



NUI Galway  
OÉ Gaillimh

# Introduction to NLP

## Introduction

Dr. Paul Buitelaar  
Data Science Institute, NUI Galway



# Learning Outcomes of this Course

Understand basic principles in automatic analysis, interpretation and transformation of textual data

Understand approaches and challenges in syntactic and semantic analysis of textual data

Understand the foundations of NLP in linguistics, statistics and machine learning

Get an insight into NLP applications, in particular information extraction, knowledge graphs and opinion mining

Awareness of ethical, legal and data privacy aspects of NLP

# Learning Outcomes of This Lecture

Understand the relevance of NLP applications in society

Get an insight into levels of complexity in human language

Get an insight into state of the art directions in NLP

Understand core aspects of NLP architectures and applications

Get useful pointers to data, research and communities in the field

# Overview

Organizational issues

Why NLP

Human language

Some core notions in NLP

Architectures

Applications

Useful pointers



# Overview

## Organizational issues

Why NLP

Human language

Some core notions in NLP

Architectures

Applications

Useful pointers



# Introduction to NLP - Course Overview

- 01 Introduction
- 02 Linguistic Concepts
- 03 Vector Space Model & Word Embeddings
- 04 Language Modelling
- 05 Tagging & Hidden Markov Models
- 06 Probabilistic Parsing
- 07 Semantic Analysis
- 08 Information Extraction
- 09 Knowledge Graphs & Chatbots
- 10 Opinion Mining, Ethics & Data Privacy
- 11 Invited Industry Talk
- 12 Summary and Q&A

# The Team

Paul Buitelaar & John McCrae - lectures



Omnia Zayed - labs

Nivranshu Pasricha - labs



Priya Rani - labs

# Recommended Reading

## Lectures

Daniel Jurafsky, James H. Martin. *Speech and Language Processing - An Introduction to Natural Language Processing, Computational Linguistics and Speech Recognition*. Pearson Prentice Hall. 3<sup>rd</sup> edition draft  
<https://web.stanford.edu/~jurafsky/slp3/ed3book.pdf>

## Labs

Steven Bird, Ewan Klein, and Edward Loper. *Natural Language Processing with Python – Analyzing Text with the Natural Language Toolkit*. O'Reilly. Contents at <http://www.nltk.org/book/>

## Other

Yoav Goldberg. *Neural network methods for natural language processing*. Synthesis Lectures on Human Language Technologies. Morgan & Claypool -- chapters 6,7,8

Chris Manning and Hinrich Schuetze. *Foundations of Statistical Natural Language Processing*. MIT Press. Cambridge. Teaching materials at <http://nlp.stanford.edu/fsnlp/>

# Overview

Organizational issues

**Why NLP**

Human language

Some core notions in NLP

Architectures

Applications

Useful pointers



# What is Natural Language Processing

scientific study and technology development concerned with

**algorithms, methods and systems**

for the

**automatic analysis, interpretation and transformation of human language**

# Why Natural Language Processing

## Search Engines

Process web queries and retrieve relevant information

## Machine Translation

Process input words or sentences to be translated

## Dialog Systems (Chatbots / Digital Assistants)

Process speech and interpret it into task-specific representations and actions

## Other applications: Opinion Mining, Knowledge Graphs

# Search Engines

Google trump

All News Images Videos Maps More Settings Tools

About 264,000,000 results (0.57 seconds)

Top stories



Energy agency rejects Trump plan to prop up coal and nuclear power plants  
The Guardian  
10 hours ago



Russia probe: Trump lawyers 'in talks over Mueller interview'  
BBC.com  
13 hours ago



Trump's plan to help his rural base? Just Obama's leftovers  
Salon  
1 hour ago

→ More for trump

Donald J. Trump (@realDonaldTrump) · Twitter  
<https://twitter.com/realDonaldTrump> 



We are fighting for our farmers, for our country, and for our GREAT AMERICAN FLAG. We want our flag respected - and we want our NATIONAL ANTHEM respected also!  
pic.twitter.com/16eOLXg...



In every decision we make, we are honoring America's PROUD FARMING LEGACY. Years of crushing taxes, crippling regs, & corrupt politics left our communities hurting, our economy stagnant, & millions of...



We have been working every day to DELIVER for America's Farmers just as they work every day to deliver FOR US. #AFBF18  
pic.twitter.com/QDH7fvF...



Donald Trump  45th U.S. President

Donald John Trump is the 45th and current President of the United States, in office since January 20, 2017. Before entering politics, he was a businessman and television personality. Trump was born and grew up in the New York City borough of Queens. [Wikipedia](#)

Born: June 14, 1946 (age 71), Jamaica Hospital Medical Center, New York City, New York, United States

Height: 1.88 m

Net worth: 3.1 billion USD (2017) [Forbes](#)

Spouse: Melania Trump (m. 2005), Marla Maples (m. 1993–1999), Ivana Trump (m. 1977–1992)

Education: Wharton School of the University of Pennsylvania (1968), [MORE](#)

Quotes  View 7+ more

*What separates the winners from the losers is how a person reacts to each new twist of fate.*

*All of the women on 'The Apprentice' flirted with me — consciously or unconsciously. That's to be expected. A sexual dynamic is always present between people, unless you are asexual.*

*Sometimes by losing a battle you find a new way to win the war.*

# Search Engines



## US Presidents list - PresidentsUSA.net

<https://www.presidentsusa.net/presvplist.html> ▾

Listing of Presidents in order and their terms in office and Vice Presidents of the United States.

Political Parties of the Presidents · Thomas Jefferson · Donald Trump · George Bush

## List of Presidents of the United States - Wikipedia

[https://en.wikipedia.org/wiki/List\\_of\\_Presidents\\_of\\_the\\_United\\_States](https://en.wikipedia.org/wiki/List_of_Presidents_of_the_United_States) ▾

Of the elected presidents, four died in office of natural causes (William Henry Harrison, Zachary Taylor, Warren G. Harding, and Franklin D. Roosevelt), four were assassinated (Abraham Lincoln, James A. Garfield, William McKinley and John F. Kennedy), and one resigned (Richard Nixon).

[List of Presidents of the United ...](#) · Lifespan timeline of Presidents ...

## Presidents of the United States (POTUS)

[www.ipl.org/div/potus/](http://www.ipl.org/div/potus/) ▾

They are listed at the bottom of the page. George Washington, 1789-1797. John Adams, 1797-1801. Thomas Jefferson, 1801-1809. James Madison, 1809-1817. James Monroe, 1817-1825. John Quincy Adams, 1825-1829. Andrew Jackson, 1829-1837. Martin Van Buren, 1837-1841.



# Search Engines

Google how many calories in a banana

All Images News Videos Shopping More Settings Tools

About 28,400,000 results (0.49 seconds)

Banana / Energy Amount

89 calories

Type Quantity  
Bananas 100 grams

Sources include: USDA Feedback

People also ask

Are bananas fattening for you?  
How bad are bananas for you?  
How many calories should you eat in a day?  
Is a banana good for you?

Feedback

Can You Eat Bananas If You Want to Lose Weight? | LIVESTRONG.COM  
<https://www.livestrong.com> › Food and Drink ▾  
At 105 calories, a medium banana has 0 cholesterol and just .4 grams of fat.

**Banana**  
Fruit



The banana is an edible fruit – botanically a berry – produced by several kinds of large herbaceous flowering plants in the genus Musa. In some countries, bananas used for cooking may be called plantains, in contrast to dessert bananas. [Wikipedia](#)

**Nutrition Facts**  
Bananas

Amount Per 100 grams	% Daily Value*
Calories 89	
Total Fat 0.3 g	0%
Saturated fat 0.1 g	0%
Polyunsaturated fat 0.1 g	
Monounsaturated fat 0 g	
Cholesterol 0 mg	0%
Sodium 1 mg	0%

# Machine Translation

The screenshot shows the Google Translate interface. On the left, the input text is "it is extremely easy to translate correctly". On the right, the translated text is "Es ist extrem einfach X korrekt zu übersetzen". A red 'X' mark is placed over the word "korrekt" in the German sentence, indicating an error. Below the text boxes are various interactive icons like microphones and keyboards.

