

# Information Extraction

Named Entity Recognition, Relation Extraction

Gustavo Gonçalves, Research Assistant @ LTI/CMU and NOVA LINCS  
ggoncalv@xcs.cmu.edu - remove x chars

# About me

- I'm a Research Assistant at the **Web and Media Search Lab @ NOVA LINCS**, and at the **Language and Technology Institute @ Carnegie Mellon University**. I'm also pursuing my Dual PhD Degree. Check the CMU Portugal Program if you need funding
- I've been doing research since 2017, and my PhD topic is centered on improving **conversational search sessions with knowledge-aware representations**.
- My research interests include: Conversational Search; Knowledge-Representations; Machine Learning; and Socia-Media Analysis.

# Lecture Outline

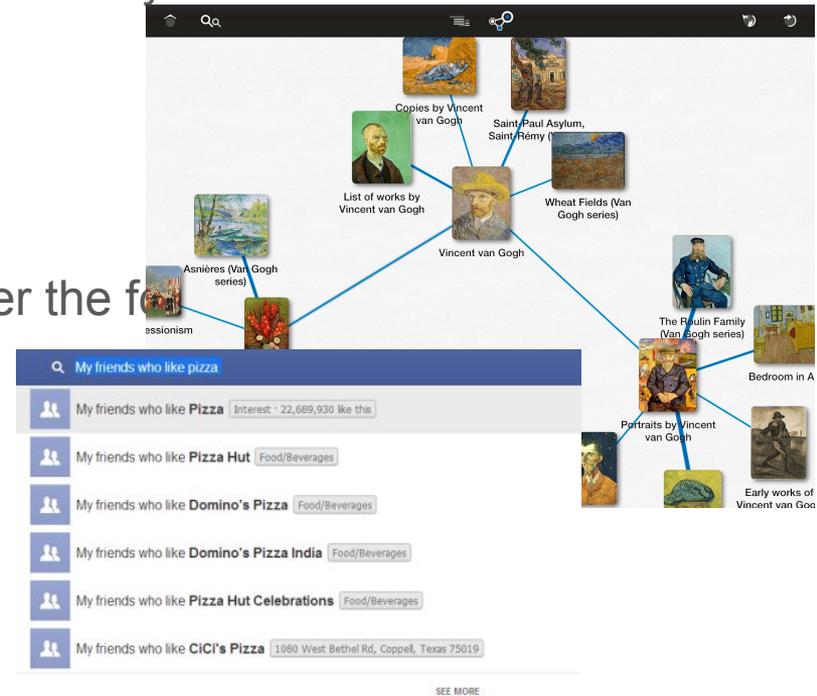
- Knowledge Graphs;
- What is Information Extraction?;
- Named Entity Recognition;
- Entity Linking;
- Knowledge Driven Information Retrieval;

# The story so far

- Relevance based models
  - Based on the query-related documents (initial search results)
- Statistical correlations
  - *Term correlations across documents*
    - *(we will revisit this in the recommendation lecture)*
  - Term correlations across term's neighborhood (word embeddings)
- Knowledge-bases expansion
  - Linguistic thesaurus: e.g. MedLine: physician, syn: doc, doctor, MD, medico
  - Can be query rather than just synonyms

# What are Knowledge Bases?

- Knowledge bases are graphs that have a manually curated source.
  - The most famous example probably is Wikipedia!
- Graph nodes are entities, and edges are relationships among entities;
- These relationships are usually stored under the form of triples
  - Example: Painted(Van Gogh, The Starry Night)



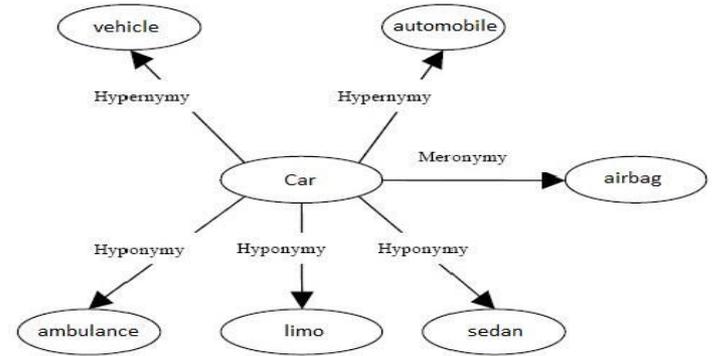
[https://en.wikipedia.org/wiki/Facebook\\_Graph\\_Search](https://en.wikipedia.org/wiki/Facebook_Graph_Search)

