

Natural Language Processing

Lecture 01

Dirk Hovy

dirk.hovy@unibocconi.it

 @dirk_hovy

Text is an exploding data source

Exabytes = 1M TB

- You read ~9000 words per day
- = 200.000.000 words in a lifetime
- = 0.4 GB of data
- 44 billion GB of new data each day

60-80% GROWTH/YEAR

UNSTRUCTURED DATA

STRUCTURED DATA

Source: IDC

NLP is booming



\$136.000.000

\$5.400.000.000

2016

2017

2018

2019

2020

2021

2022

2023

2024

2025

Machine Translation



Text Generation



In a shocking finding, scientist discovered a herd of unicorns living in a remote, previously unexplored valley, in the Andes Mountains. Even more surprising to the researchers was the fact that the unicorns spoke perfect English.

The scientist named the population, after their distinctive horn, Ovid's Unicorn. These four-horned, silver-white unicorns were previously unknown to science.

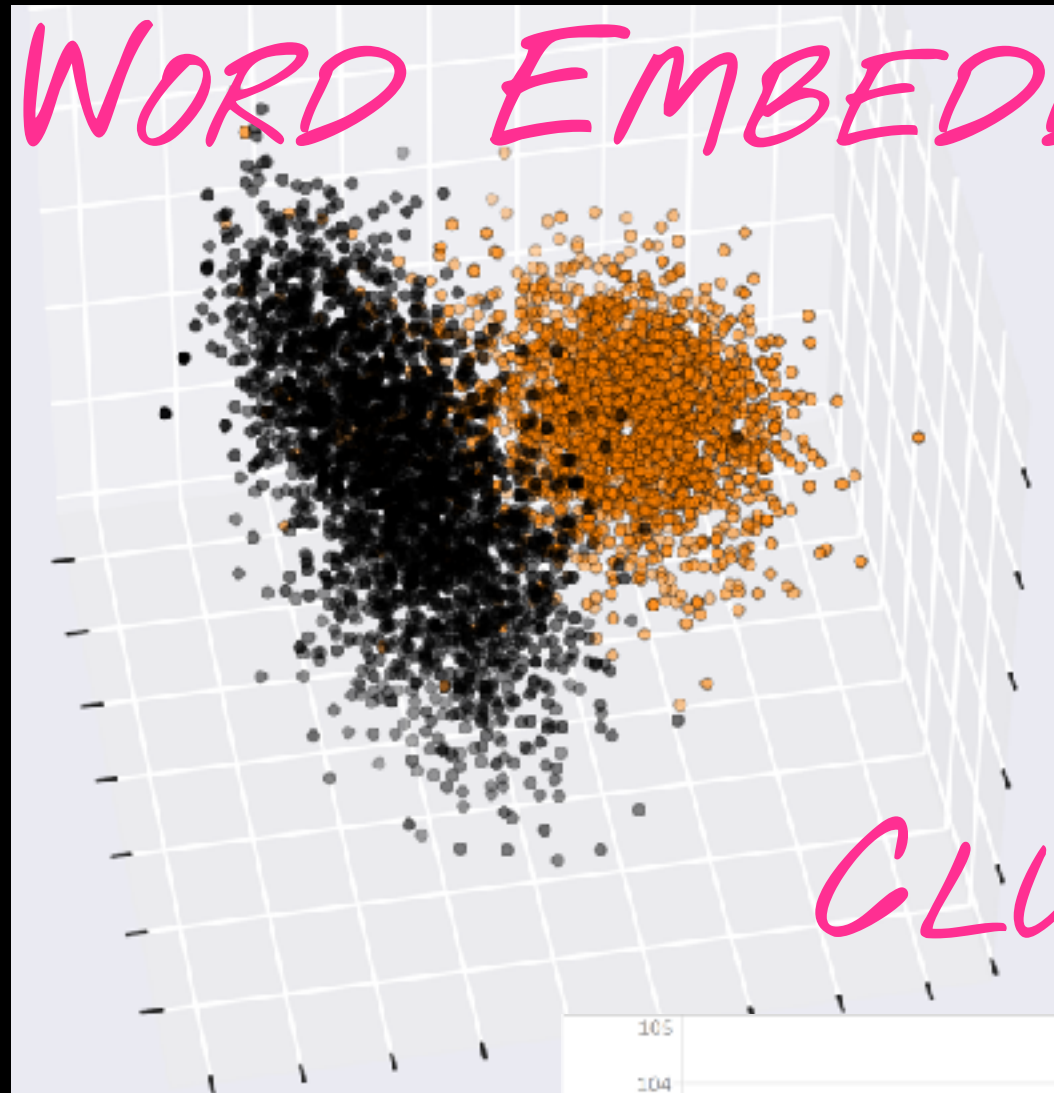
Now, after almost two centuries, the mystery of what sparked this odd phenomenon is finally solved.