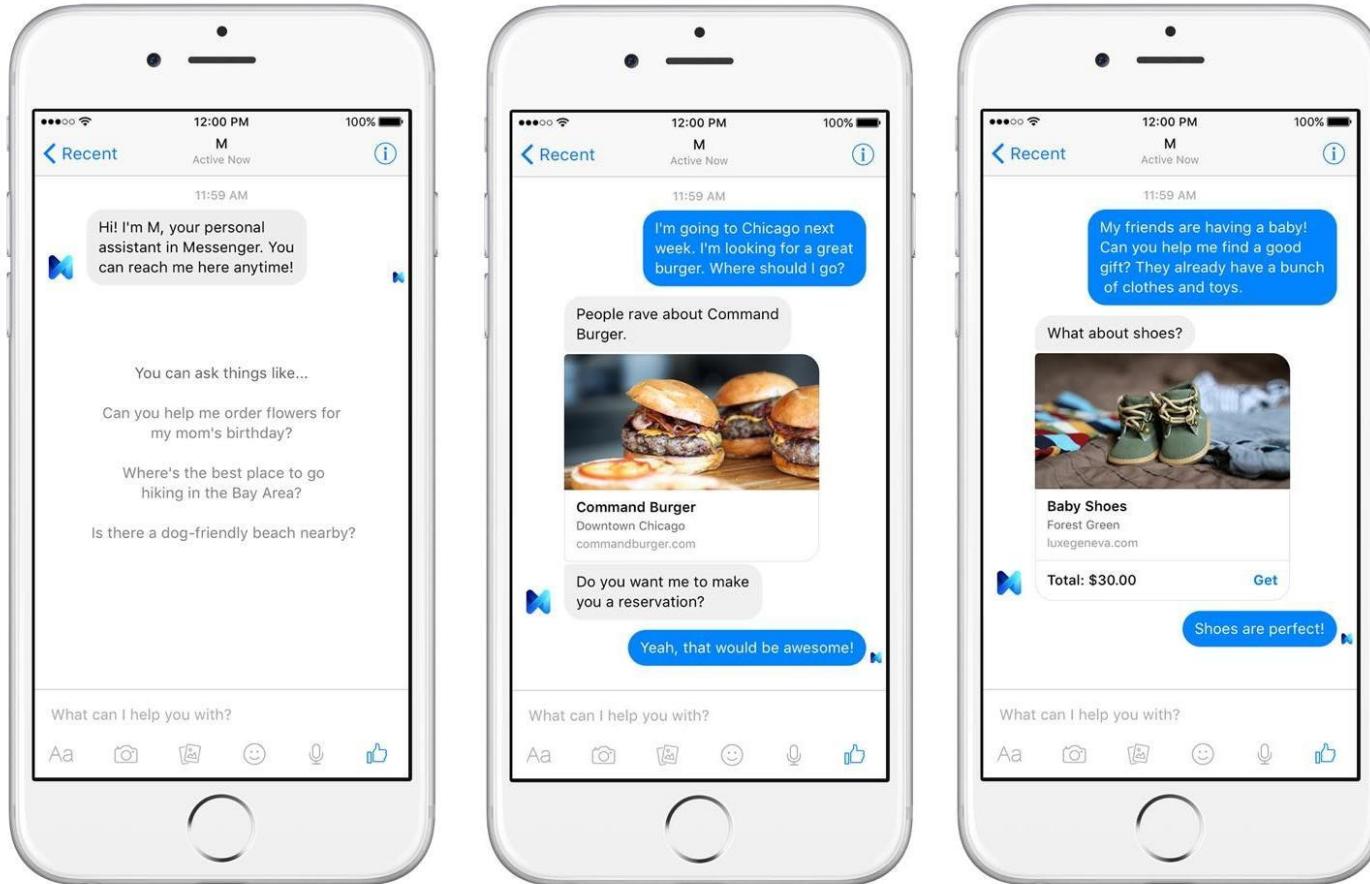
The screenshot shows the Google Translate interface. On the left, the input text is "It is extremely easy to translate correctly". On the right, the translated text is "Es ist extrem einfach, korrekt zu übersetzen". Below the text boxes are various interactive icons like microphones and keyboards. At the bottom right, there is a link "Suggest an edit".

# Chatbots



NUI Galway  
OÉ Gaillimh

# Chatbots - Dialog Systems



# Chatbots - Dialog Systems



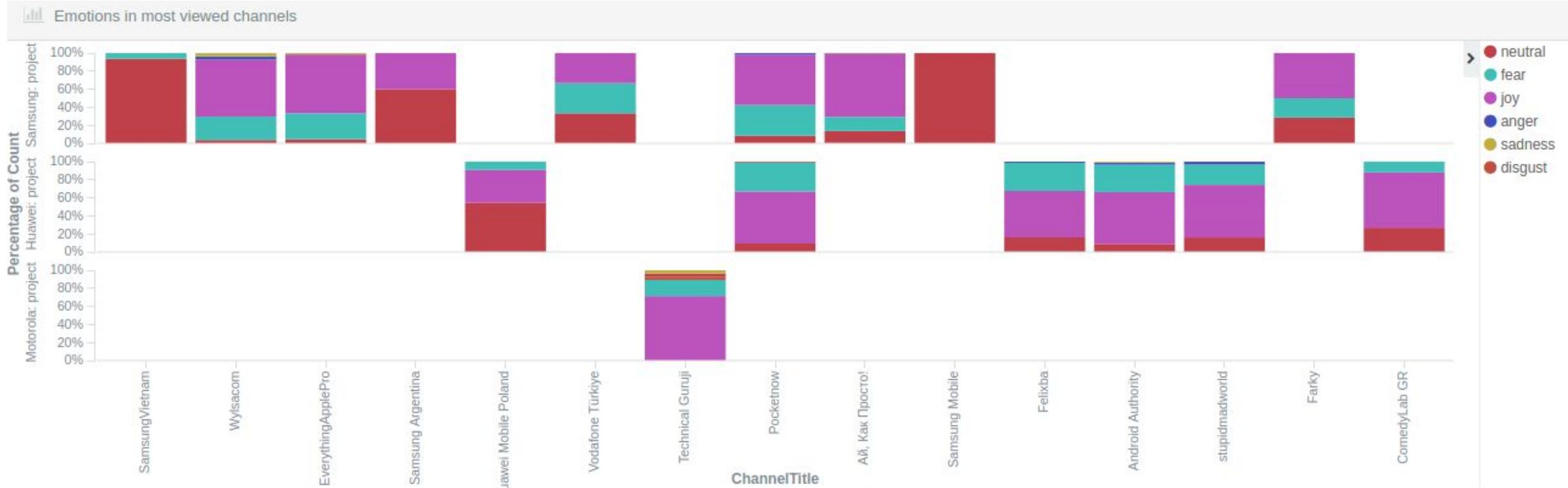
# Opinion Mining



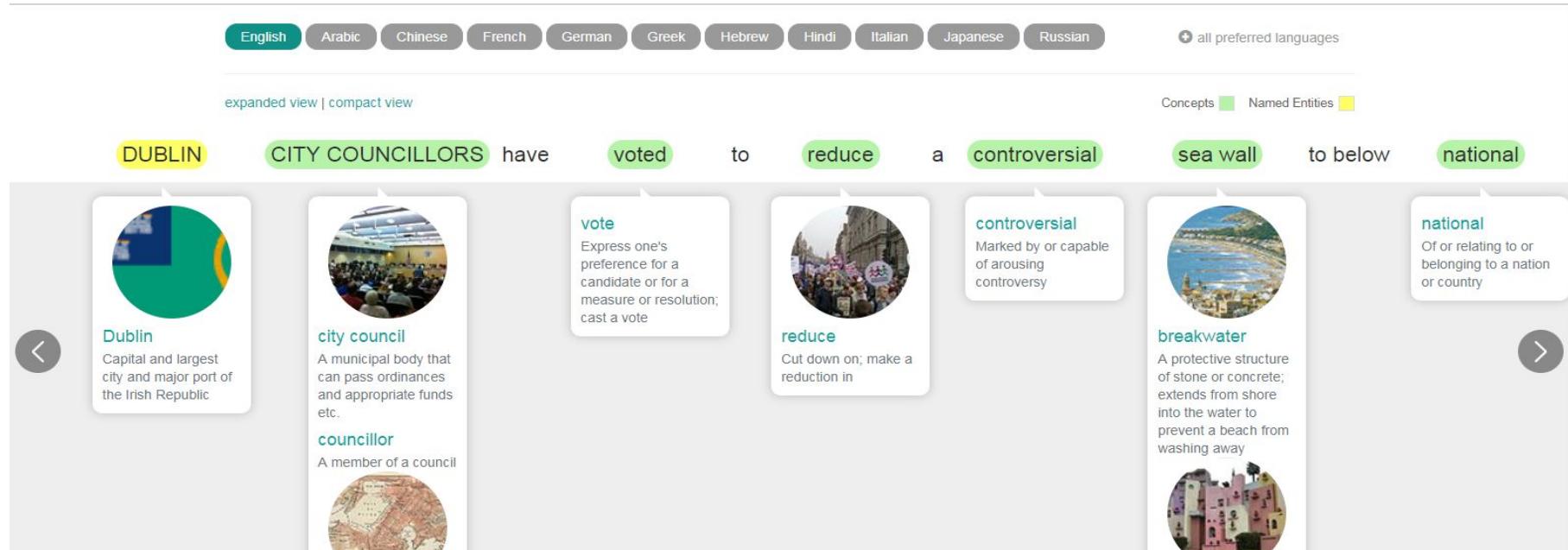
- ▶ Text analytics is now a requirement for large companies
- ▶ Multiple sources of unstructured data
- ▶ "Voice of the Customer" on steroids
- ▶ Ability to understand why NPS metrics are going up or down
- ▶ Ability to respond quickly to service issues
- ▶ Compelling ROI for recent deployments