<https://www.designnominees.com/apps/learn-discovery-mindmap-of-wikipedia>

# WordNet: A lexical database

“WordNet® is a large lexical database of English. Nouns, verbs, adjectives and adverbs are grouped into sets of cognitive synonyms (synsets), each expressing a distinct concept.

Synsets are interlinked by means of conceptual-semantic and lexical relations.”



<https://wordnet.princeton.edu/>

“WordNet interlinks specific senses of words. As a result, words that are found in close proximity to one another in the network are semantically disambiguated. Second, WordNet labels the semantic relations among words, whereas the groupings of words in a thesaurus does not follow any explicit pattern other than meaning similarity.”

# ImageNet: A visual taxonomy

- Selected words of WordNet are illustrated in ImageNet.

<http://image-net.org/explore.php>

- Currently, there are over 14.000 concepts illustrated.
- Roughly 1.000 concepts are used by VOC.
- Great impact in advancing the state of the art.

The screenshot shows the ImageNet website interface. At the top, the logo 'IMAGENET' is displayed next to a search bar containing the text '14,197,122 images, 21841 synsets selected'. To the right of the search bar is a green 'SEARCH' button. Further right are links for 'Home', 'About', 'Explore', and 'Download'. Below the search bar, the text 'Not logged in. Login | Signup' is visible.

The main content area is titled 'Sport, athletics' with the subtitle 'An active diversion requiring physical exertion and competition'. To the right of this title, it shows '1888 pictures' and '92.64% Popularity Percentile'. There is also a 'WordNet IDs' icon.

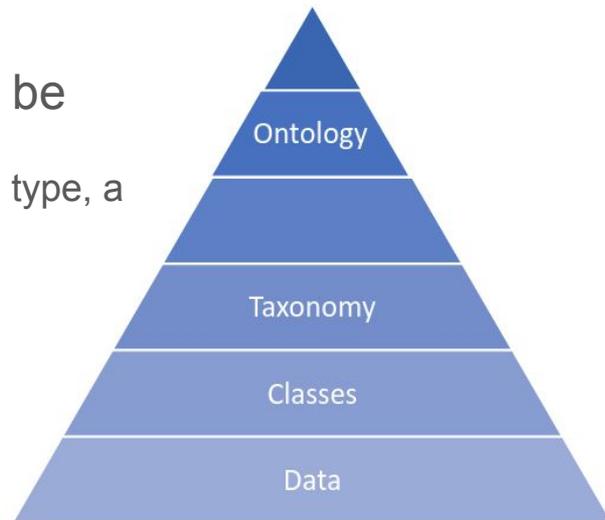
Below the title, there are three tabs: 'Treemap Visualization', 'Images of the Synset', and 'Downloads'. The 'Images of the Synset' tab is active, showing a grid of small images categorized under various sub-headers like 'Athletic', 'Contact', 'Outdoor', 'Water', 'Blood', 'Racing', 'Gymnast', 'Sliding', 'Cycling', 'Team', 'Skating', 'Archery', 'Judo', 'Rowing', 'Funambulism', 'Riding', 'Track', 'Rock', and 'Skiing'.

On the left side of the image grid, there is a list of synsets with their corresponding counts in brackets. The list includes: 'ImageNet 2011 Fall Release (32326)', 'plant, flora, plant life (4486)', 'geological formation, formation (17)', 'natural object (1112)', 'rock, stone (30)', 'asterism (0)', 'carpet (0)', 'black body, blackbody, full radiator (1)', 'consolidation (0)', 'mechanism (12)', 'body, organic structure, physical nest (7)', 'plant part, plant structure (681)', 'body (93)', 'cocoon (0)', 'sample (12)', 'covering, natural covering, cover (2)', 'tangle (2)', 'universe, existence, creation, constellation (0)', 'celestial body, heavenly body (3)', 'body, dead body (6)', 'extraterrestrial object, extraterrestrial (178)', 'sport, athletics (178)', 'rowing, row (2)', 'funambulism, tightrope walking (0)', 'judo (0)', 'blood sport (10)', 'gymnastics, gymnastic exercise (16)', and 'water sport (16)'. The 'sport, athletics (178)' entry is highlighted in blue.