Dr. Jorge Pérez, an evolutionary biologist from the University of La Paz, and several companions, were exploring the Andes Mountains when they found a small valley, with no other animals or humans. Pérez noticed that the valley had what appeared to be a natural fountain, surrounded by two peaks of rock and silver snow.

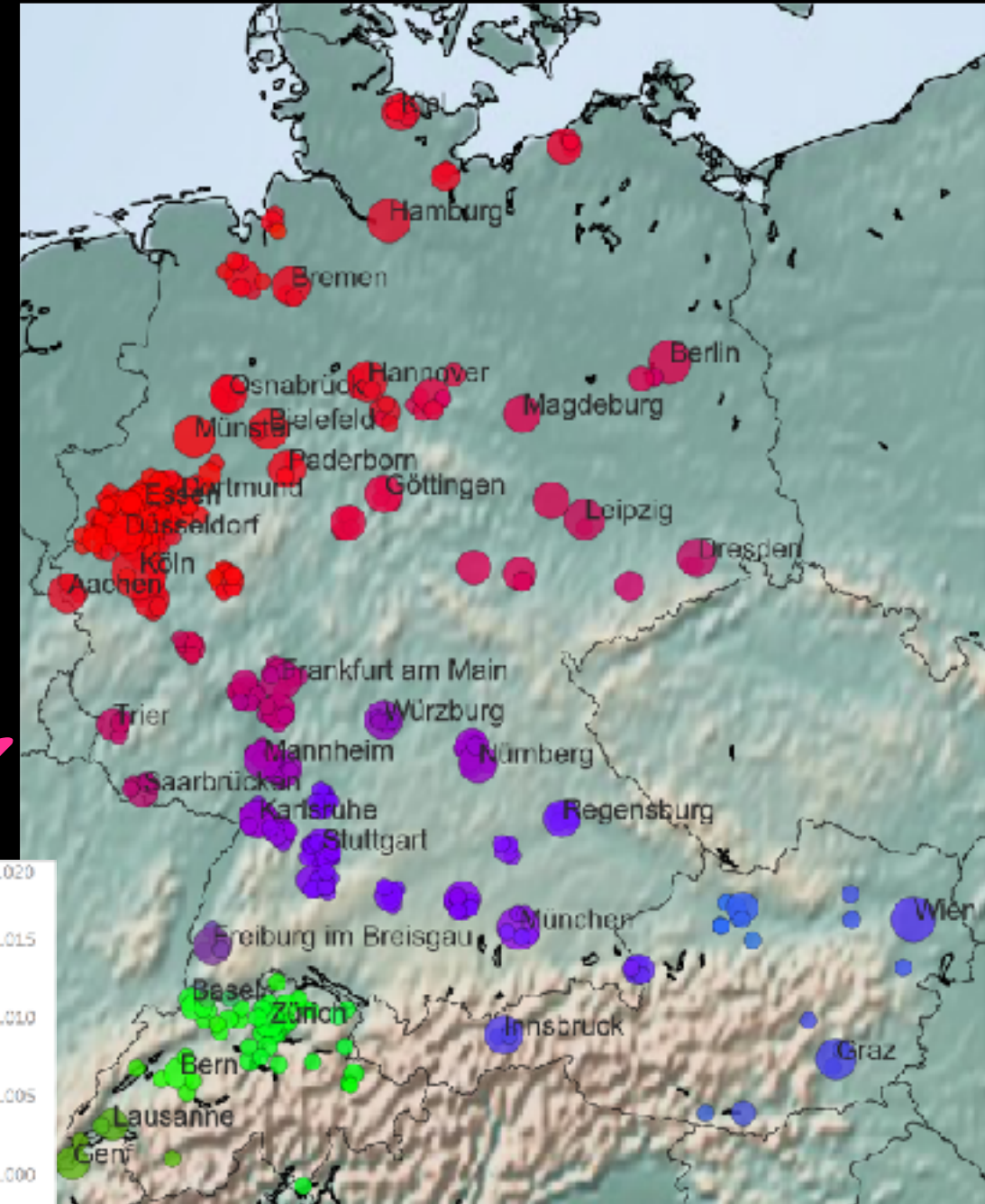
Pérez and the others then ventured further into the valley. "By the time we reached the top of one peak, the water looked blue, with some crystals on top," said Pérez.