# Opinion Mining



# Knowledge Graphs



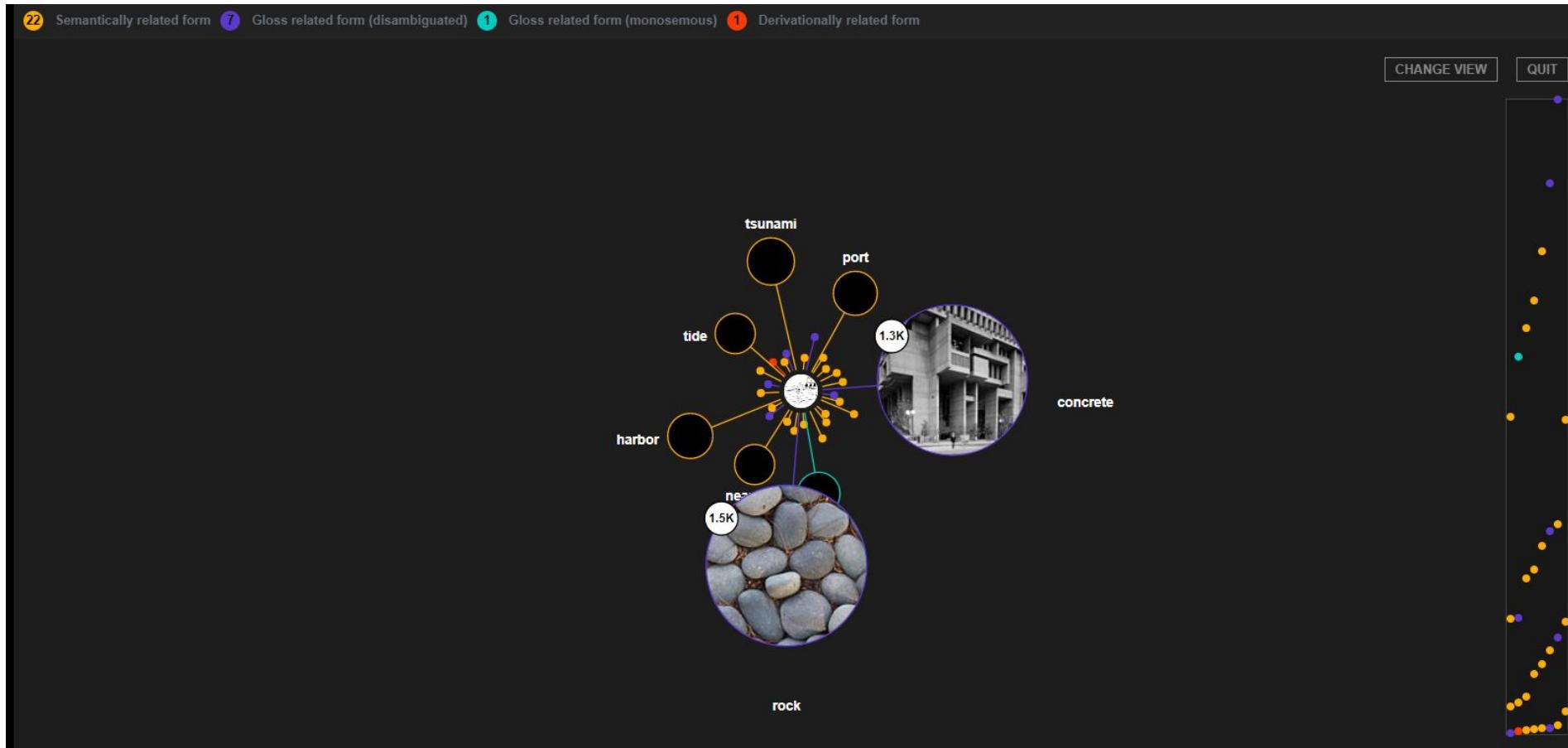
Babelfy



NUI Galway  
OÉ Gaillimh

<http://babelfy.org/>

# Knowledge Graphs



# Overview

Organizational issues

Why NLP

**Human language**

Some core notions in NLP

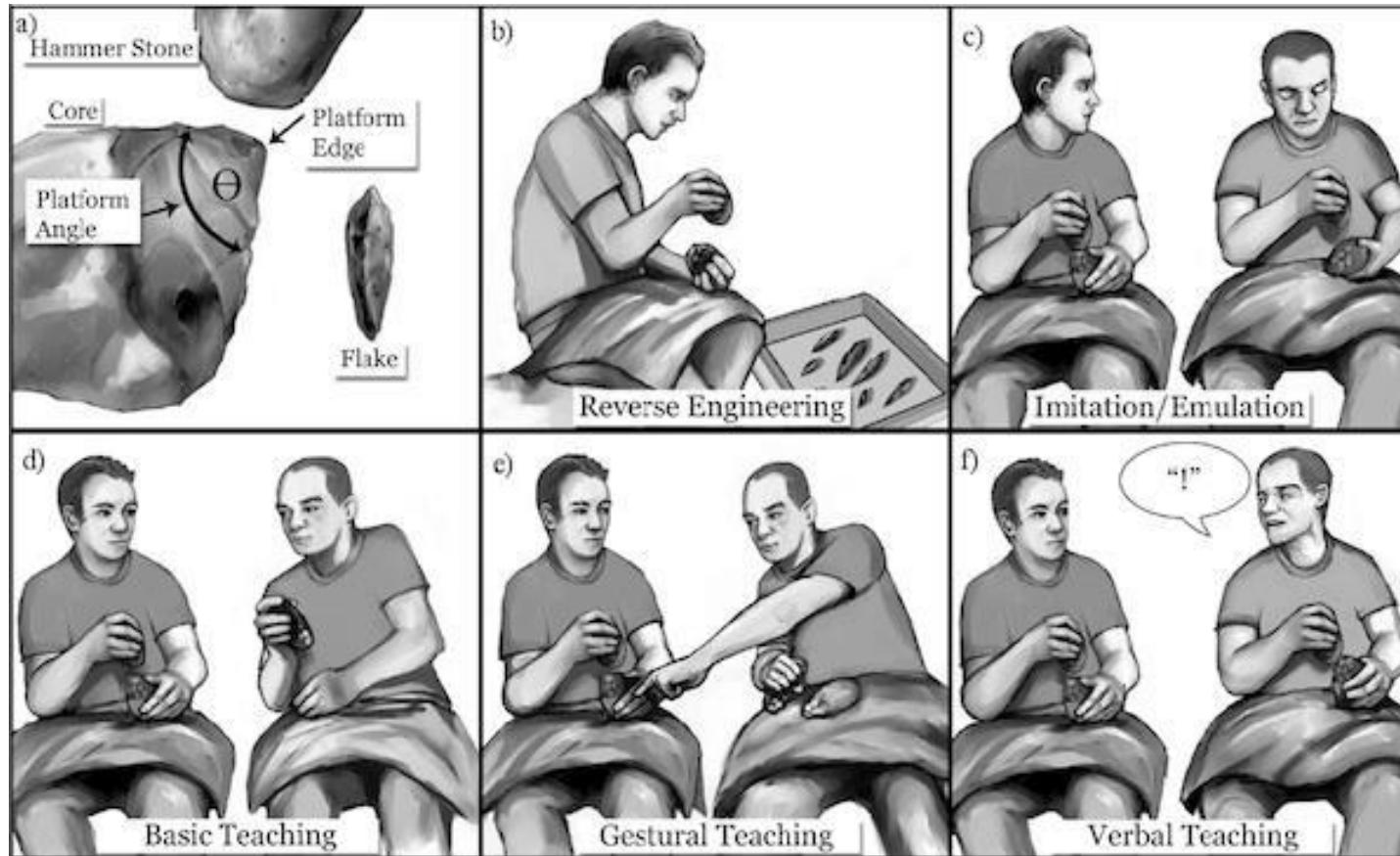
Architectures

Applications

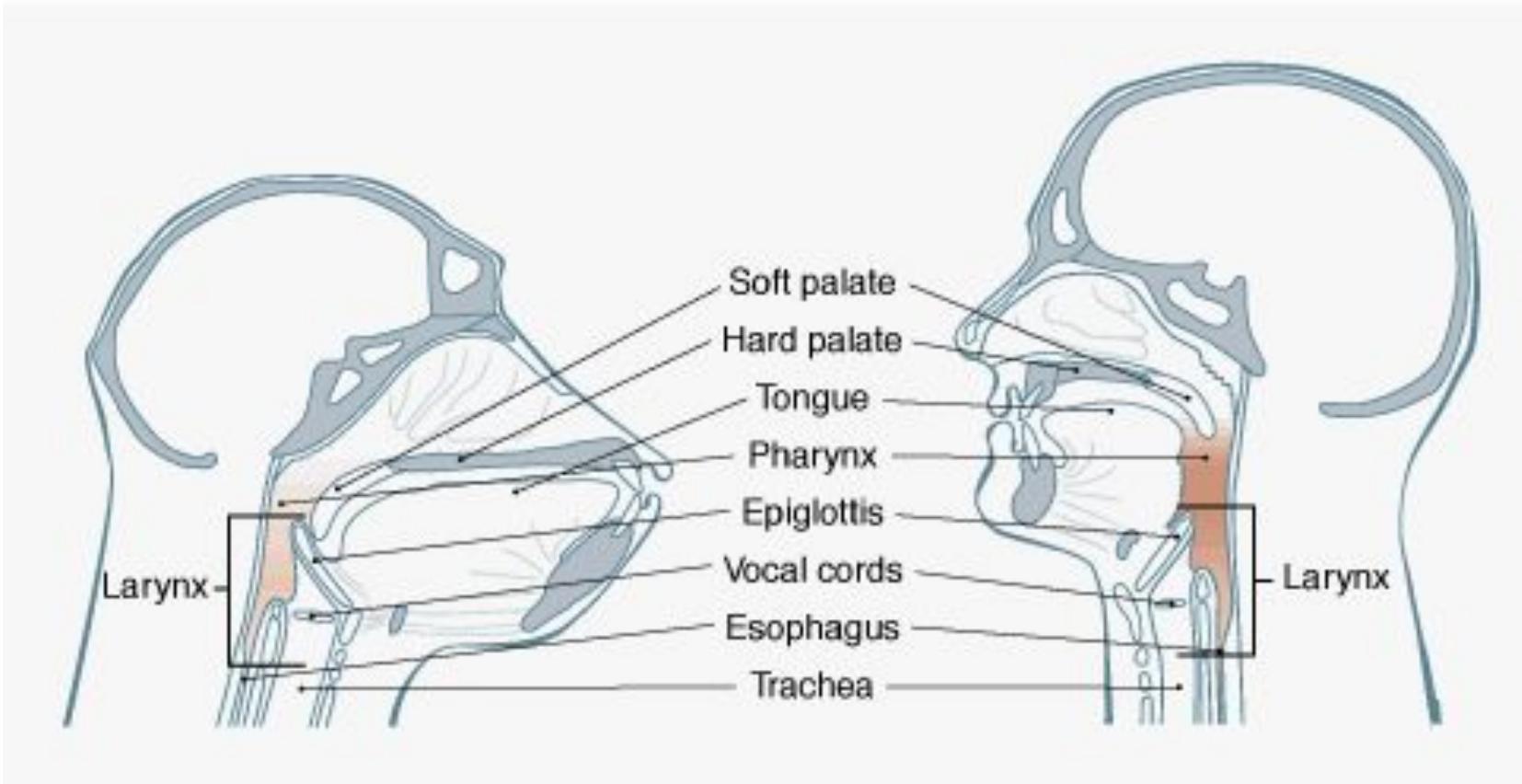
Useful pointers



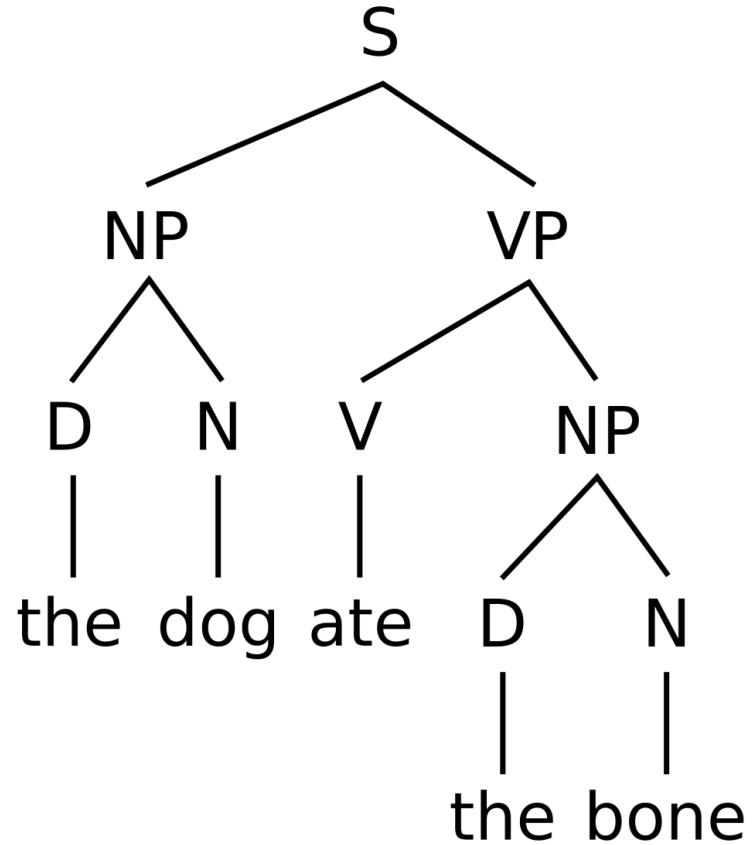
# Human Language as a Social Tool



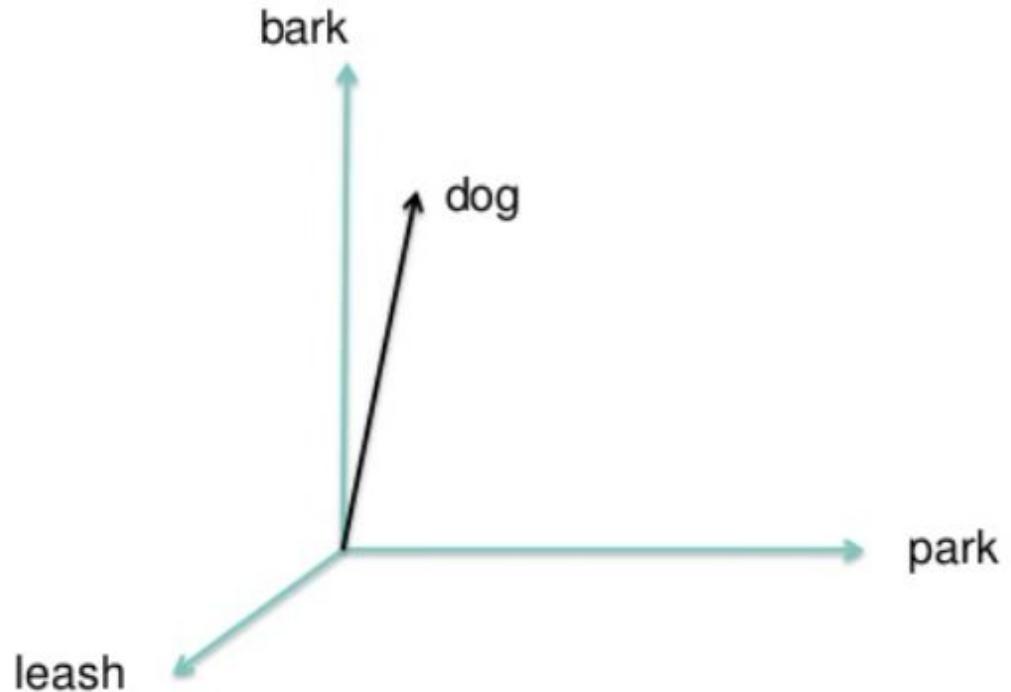
# Human Language as a Biological System



# Human Language as a Logical System



# Human Language as a Mathematical System



# Human Language and Politics



NUI Galway  
OÉ Gaillimh

# Language Families

NLP applications may be easier to transfer within language families



Khoisan

Niger-Kordofanian

Nilo-Saharan

Afro-Asiatic

Dravidian

Kartvelian

Eurasian:

Indo-European

Uralic

Altaic

Korean-Japanese-Ainu

Gilyak

Chukchi-Kamchatkan

Eskimo-Aleut

Dene-Caucasian

Austric

Indo-Pacific

Australian

Amerind

# Overview

Organizational issues

Why NLP

Human language