At the bottom of the page, there is a copyright notice: '© 2010 Stanford Vision Lab, Stanford University, Princeton University support@image-net.org Copyright infringement'.

# From data to information

- A taxonomy is concerned with classifying and organizing hierarchically concepts of a specific domain.
- It is important to identify the list of items that need to be detected.
  - These items are domain specific, and can be a topic, a scene type, a visual object or a named entity.
  - They are normally associated to a class in a supervised learning task.



# What is Information Extraction (IE)?

- We want to work with Things, not Strings!
- IE systems extract factual, and clear data;
  - What? Where? Who? When?
- We can roughly think of IE systems as entity recognition and relation extraction machines
  - For Example:

**Tesla, Inc.** (formerly **Tesla Motors, Inc.**) is an American [electric vehicle](#) and [clean energy](#) company based in [Palo Alto, California](#).<sup>[9]</sup>

Location(“Tesla, Inc.”, “Palo Alto, California”)

# Classification, detection, linking



**Detection:**

Statue

**Linking:**

Statue of Liberty

[https://en.wikipedia.org/wiki/Statue\\_of\\_Liberty](https://en.wikipedia.org/wiki/Statue_of_Liberty)

**Classification:**

Sea side

Statue

City

Sky

**Linking:**

New York City

[https://en.wikipedia.org/wiki/New\\_York\\_City](https://en.wikipedia.org/wiki/New_York_City)

# You use IE every day!

- Rule-based and Machine Learning IE systems help us daily with convenient tasks;
- These tasks can go from low-level information extraction, such as with regular expressions;

en.wikipedia.org > wiki > Tesla, Inc

[Tesla, Inc. - Wikipedia](#)

**Tesla**, Inc. (formerly **Tesla** Motors, Inc.) is an American electric vehicle and clean energy **company** based in Palo Alto, California.

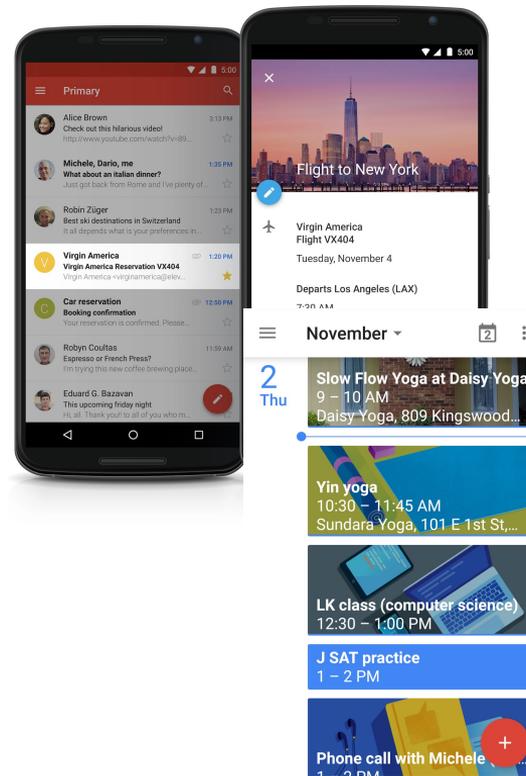
**Products:** Electric vehicles; **Tesla** batteries; Sol... **Total assets:** US\$34.309 billion (2019)

**Number of employees:** 48,016 (2019) **Total equity:** US\$6.618 billion (2019)

[History of Tesla, Inc.](#) · [Tesla Model Y](#) · [Tesla Powerwall](#) · [Tesla Megapack](#)

Google

tesla company



# Named Entity Recognition (NER)

- Is one of the sub-tasks in Information Extraction:  
Identify and classify names in text, for example:
  - The decision by the independent MP **Andrew Wilkie** to withdraw his support for the minority **Labor** government sounded dramatic but it should not further threaten its stability. When, aWer the **2010** election, **Wilkie**, **Rob Oakeshott**, **Tony Windsor** and the **Greens** agreed to support **Labor**, they gave just two guarantees: confidence and supply.