The Goals

WORD EMBEDDINGS



CLUSTERING

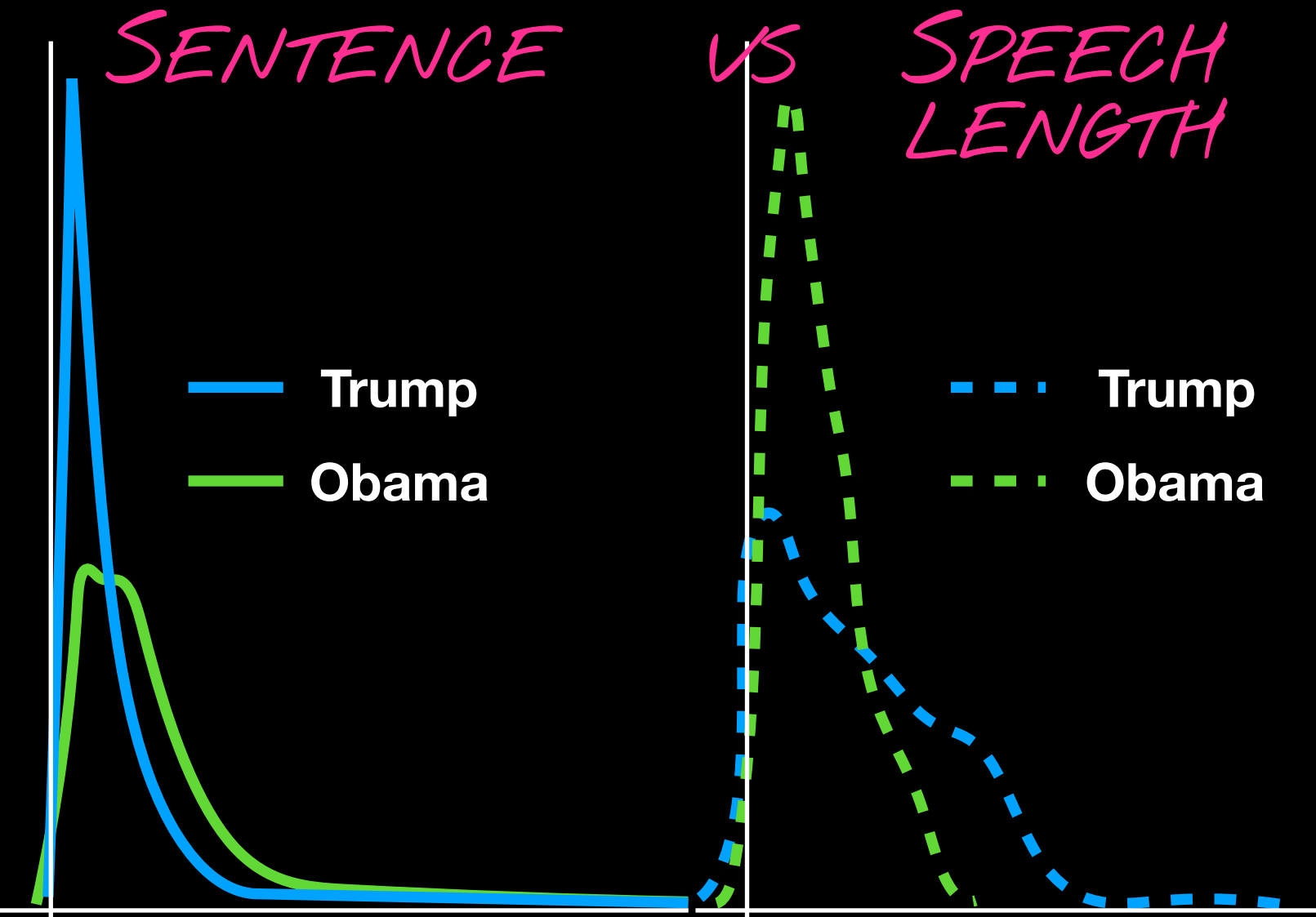


VISUALIZATION

SENTIMENT ANALYSIS



Student Projects



TALK SIMILARITY