**Some core notions in NLP**

Architectures

Applications

Useful pointers



# Ambiguity

ambiguity is at the core of human language complexity

human language is

**non-deterministic:** *one word/phrase expressing different meanings*

**redundant:** *different words/phrases expressing the same meaning*

therefore **ambiguous**

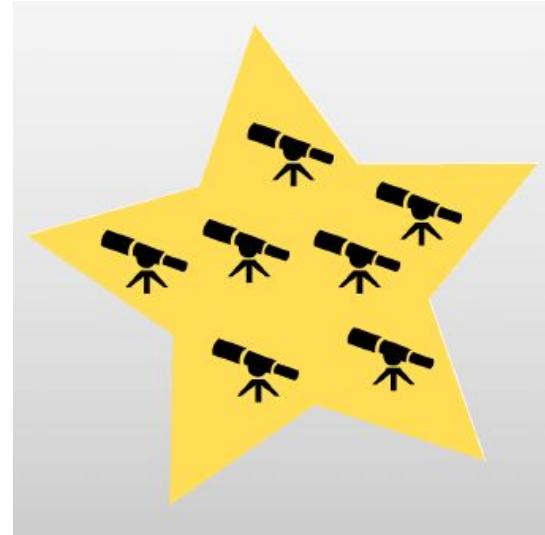


# Lexical Ambiguity



*savings bank* vs. *river bank*

# Syntactic Ambiguity



*astronomers saw stars with telescopes*



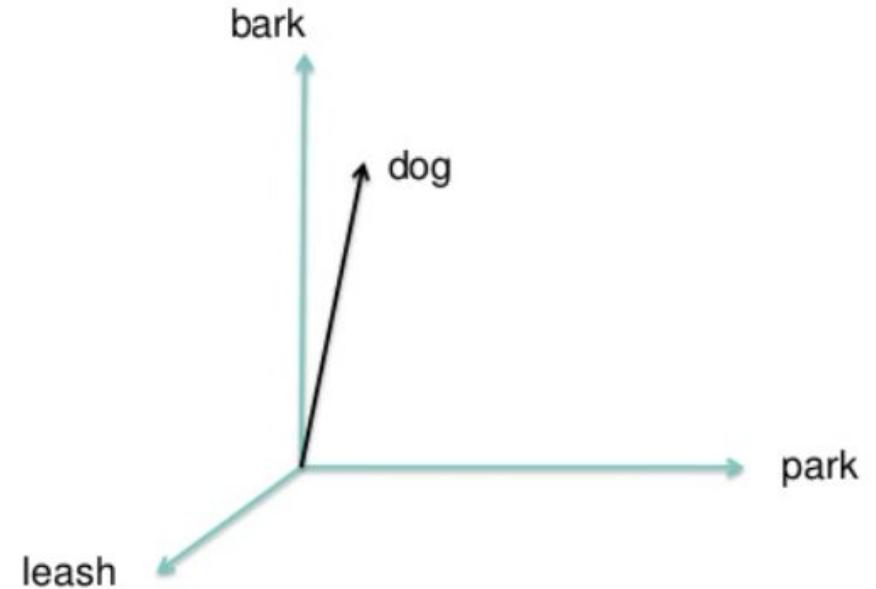
NUI Galway  
OÉ Gaillimh

# Vector Space Model & Word Embeddings

**Distributional Semantics:** words are similar in meaning if they share similar contexts

Represent words using **vectors (embeddings)**

Based on **language modeling**



# Language Modeling

**Frequency and co-occurrence** are central to language modelling

Example: n-grams from the Corpus of Contemporary American English

frequency	word1	word2	word3
31891	much	of	the
13261	much	of	a
8000	much	more	than
7396	much	as	i
5650	much	the	same
5633	much	of	it
4229	much	better	than
4191	much	as	the