Person  
Date  
Location  
Organization

- We've seen this kind of task in Machine Learning, it's a multi-label classification task. What are the features and evaluation metrics?

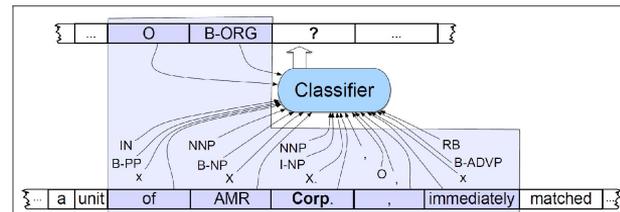
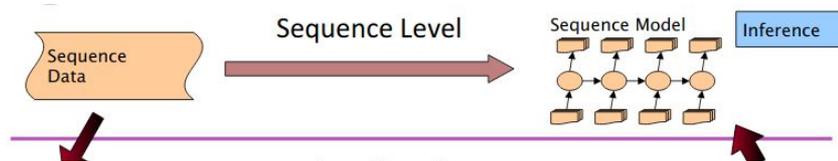


Figure 13.7 Named entity recognition as sequence labeling. The features available to the classifier during training and classification are those in the boxed area.

# Machine Learning approach to NER

## Testing

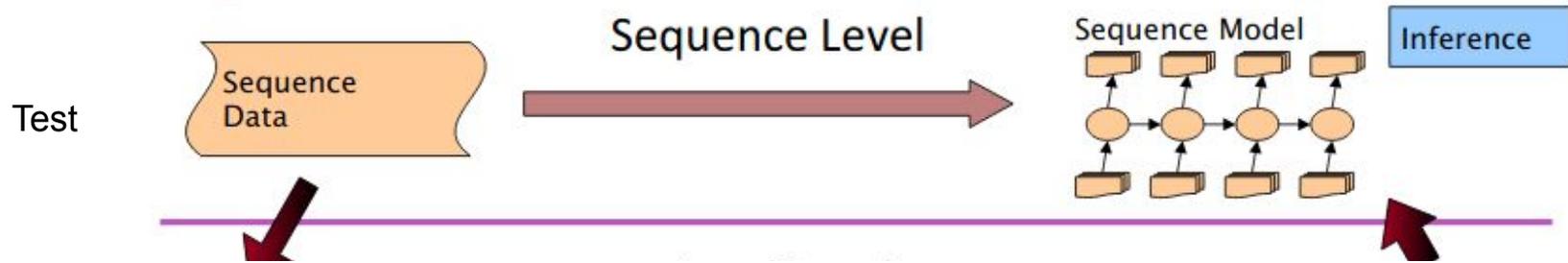
1. Receive a set of testing documents
2. Run **sequence model** inference to label each token
3. Appropriately output the recognized entities