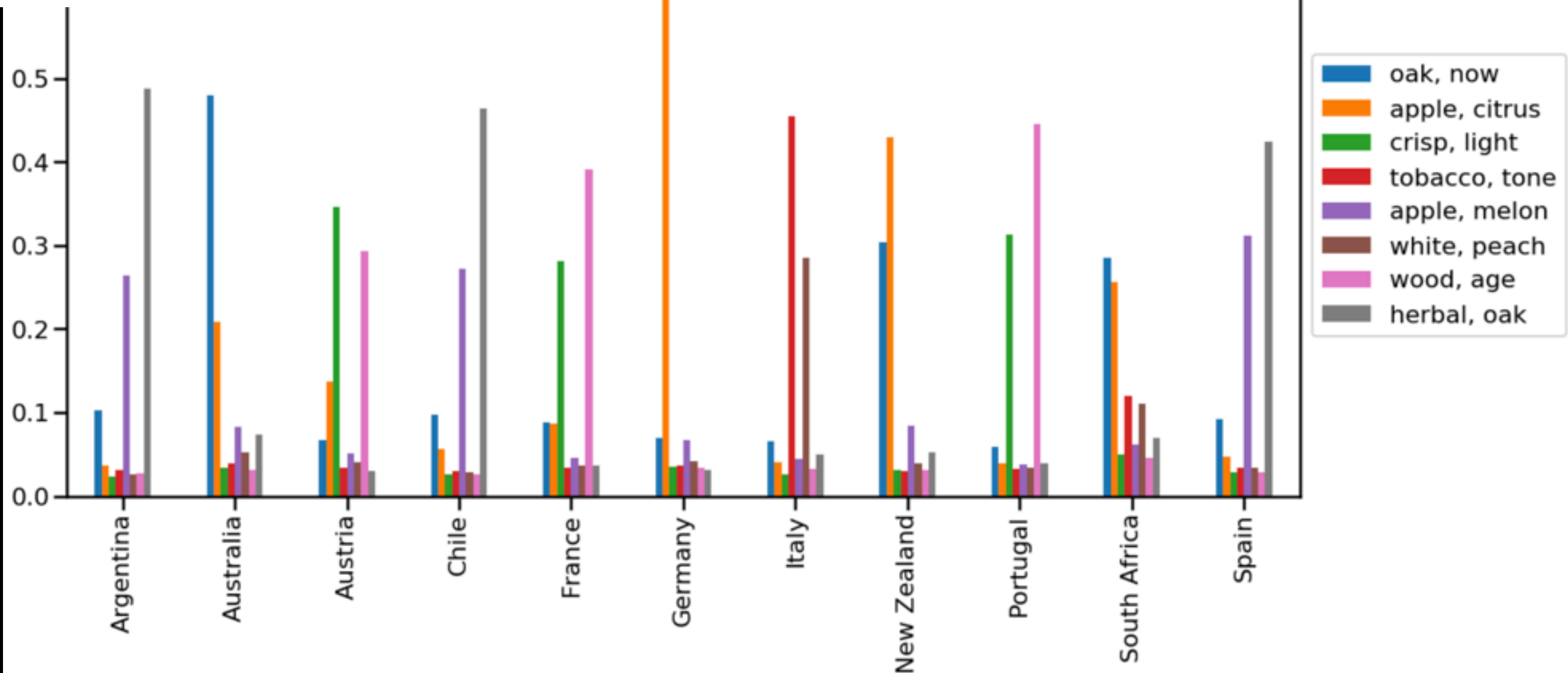
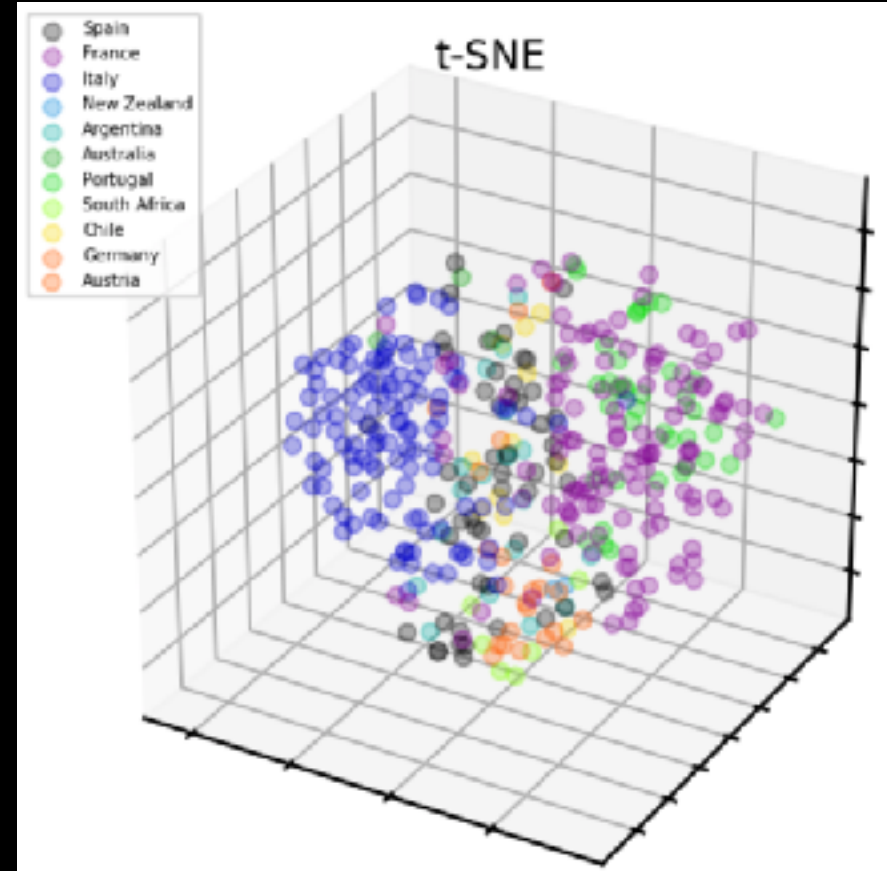
*CANDIDATE LMS FINISH
A SENTENCE STARTING
WITH "AMERICA"*

"America first — America first ..."