# Overview

Organizational issues

Why NLP

Human language

Some common notions in NLP

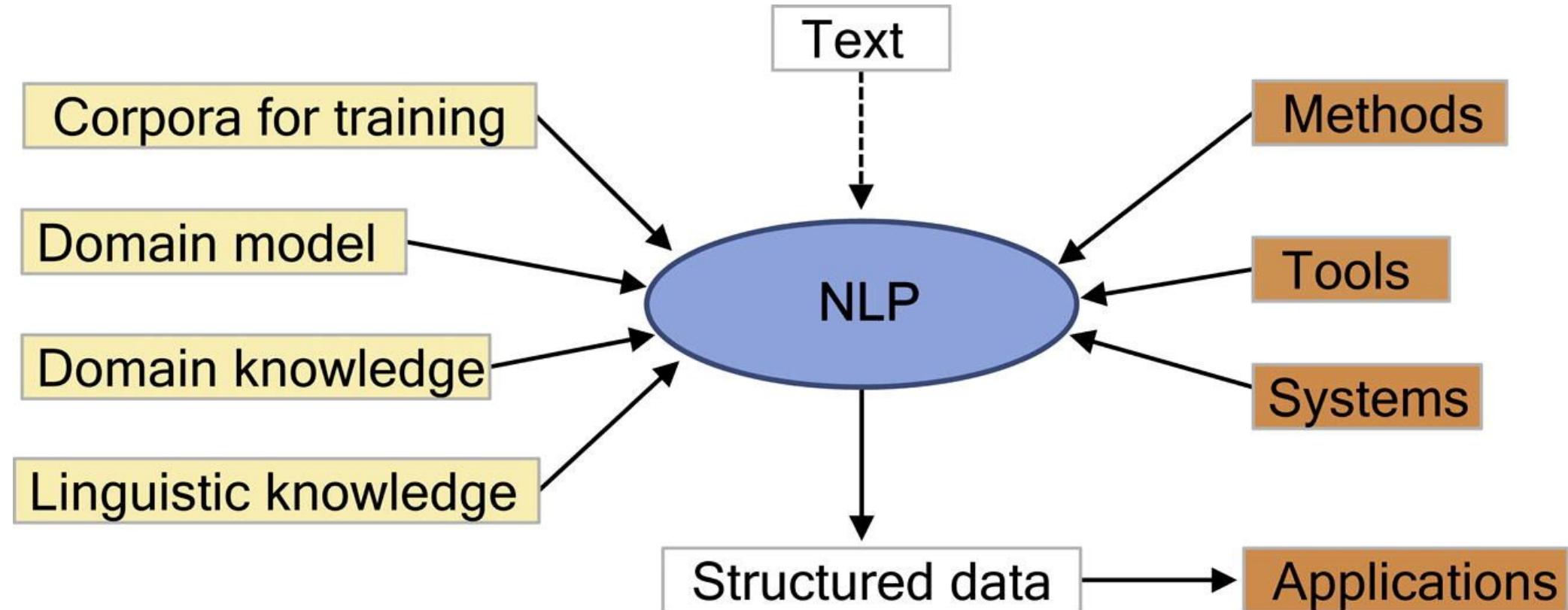
**Architectures**

Applications

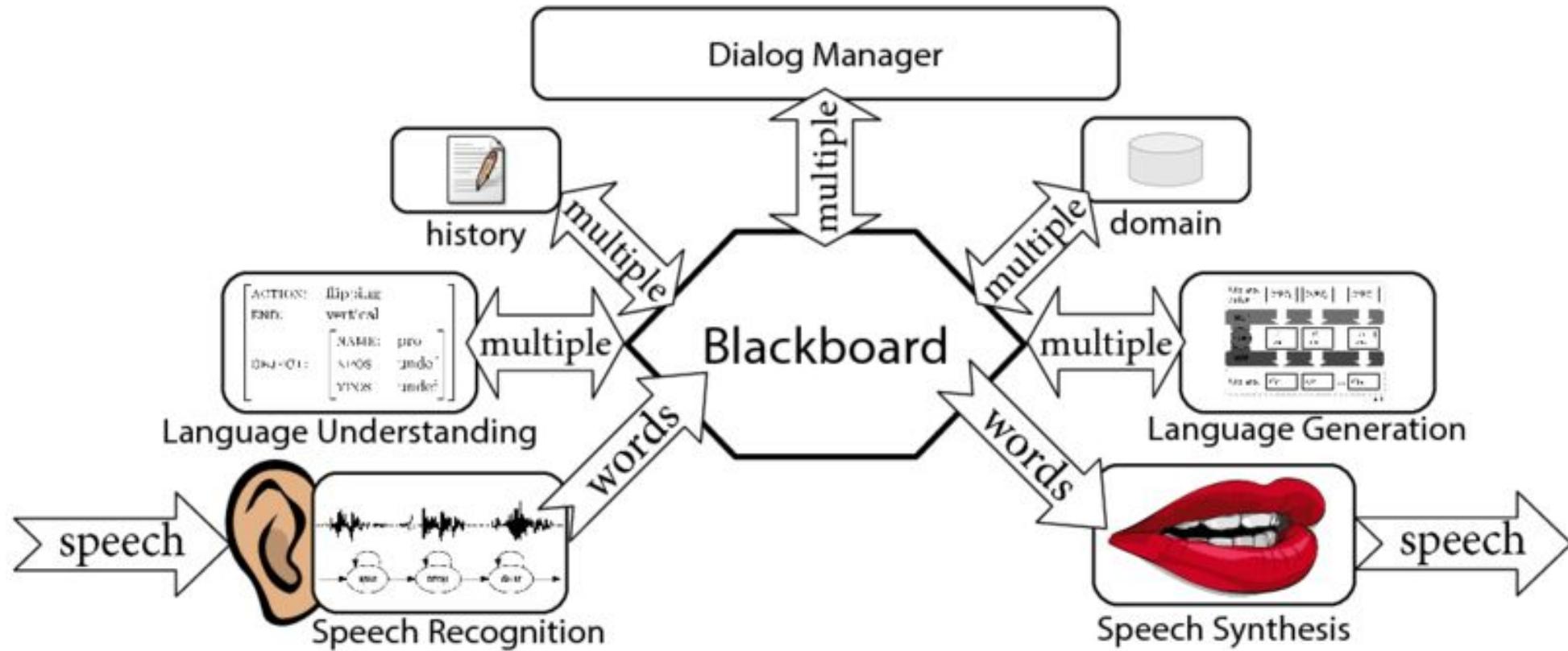
Useful pointers



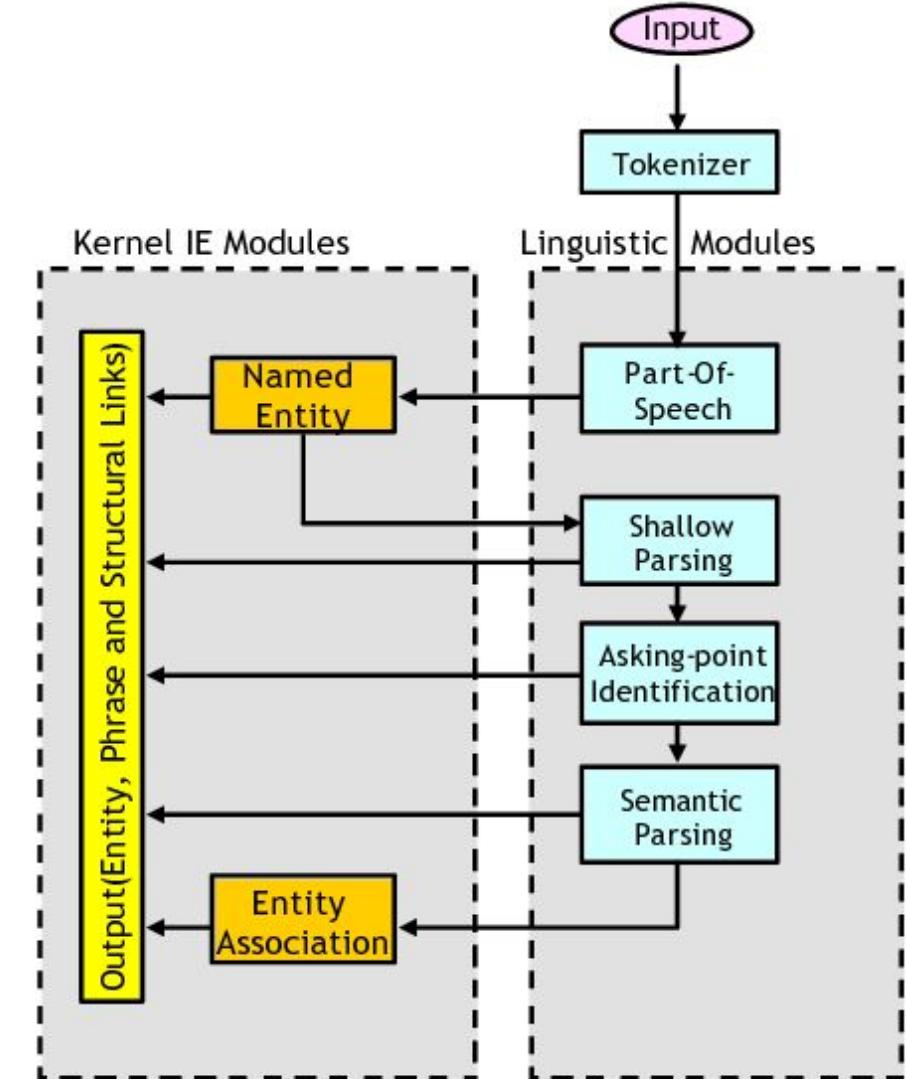
# NLP Architecture - General Overview



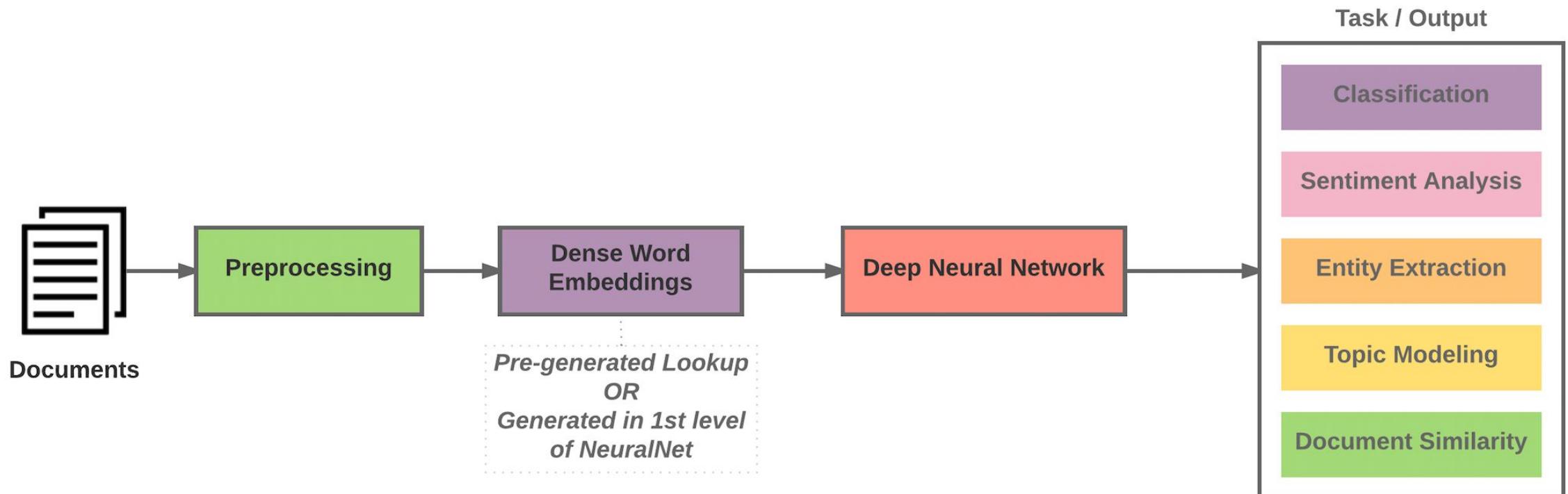
# Blackboard Architecture



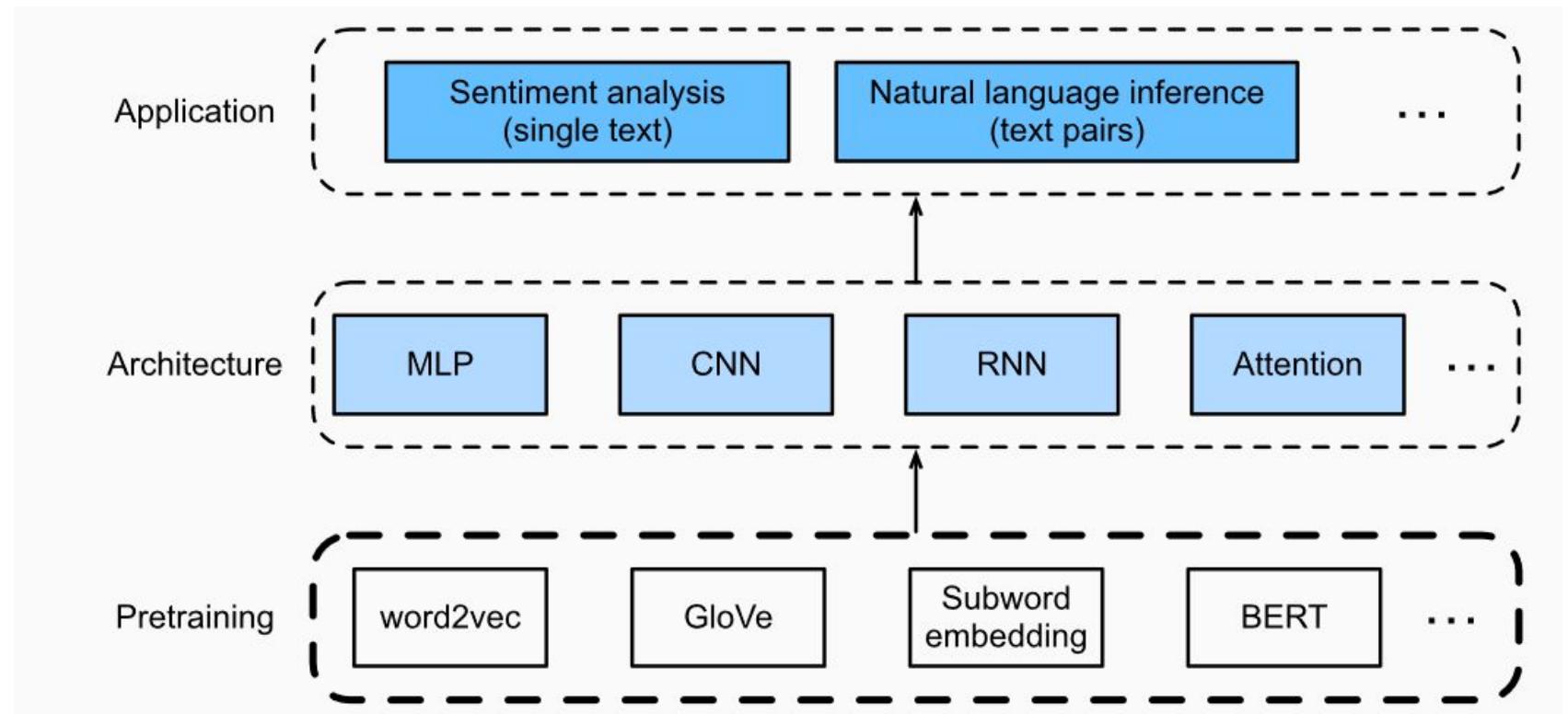
# Pipeline Architecture



# Deep Learning Architecture ‘End-to-End’



# Deep Learning - Pre-Trained Models



# Overview

Organizational issues