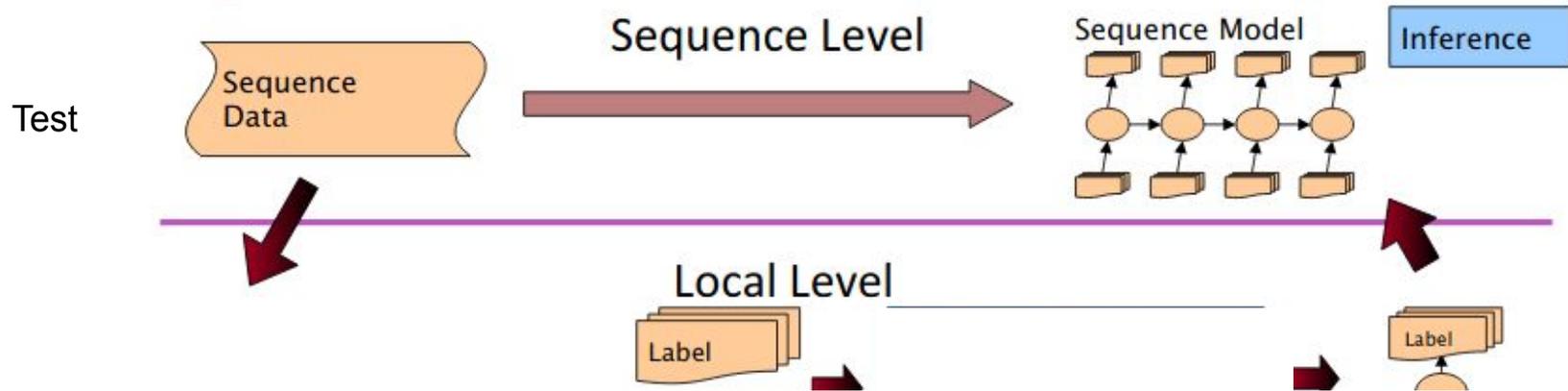
## Training

1. Collect a set of representative training documents
2. Label each token for its **entity class or other (O)**
3. Design feature extractors appropriate to the text and classes
4. Train a **sequence classifier** to predict the labels from the data

# Language is Sequential, and so is IE

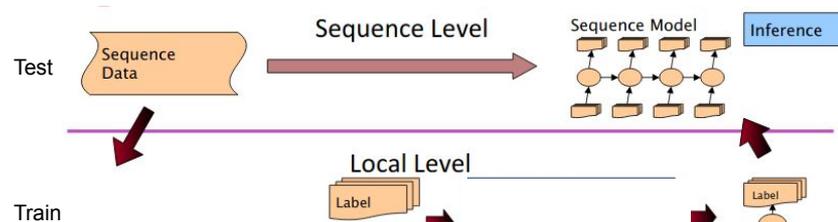


# Language is Sequential, and so is IE

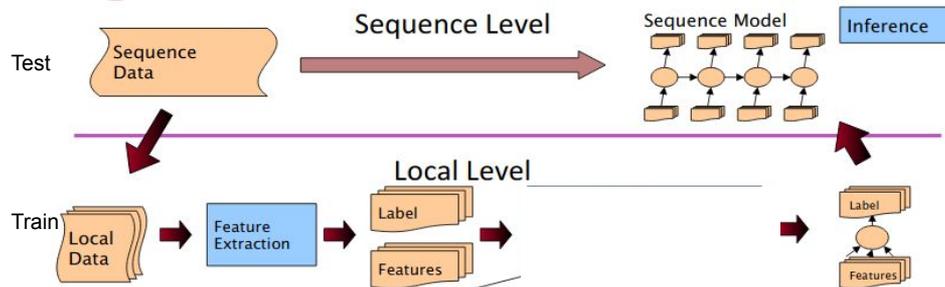


# NER - Classes (Labels)

	IO encoding	IOB encoding
Fred	PER	B-PER
showed	O	O
Sue	PER	B-PER
Mengqiu	PER	B-PER
Huang	PER	I-PER
's	O	O
new	O	O
painting	O	O



# NER - Features



- Words
  - Current words
  - Previous/next word (context)
- Other kinds of inferred linguistic classification
  - Part-of-speech (POS) tags
- Label context
  - Previous (and perhaps next) label
- Word Substrings
  - Croatian surnames (many end with “ić”)
- Word shapes
  - Length, capitalization, numerals, greek letters, word punctuation
  - Examples: SARS-CoV-2 = XXXX-XxX-d / COVID-19 = XXXXX-dd

# Important Sequence Modeling Problems in NLP

- Many problems in NLP have data which is a sequence of characters, words, phrases, lines, or sentences ...
- We can think of our task as one of labeling each item

<b>VBG</b>	<b>NN</b>	<b>IN</b>	<b>DT</b>	<b>NN</b>	<b>IN</b>	<b>NN</b>
Chasing	opportunity	in	an	age	of	upheaval

**POS tagging**

<b>PERS</b>	<b>O</b>	<b>O</b>	<b>O</b>	<b>ORG</b>	<b>ORG</b>
Murdoch	discusses	future	of	News	Corp.