"America was actually on track to top
\$ 1 trillion in spending over the coming
decade -- because the freedom and
dignity --"

Student Projects

WINES FROM
DIFFERENT COUNTRIES



Syllabus

Lesson	Topic
1	Intro to NLP 1
2	Linguistic Analysis
3	Information Retrieval
4	Regular Expressions
5	Discrete Representations and TFIDF
6	Discrete Representations
7	Continuous Representations
8	Word2Vec and Doc2Vec
9	Topic models 1
10	Topic models 2
11	Dimensionality reduction and Clustering, Visualization
12	Latent Dimensions
13	Text classification
14	Example: Sentiment Classification
15	The Perceptron
16	Multilayer Perceptrons
17	Structured Prediction
18	The Structured Perceptron
19	Convolutional and Recurrent Neural Networks
20	LSTMs in keras
21	CNN in keras
22	The Transformer and BERT
23	Ethics in NLP
24	Final Project Presentations

BASICS

EXPLORATION

PREDICTION

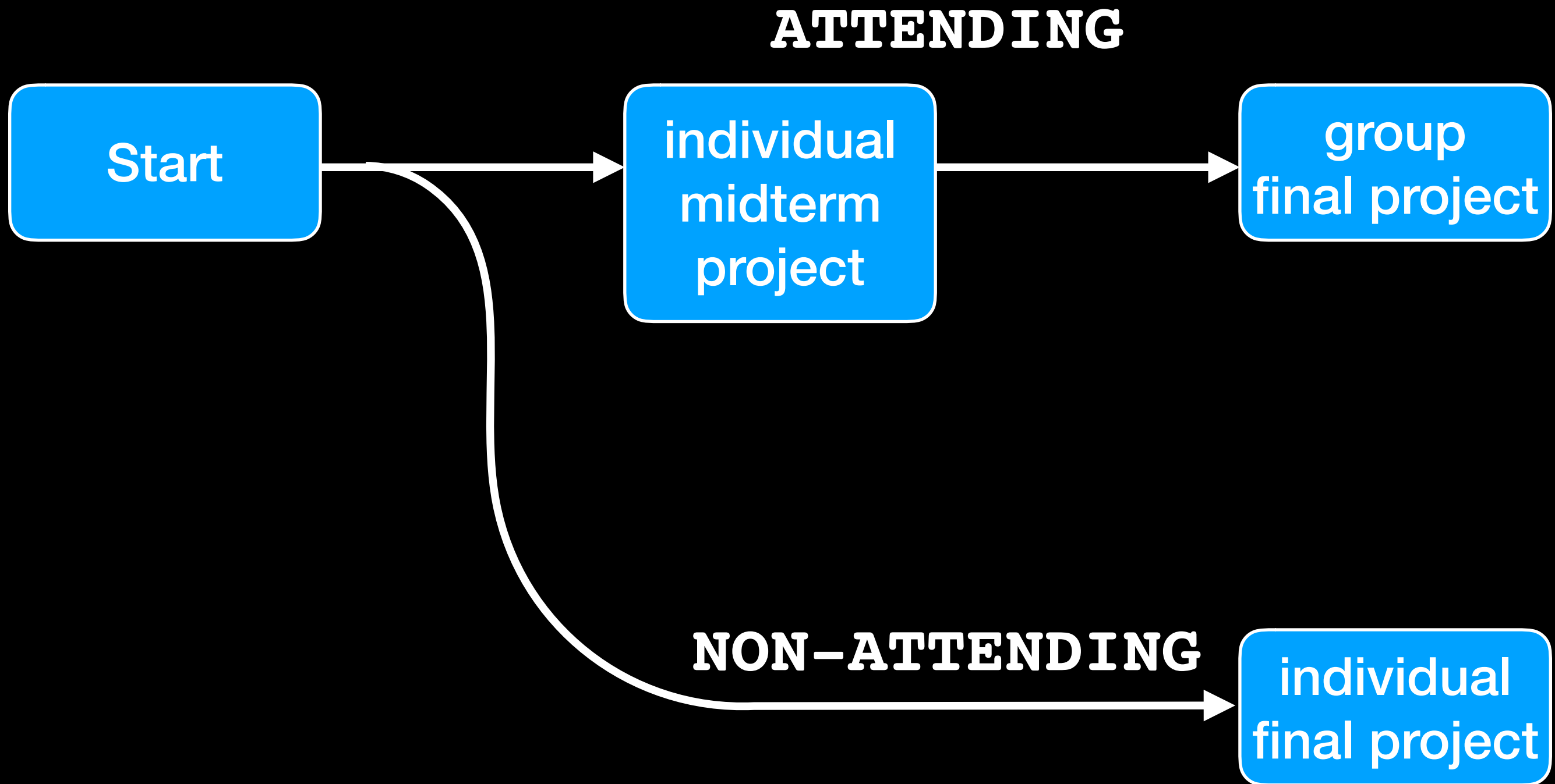
New and Improved!

