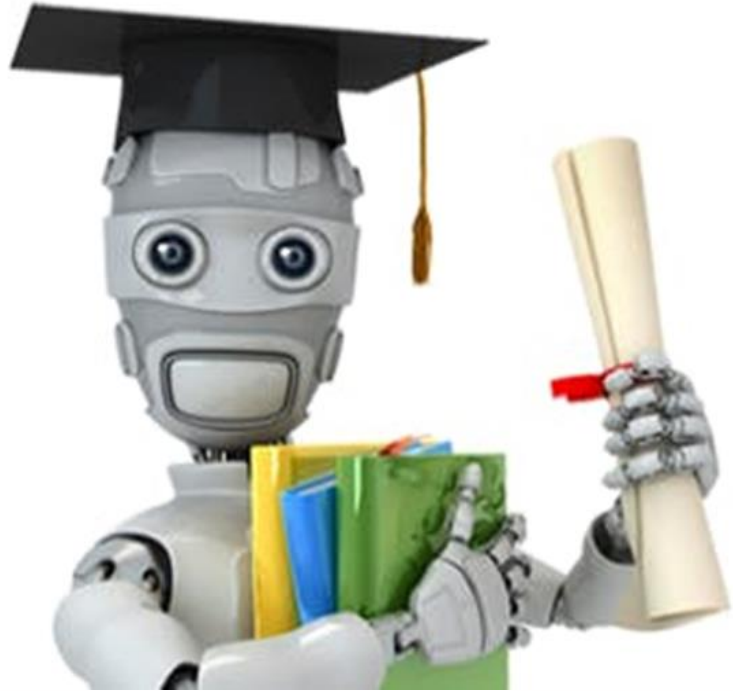
# Figurative Language Processing



Idiom (noun): a group of words established by usage as having a meaning not deducible from those of the individual words




Hate Speech


I speak fluent SARCASM

OLLSCOIL NA GAILLIMHE
UNIVERSITY OF GALWAY

38

# Figurative Language Processing



We appear to have achieved the objective in **sweeping fashion**.