**Named entity recognition**

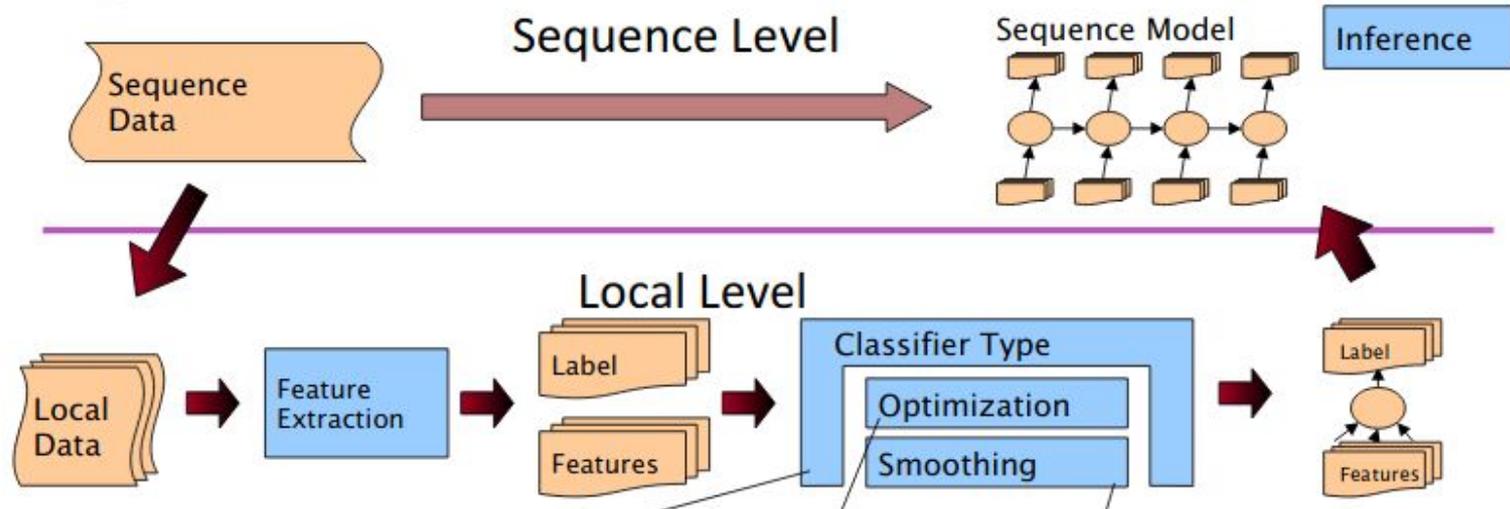
<b>B</b>	<b>B</b>	<b>I</b>	<b>I</b>	<b>B</b>	<b>I</b>	<b>B</b>	<b>I</b>	<b>B</b>	<b>B</b>
而	相	对	于	这	些	品	牌	的	价

**Word segmentation**



**Text segmentation**

# Language is Sequential, and so is IE



# Entity Linking - What is an Entity?

- Depends on who you ask!
  - For simplicity let's follow Balog's[3] definition and consider Named Entities (Persons, Locations, Dates) Real-World objects and Concepts (Emotion, Gaussian Kernel, Peace) as abstract objects.
- Entity Linking is a special useful tool to inject knowledge into other representations. These representations can be raw count based, or vectorial (embeddings)
- What are the elements involved in a ML Entity Linker?
  - NER system to detect mentions;
  - Knowledge Base to relate candidate pairs of entities (triples);
  - Feature Extractor (substrings, words, patterns);
  - Labeled Disambiguation Data (For supervised approaches);

# Disambiguation - How many Ronaldo's do you know?



- Main page
- Contents
- Current events
- Random article
- About Wikipedia
- Contact us
- Donate
- Contribute
- Help
- Learn to edit
- Community portal
- Recent changes
- Upload file
- Tools
- What links here
- Related changes
- Special pages
- Permanent link
- Page information

Not logged in [Talk](#) [Contributions](#) [Create account](#) [Log in](#)

Article [Talk](#)

Read

[View source](#)

[View history](#)



## Ronaldo



From Wikipedia, the free encyclopedia  
(Redirected from [Ronaldo \(disambiguation\)](#))

*For the documentary about the Portuguese footballer, see [Ronaldo \(film\)](#).*

**Ronaldo** is a Portuguese<sup>[1]</sup> name equivalent to the Scottish [Ronald](#). From the football super stars, it became a very common name in all Portuguese speaking countries, being also prevalent in Italy and Spanish speaking countries.