Why NLP

Human language

Some common notions in NLP

Architectures

**Applications**

Useful pointers



# NLP Applications

**Classification** such as Sentiment Analysis

**Retrieval** such as Question Answering

Natural Language **Generation** such as Dialog Systems (Chatbots)

**Transformation** such as Machine Translation

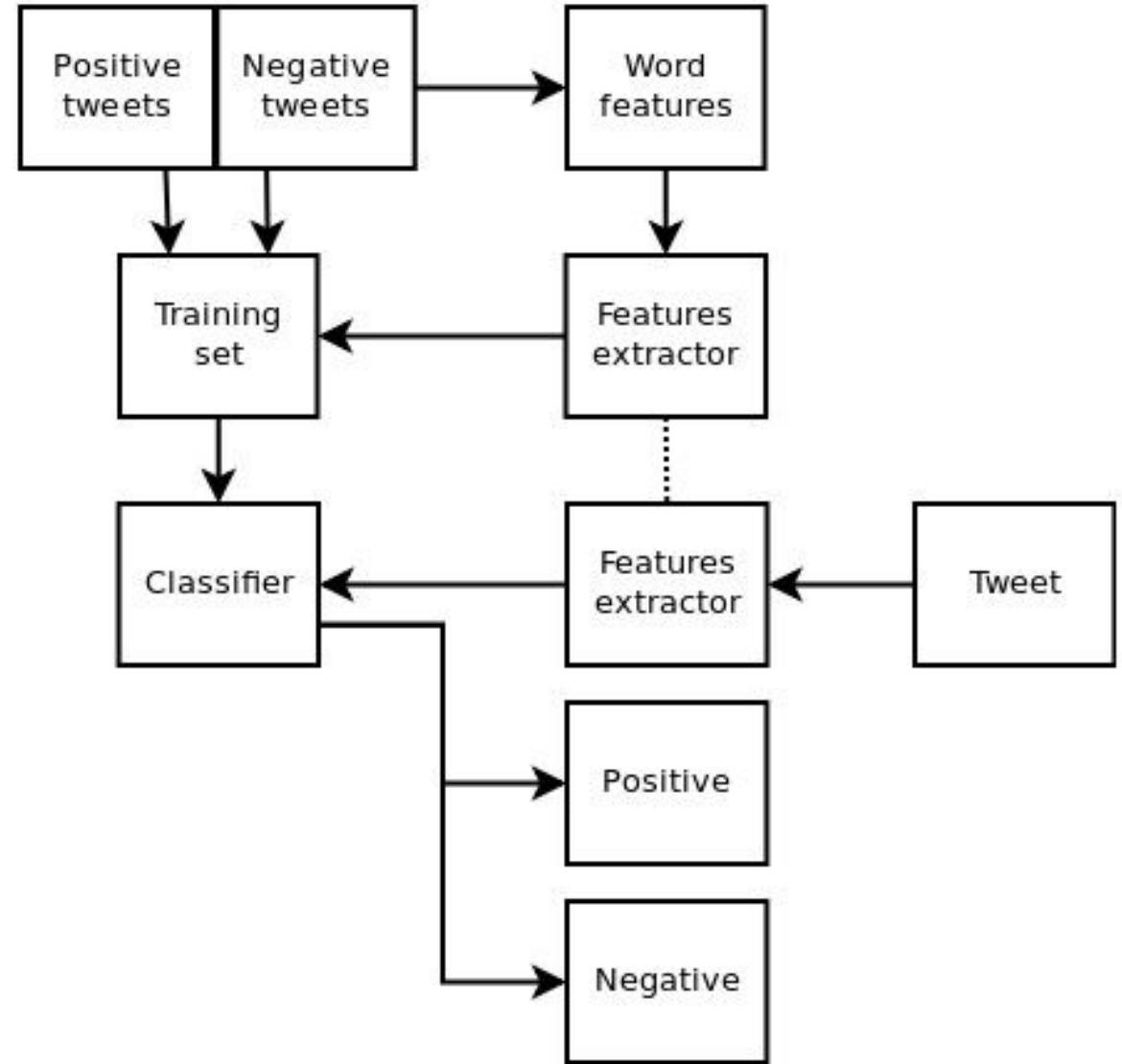


# NLP Applications - Classification

Many NLP tasks can be formulated as a **classification** problem, e.g.