Translate Tweet

# Language Generation Tasks

Machine Translation

Data-to-Text Generation

Dialog System Development

# Language Generation: Machine Translation

*English text ….*
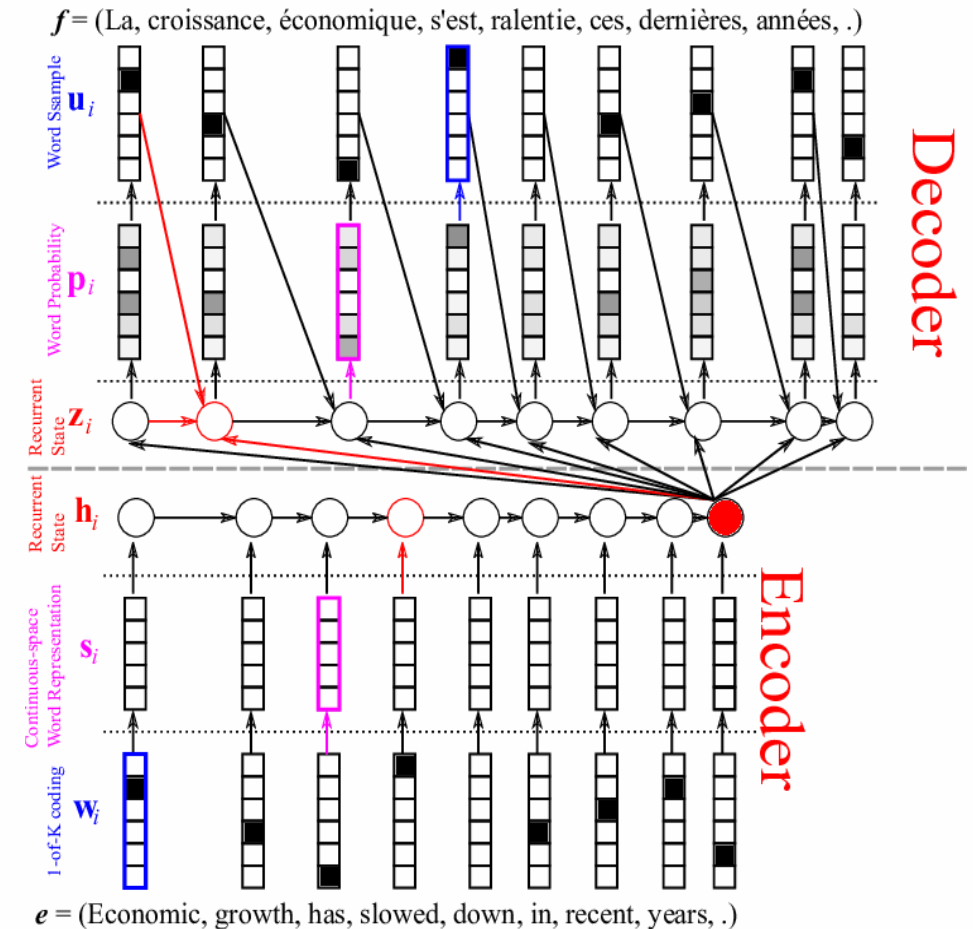
**NLP**

*Nederlandse tekst …*

# Machine Translation

gerade zu diesem Stamm gehören ||| belong just to these families ||| 0.390442 2.9025
gerade zu diesem Stamm gehören ||| belong just to these ||| 0.390442 2.9025e-15 0.2
gerade zu diesem Stamm gehören ||| them belong just to these ||| 0.390442 2.9025e-
gerade zu diesem Stamm ||| belong just to these families ||| 0.390442 1.5563e-11 0.2
gerade zu diesem Stamm ||| belong just to these ||| 0.390442 1.5563e-11 0.260295 9.
gerade zu diesem Stamm ||| them belong just to these ||| 0.390442 1.5563e-11 0.260
gerade zu diesem ||| belong just to ||| 0.390442 1.33531e-08 0.260295 4.9937e-09 2.
gerade zu diesem ||| of them belong just to ||| 0.390442 1.33531e-08 0.260295 4.79
gerade zu diesem ||| them belong just to ||| 0.390442 1.33531e-08 0.260295 4.6096
gerade zu ||| belong just to ||| 0.390442 1.8515e-05 0.260295 4.9937e-09 2.718 |||
gerade zu ||| of them belong just to ||| 0.390442 1.8515e-05 0.260295 4.79751e-13
gerade zu ||| them belong just to ||| 0.390442 1.8515e-05 0.260295 4.60968e-12 2.7



Statistical Machine Translation

Neural Machine Translation
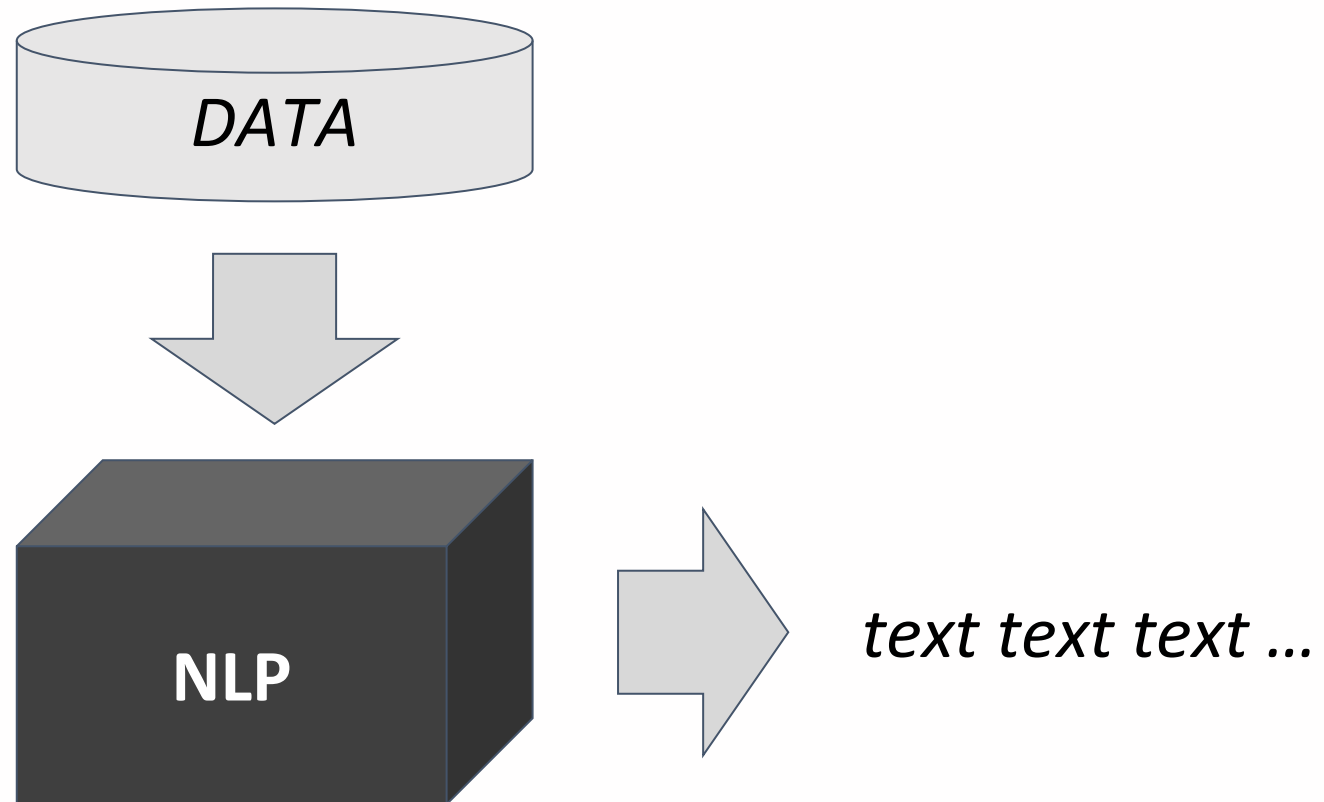
# Language Generation: Data-to-Text