### Contents [hide]

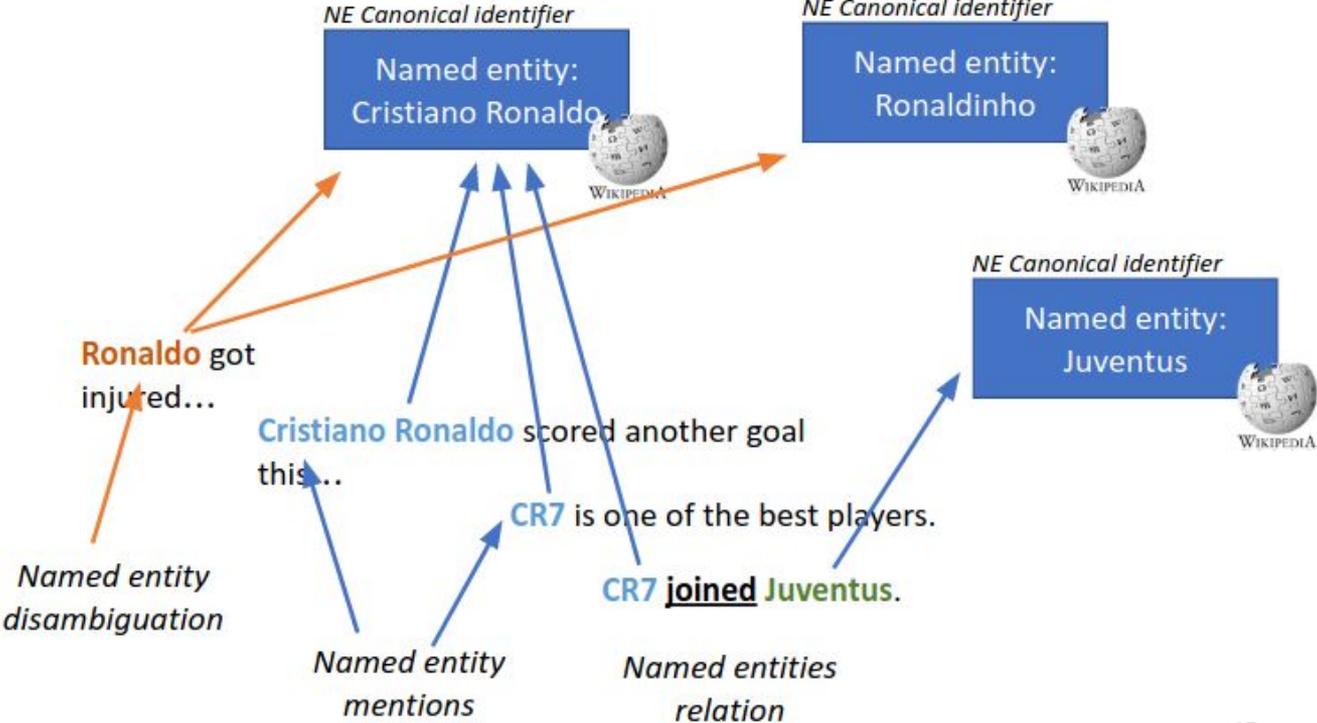
- [People](#)
- [Fictional characters](#)
- [See also](#)
- [References](#)

## People

Notable people known as Ronaldo include:

- [Ronaldo \(Brazilian footballer\)](#) (born 1976), Ronaldo Luís Nazário de Lima, was known as "Ronaldinho" in his early career to distinguish with Ronaldo Rodrigues de Jesus
- [Cristiano Ronaldo](#) (born 1985), Portuguese international footballer
- [Ronaldinho](#), full name Ronaldo de Assis Moreira (born 1980), Brazilian international footballer, also known as "Ronaldinho Gaúcho"

# Entity Disambiguation Example



# Knowledge Driven Information Retrieval

- Spoiler alert, all of you interact with entity retrieval daily.
- Entity cards are displayed by identifying and generating tabular and concise information about the entities contained in the query.
- The hypothesis behind this is that entities are important words that should be identified and treated a bit differently when possible. Users seem to enjoy it!