- |                            |   |
|----------------------------|---|
| Sentiment Analysis:        | classification of a sentence as Pos or Neg        |
| Part of Speech Tagging:    | classification of a word as a Noun, Verb, Adj,... |
| Word Sense Disambiguation: | classification of ‘bank’ as Sense-1 or Sense-2    |

# Sentiment Analysis



# NLP Applications - Retrieval

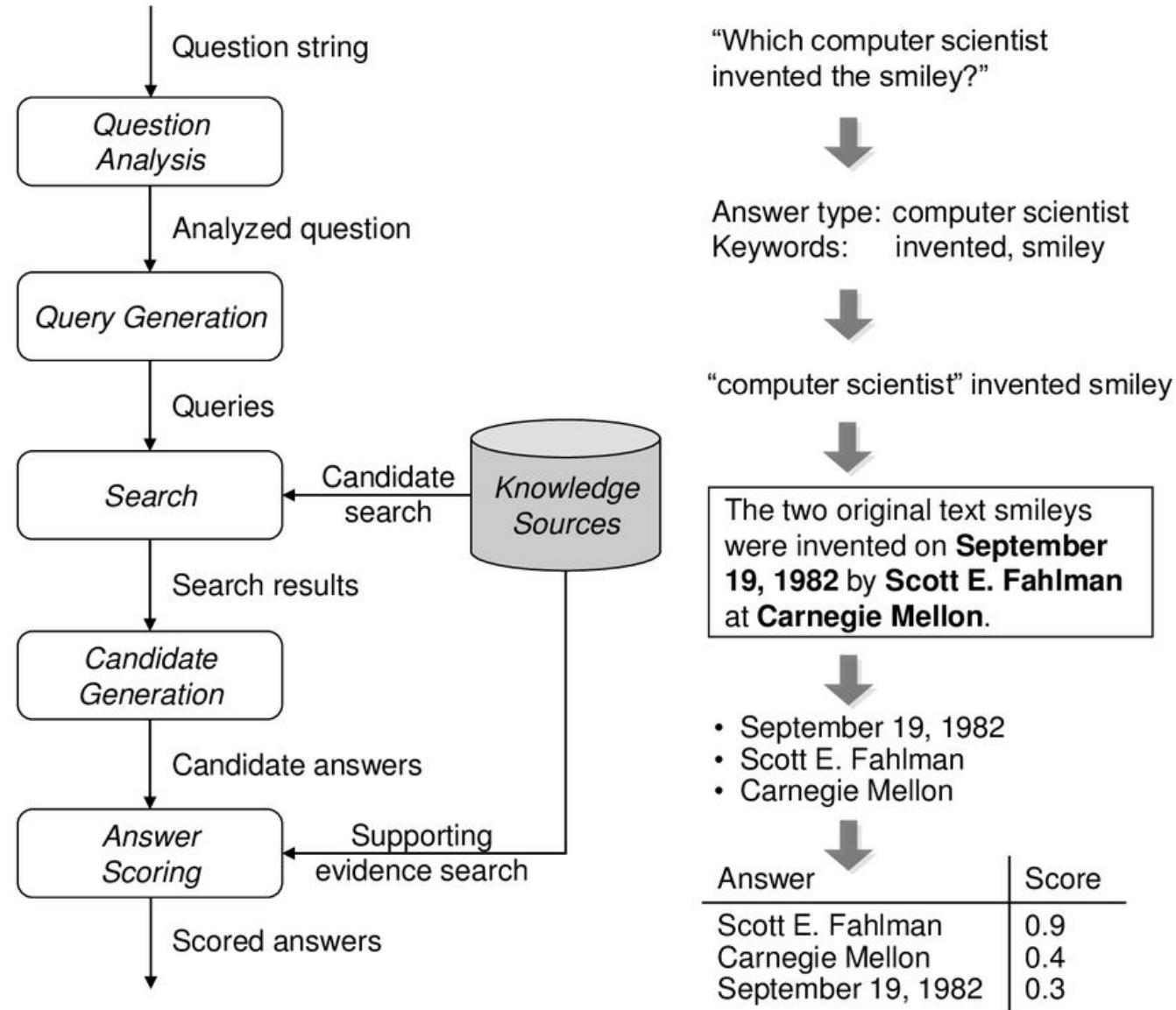
NLP tasks based on **retrieval**, e.g.

Question Answering



NUI Galway  
OÉ Gaillimh

# Question Answering



# NLP Applications - Generation

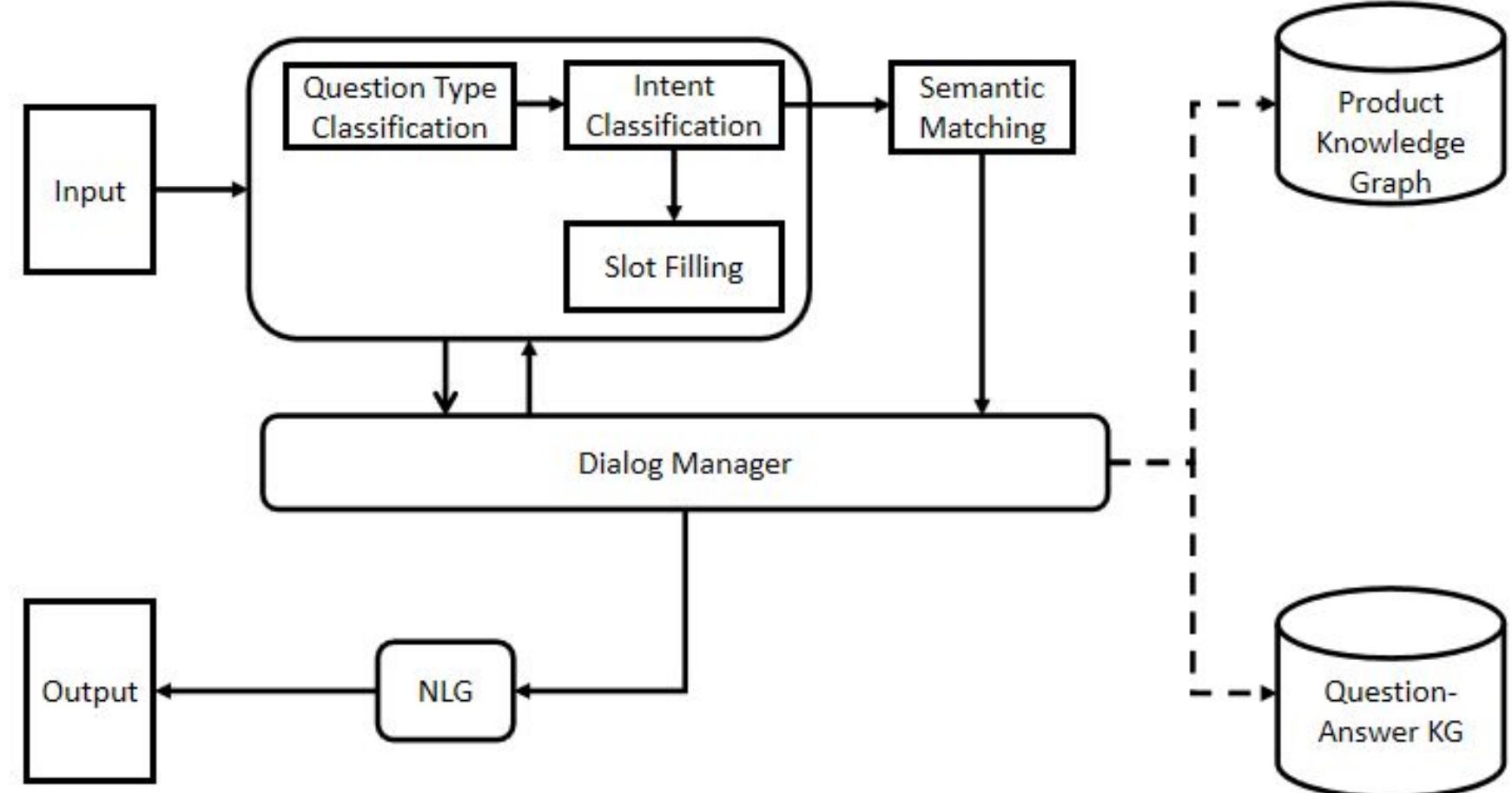
NLP tasks based on **generation**, e.g.

Dialog Systems (Chatbots)

Natural Language Generation from Data

Abstractive Document Summarization

# Dialog Systems (Chatbots)

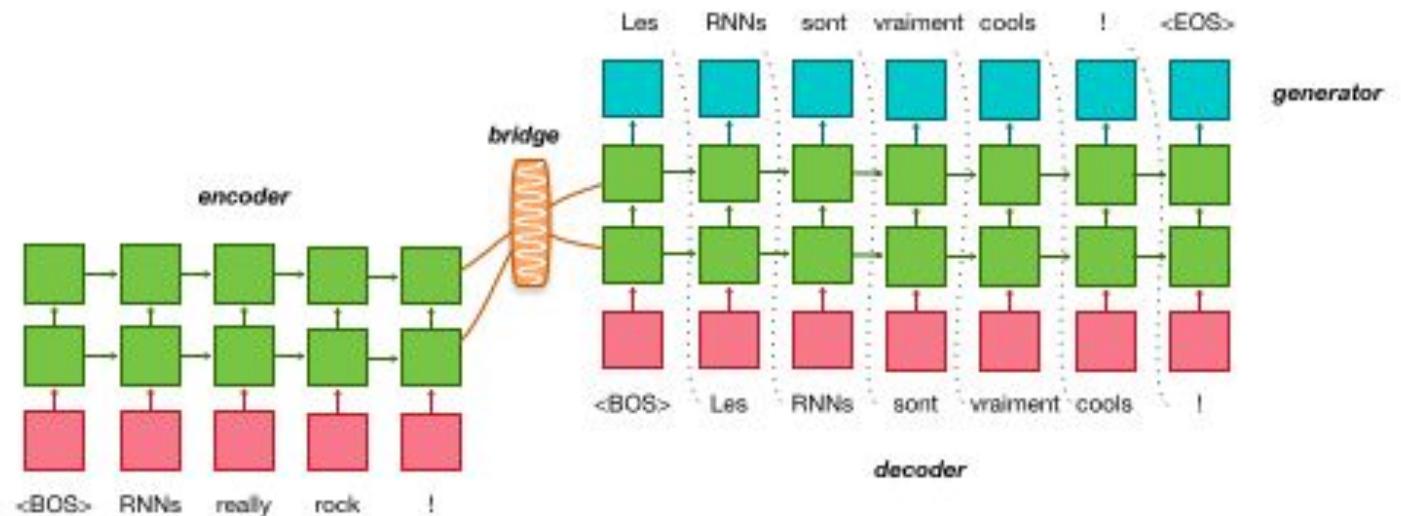


# NLP Applications - Transformation

NLP tasks based on **transformation**, e.g.