# Data-to-Text Generation



Image by Nivranshu Pasricha

OLLSCOIL NA GAILLIMHE
UNIVERSITY OF GALWAY

44

# Language Generation: Dialog Systems

DATA

**NLP**

*text text text ….*

*text text text …*

# Dialog System Development

# NLP in industry

University *of*Galway.ie

# How Google uses NLP to better understand search queries, content

Learn the role that natural language processing plays in making Google search even more semantic and context-based.
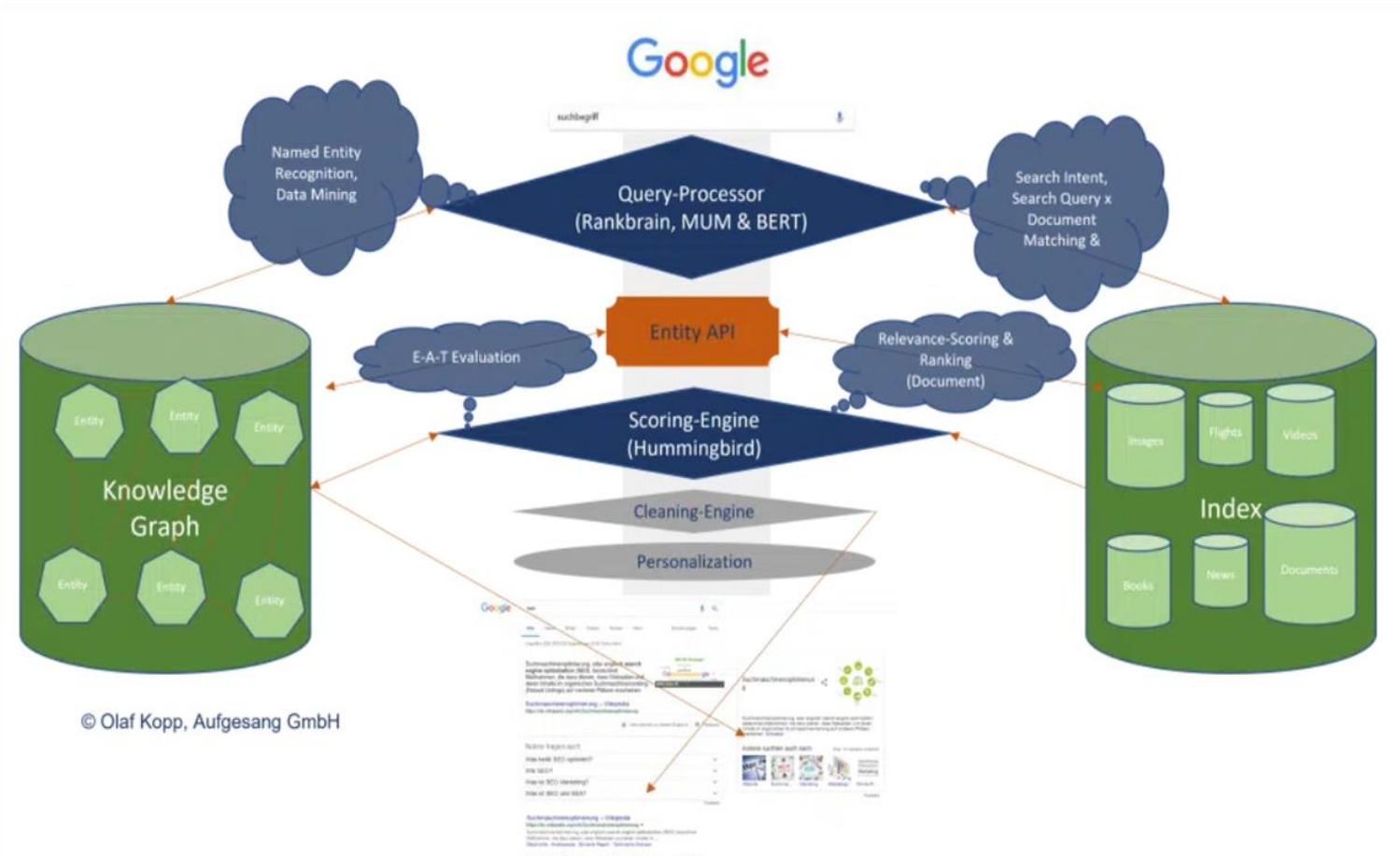
Olaf Kopp on August 23, 2022 at 6:00 am | Reading time: 10 minutes