- Lecture notes now available as book draft for download from BBoard (first part as book at Cambridge University Press)
- Additional focus on latest neural network methods

Class Structure

- Thursdays: intuition, theory, math (slides)
- Fridays: exercises and practice (Jupyter Notebooks)
- Material on BBoard or at <https://github.com/dirkhovy/NLPclass>
- Latest book draft: <http://www.dirkhovy.com/portfolio/papers/download/nlpss.pdf>

Attendance



Grading

Individual midterm project (50%): Exploration and visualization

Group final project (50%): Data annotation and prediction, visualization

Individual final project (100%): whole class

- All projects are to be handed in as runnable **Jupyter Notebooks**
- Graded on data set size, correctness of implementations, annotation quality, performance of prediction
- No point changes, only **complete regrades** (total can go down)!

How do I succeed?

- Code well
- Pay attention
- Code some more

I want more NLP!

- join our reading group (Thu at 13:00), email Tommaso Fornaciari, fornaciari@unibocconi.it
- DMI online talk series:
 - February 17th, 12:30 CET, Ciro Cattuto (U. Torino)
 - March 1st, 12:30 CET/11:30 UK, Sebastian Ruder (DeepMind)
 - March 15th, 12:30 CET, Fabiana Zollo (U. Ca' Foscari Venice)
 - March 29th, 12:30 CET, Florian Ellsaesser (Frankfurt School of Business)
 - April 12th, 12:30 CET, Dong Nguyen (U. Utrecht)
 - April 26th, 16:30 CET, Nikita Nangia (NYU)
 - May 10th or 24th, 16:30 CET, Adina Williams (Facebook)

What to do with a problem

1. Don't panic.
2. Google it. stackoverflow.com is your friend
3. Talk to your classmates
4. Ask the TA, Tommaso Fornaciari
`fornaciari@unibocconi.it`
5. Make a you@B appointment for my office hours (Mon, 18-19:30)

WARNING :

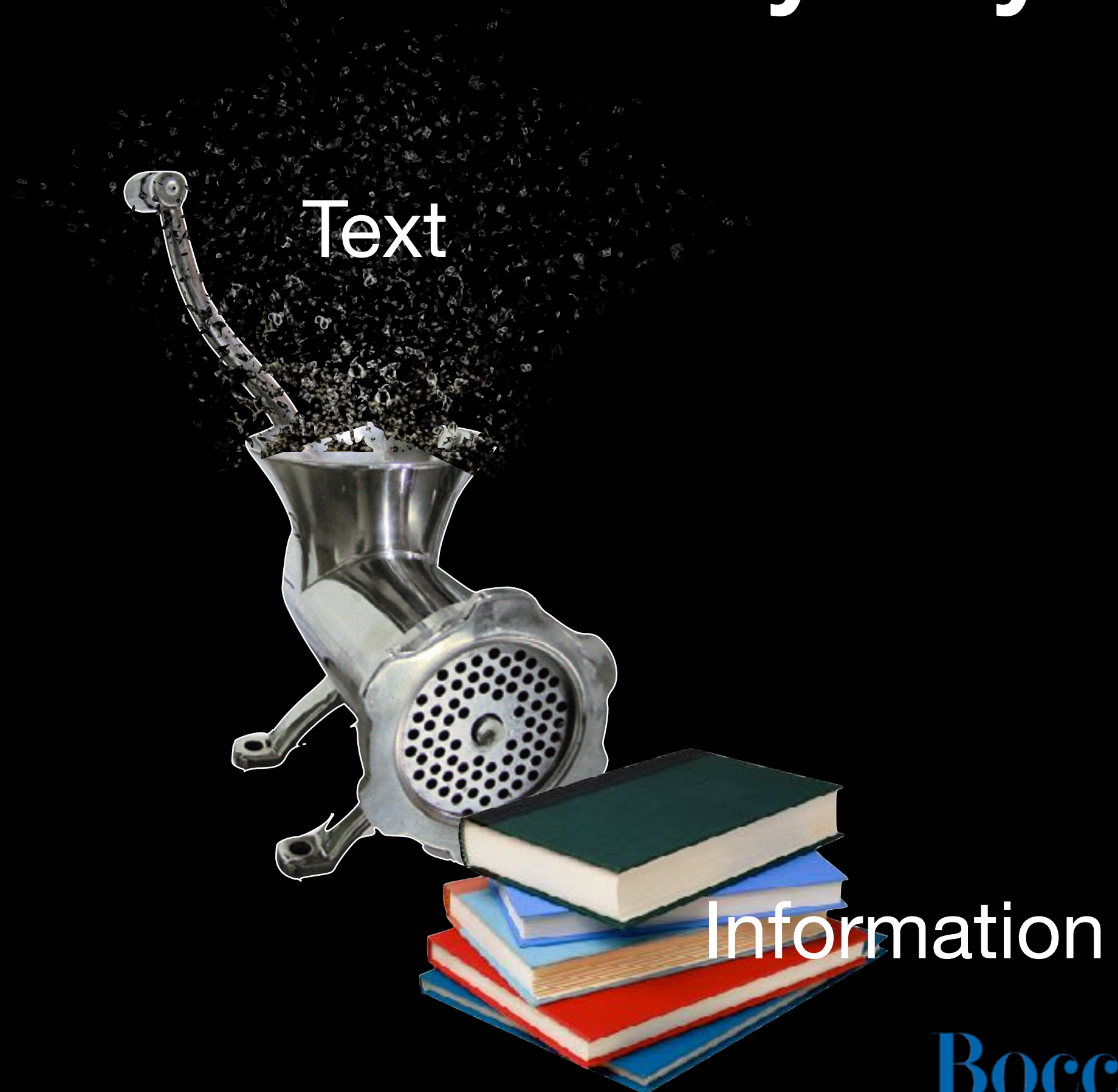
For any question we can solve with a Google search, we deduct points!