Machine Translation

# Neural Machine Translation



# Bilingual Corpus

What  
is  
the  
anticipated  
cost  
of  
collecting  
fees  
under  
the  
new  
proposal  
?  
En  
vertu  
de  
les  
nouvelles  
propositions  
,

quel  
est  
le  
coût  
prévu  
de  
perception  
de  
les  
droits  
?



# Overview

Organizational issues

Why NLP

Human language

Some core notions in NLP

Architectures

Applications

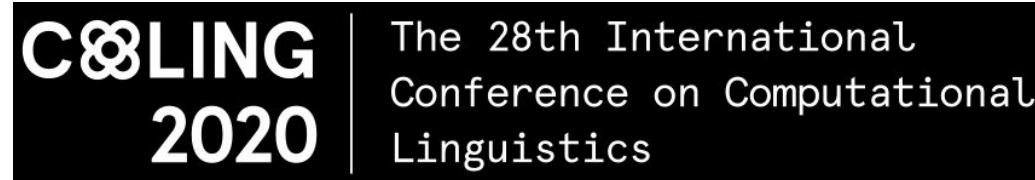
**Useful pointers**



# NLP Publications - Conferences



Association for  
Computational Linguistics



LREC 2020  
Marseille  
Palais du Pharo  
May 11-16, 2020

# NLP Publications - Journals

Transactions of the Association for Computational Linguistics

HOME ABOUT LOGIN REGISTER SEARCH CURRENT VOLUME ALL VOLUMES ANNOUNCEMENTS

Home > Volumes > Vol 3 (2015)

Vol 3 (2015)

Table of Contents

Reasoning about Quantifiers in Natural Language  
Subroto Roy, Tim Vieira, Dan Roth

Cross-Document Co-Reference Resolution using Sample-Based Clustering with Knowledge Enrichment  
Sourov Dutta, Gerhard Weikum

Efficient Inference and Structured Learning for Semantic Role Labeling  
Oscar Tackström, Kuzman Ganchev, Dipanjan Das

SPIRTE: Generalizing Topic Models with Structured Priors  
Michael J. Paul, Marc Dredze

A Sense-Topic Model for Word Sense Induction with Unsupervised Data Enrichment  
Jing Wang, Mohit Bansal, Kevin Gimpel, Brian D. Ziebart, Clement T. Yu

Which Step Do I Take First? Troubleshooting with Bayesian Models  
Armine Louis, Mirella Lapata

Gappy Pattern Matching on GPUs for On-Demand Extraction of Hierarchical Translation Grammars  
Hu He, Jimmy Lin, Adam Lopez

Erratum: "Exploring Compositional Architectures and Word Vector Representations for Prepositional Phrase Attachment"  
Yonatan Belinkov, Tao Lei, Regina Barzilay, Amir Globerson

A Bayesian Model of Grounded Color Semantics  
Brian McMahan, Matthew Stone

Exploiting Parallel News Streams for Unsupervised Event Extraction  
Congle Zhang, Stephen Soderland, Daniel S. Weld

Unsupervised Declarative Knowledge Induction for Constraint-Based Learning of Information Structure in Scientific Documents  
Yufan Guo, Roi Reichart, Anna Korhonen

Entity Disambiguation with Web Links  
Andrew Chisholm, Ben Hachey

An Unsupervised Method for Uncovering Morphological Chains  
Karthik Narasimhan, Regina Barzilay, Tommi Jaakkola

PDF (PRESENTED AT NAACL 2015) 1-13

PDF 15-28

PDF (PRESENTED AT ACL 2015) 29-41

PDF (PRESENTED AT NAACL 2015) 43-57

PDF (PRESENTED AT NAACL 2015) 59-71

PDF 73-85

PDF (PRESENTED AT NAACL 2015) 87-100

ERRATUM PDF ORIGINAL PAPER 101-101

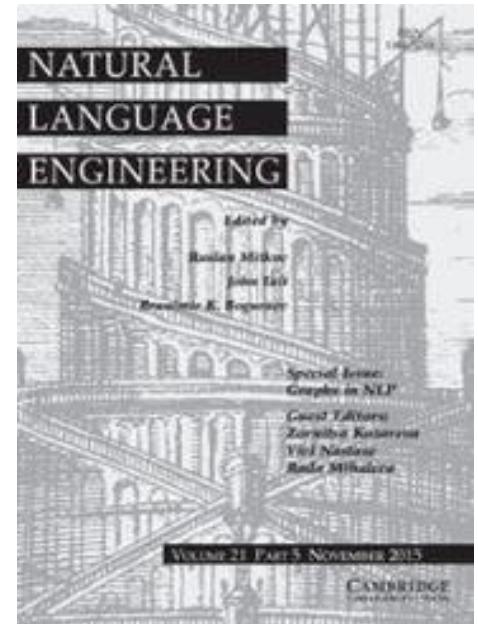
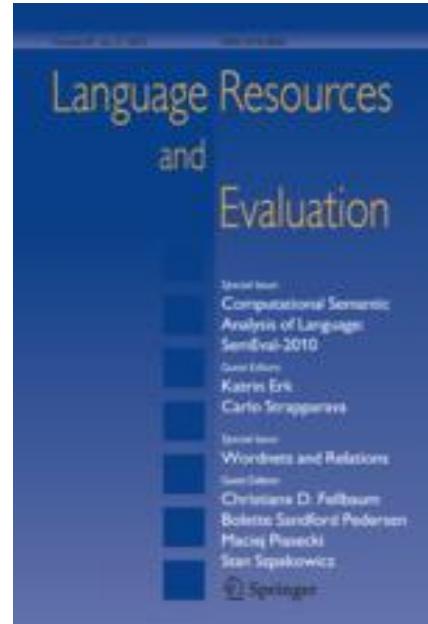
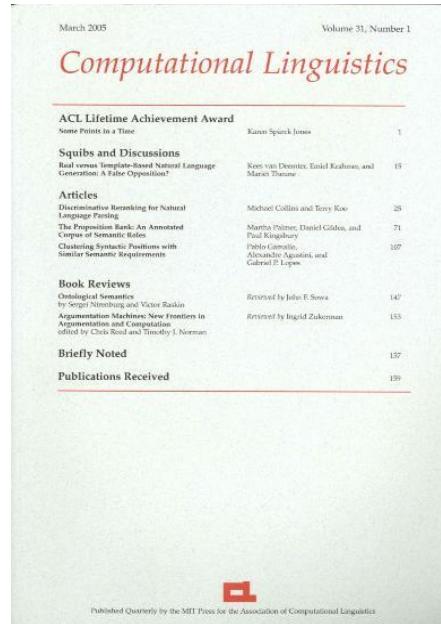
PDF (PRESENTED AT NAACL 2015) 103-115

PDF (PRESENTED AT ACL 2015) 117-129

PDF (PRESENTED AT NAACL 2015) 131-143

PDF (PRESENTED AT NAACL 2015) 145-156

PDF (PRESENTED AT ACL 2015) 157-167



# NLP Publications - Ranking

Categories > Engineering & Computer Science > Computational Linguistics ▾

Publication	<a href="#">h5-index</a>	<a href="#">h5-median</a>
1. Meeting of the Association for Computational Linguistics (ACL)	<a href="#">106</a>	168
2. Conference on Empirical Methods in Natural Language Processing (EMNLP)	<a href="#">88</a>	157
3. Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (HLT-NAACL)	<a href="#">61</a>	104
4. International Workshop on Semantic Evaluation	<a href="#">49</a>	89
5. Transactions of the Association for Computational Linguistics	<a href="#">47</a>	80
6. International Conference on Language Resources and Evaluation (LREC)	<a href="#">45</a>	71
7. International Conference on Computational Linguistics (COLING)	<a href="#">41</a>	68
8. Conference of the European Chapter of the Association for Computational Linguistics (EACL)	<a href="#">36</a>	62
9. Computer Speech & Language	<a href="#">36</a>	48
10. Conference on Computational Natural Language Learning (CoNLL)	<a href="#">34</a>	51
11. Workshop on Machine Translation	<a href="#">27</a>	45
12. IEEE Spoken Language Technology Workshop (SLT)	<a href="#">26</a>	43
13. Annual Meeting of the Special Interest Group on Discourse and Dialogue (SIGDIAL)	<a href="#">26</a>	39
14. Computational Linguistics	<a href="#">25</a>	51
15. Language Resources and Evaluation	<a href="#">23</a>	34

# NLP Forums - Mailing Lists

## Corpora List

<http://clu.uni.no/icame/corpora/>

## Linguist List

<http://www.linguistlist.org/>

# ACL - Professional Association for NLP

**Association for Computational Linguistics (ACL)**

<http://www.aclweb.org/>

**ACL Anthology:** archive of NLP conference publications

<https://www.aclweb.org/anthology/>

# Lab of this Week

Intro to NLTK - labs will use NLTK in Python

<https://www.nltk.org/>



NUI Galway  
OÉ Gaillimh

# Advanced Topics in NLP - CT5121

Optional Semester 2 course by Dr. John McCrae & Dr. Mihael Arčan

Word Embeddings

Recurrent Neural Network (RNN) & Long Short-Term Memory (LSTM)

Machine Translation

Textual Similarity

Topic Models

Linguistic Linked Open Data

Lexicography and Under-resourced Languages



NUI Galway  
OÉ Gaillimh

QA