Natural language processing opened the door for semantic search on Google.

SEOs need to understand the switch to entity-based search because this is the future of Google search.

OLLSCOIL NA GAILLIMHE
UNIVERSITY OF GALWAY

# Google Index & Knowledge Graph



© Olaf Kopp, Aufgesang GmbH

# "Use of NLP in Google Search"



According to Olaf Kopp of Aufgesang GmbH (article in Search Engine Land), Google Search uses NLP for the following:

- Interpretation of search queries.

- Classification of subject and purpose of documents.

- Entity analysis in documents, search queries and social media posts.

- Generating featured snippets and answers in voice search.

- Interpretation of video and audio content.

- Expansion and improvement of the Knowledge Graph.



OLLSCOIL NA GAILLIMHE
UNIVERSITY OF GALWAY

https://searchengineland.com/how-google-uses-nlp-to-better-understand-search-queries-content-387340

Which NLP tasks are needed for each of these steps?

# Interpretation of search queries

*"calories in bar of chocolat"*

- Relation Extraction, Parsing, Part-of-Speech, WSD

*"Biden"*

- Entity Linking

*"great movies" "scary movie"*

- Sentiment / Emotion Analysis

# Classification of subject, purpose of documents

*legal, health, … documents*

- Text Classification, Probability, Language Modeling, Word Embeddings

*documents about concept XYZ*

- Concept Extraction

# Entity analysis in documents, queries and posts

*identify names of people, locations, things*

- NER, Entity Linking, Concept/Taxonomy Extraction

*identify relations between people, locations, things*

- Relation Extraction, Parsing, Semantic Role Labeling

*resolution of pronouns*

- Coreference Resolution

# Generating featured snippets and answers

*summarize one or more retrieved documents*

- Text Summarization

*generate a specific answer*

- Natural Language Generation, Dialog

*translate a text*

- Machine Translation

# Interpretation of video and audio content

*retrieve relevant videos for a search query*

- Multimodal Analysis

*speech interaction*

- Speech-to-Text, Dialog

OLLSCOIL NA GAILLIMHE
UNIVERSITY OF GALWAY

# Expansion/improvement of Knowledge Graph

*include new entities*

- Entity Linking, Concept Extraction, Taxonomy Extraction/Extension

*include new or update existing relations*

- Relation Extraction, Knowledge Graph Completion

# Lab of this week



Intro to base NLP tools and methods

Tools:

- pandas
- Hugging Face Datasets
- scikit-learn
- Hugging Face Transformers

Dataset:

- IMDB reviews: https://huggingface.co/datasets/imdb