Let's start!

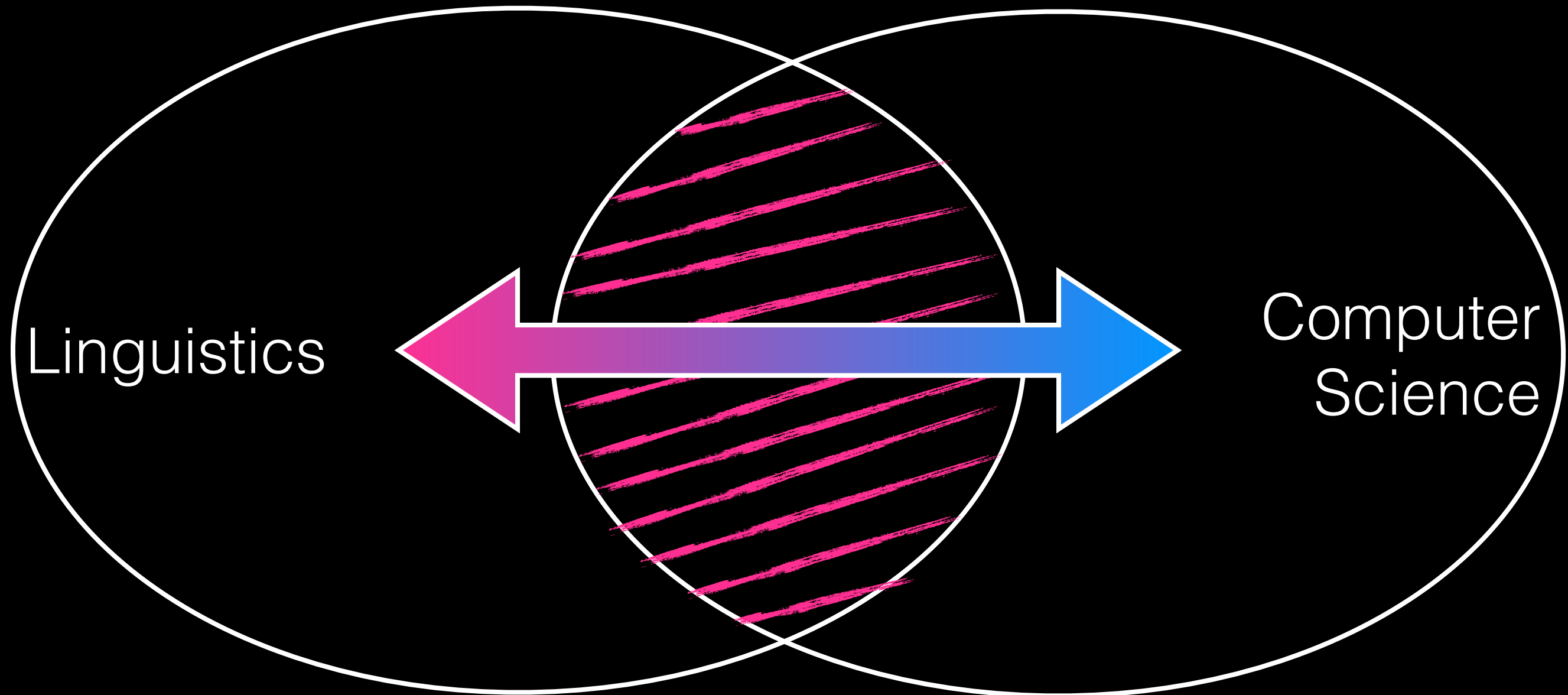
Today's Goals

- Understand where NLP comes from
- Learn about the different steps of preprocessing