All-time (List of Grand Slam men's singles champions)

Rank	Player	Total	Years
1	Roger Federer	20	2003–2018
2	Rafael Nadal	17	2005–2018
3	Pete Sampras	14	1990–2002
3	Novak Djokovic	14	2008–2018
5	Roy Emerson	12	1961–1967

Annotations: (A) Column Type Identification: 'Person' below the table. (B) Entity Linking: 'http://dbpedia.org/page/Rafael\_Nadal' pointing to the Nadal row. (C) Relation extraction: '<Peter\_Sampras, careerYears, 1990-2002>' pointing to the Sampras row.

Figure 3: Illustration of table interpretation: (A) Column Type Identification. (B) Entity Linking. (C) Relation extraction.

# Use case - Medical Search

- Domain specific terminologies are curated by domain experts and are designed with specific tasks and workflows in mind.
- In the medical domain, the SNOMED-CT is intended to describe medical conditions, procedures, admin, etc. <http://browser.ihtsdotools.org/>

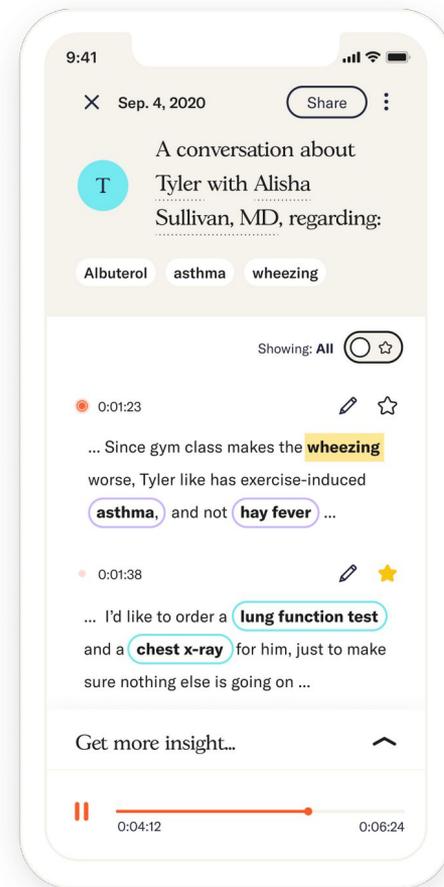
# Use case - Medical Search

- Domain specific terminologies are curated by domain experts and are designed with specific tasks and workflows in mind.
- In the medical domain, the SNOMED-CT is intended to describe medical conditions, procedures, admin, etc. <http://browser.ihtsdotools.org/>



# Solving real world problems

- We are only getting to the good stuff!
- By detecting patterns and relating information we are able to make powerful inferences;
- This app is an example of product that uses IE and very possibly taxonomies to create a network of concepts that are discussed during a doctor's appointment.



# Summary

- In this class we discussed a complementary approach to NLP and IR using IE.
- You now have under your belt another tool that can be used in your work.
- This completes the third possible representation when working with NLP, the three representations that we have learned so far are:
  - Count Based Language Models (Using the Vector Space Model);
    - Coarse
    - Cheap
  - Co-occurrence based language models (word embeddings);
    - Complete
    - Noisy
  - Manually curated knowledge bases (at the very least a controlled vocabulary);
    - Sparse
    - Reliable

# Paper and Book References

[1] Daniel Jurafsky & James H. Martin. Speech and Language Processing. Information Extraction - Chapter 18. 2019

<https://web.stanford.edu/~jurafsky/slp3/18.pdf>

[2] Shuo Zhang and Krisztian Balog. Web Table Extraction, Retrieval, and Augmentation: A Survey. ACM Trans. Intell. Syst. Technol. 2020

<https://arxiv.org/pdf/2002.00207.pdf>

[3] Krisztian Balog. Entity Retrieval. Encyclopedia of Database Systems, Second Edition. 2018 <https://link.springer.com/book/10.1007%2F978-3-319-93935-3>