- Understand the use of
 - parts of speech,
 - parsing, and
 - named entities

So, what's NLP anyway?

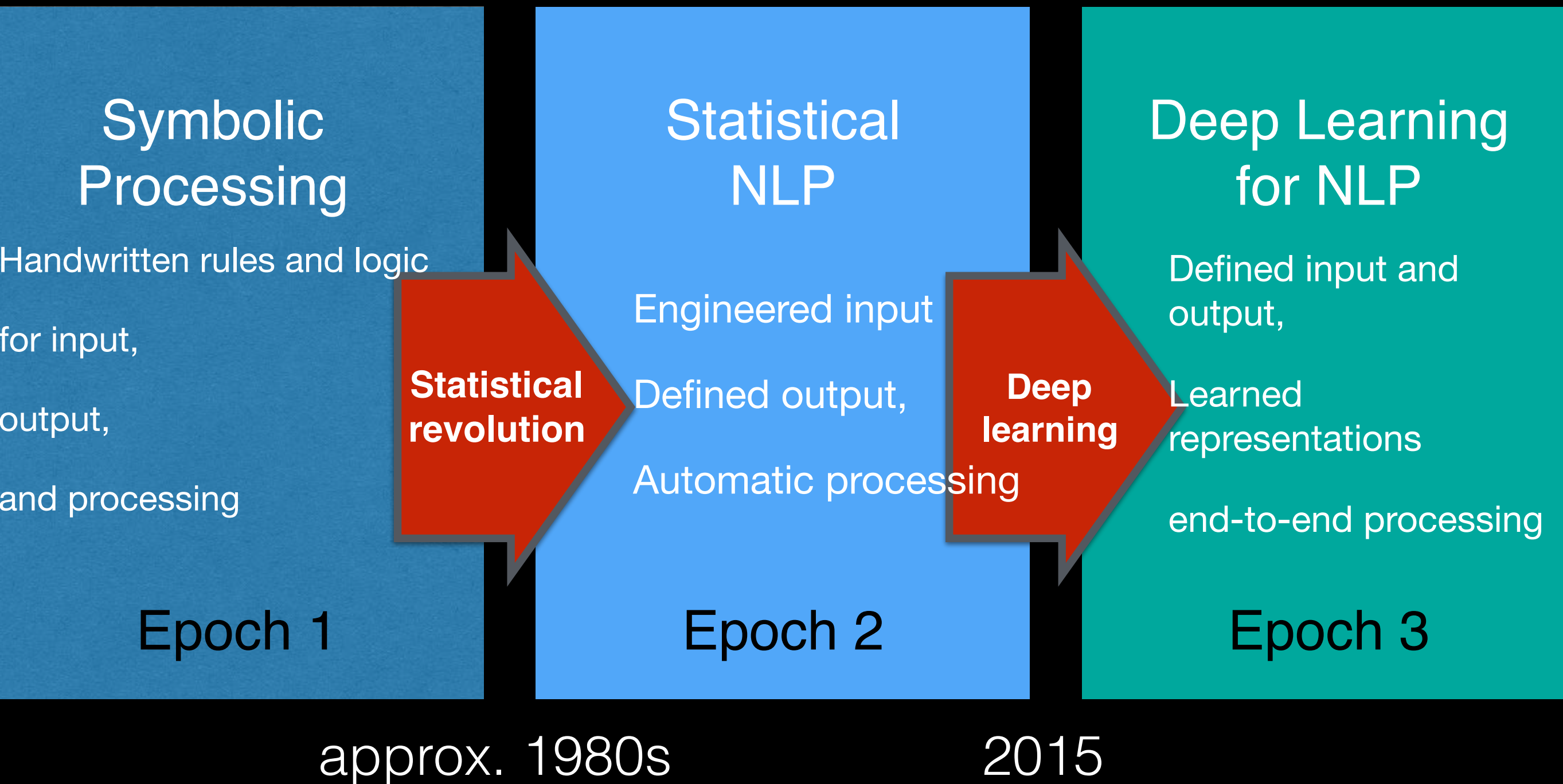


The two sides of NLP



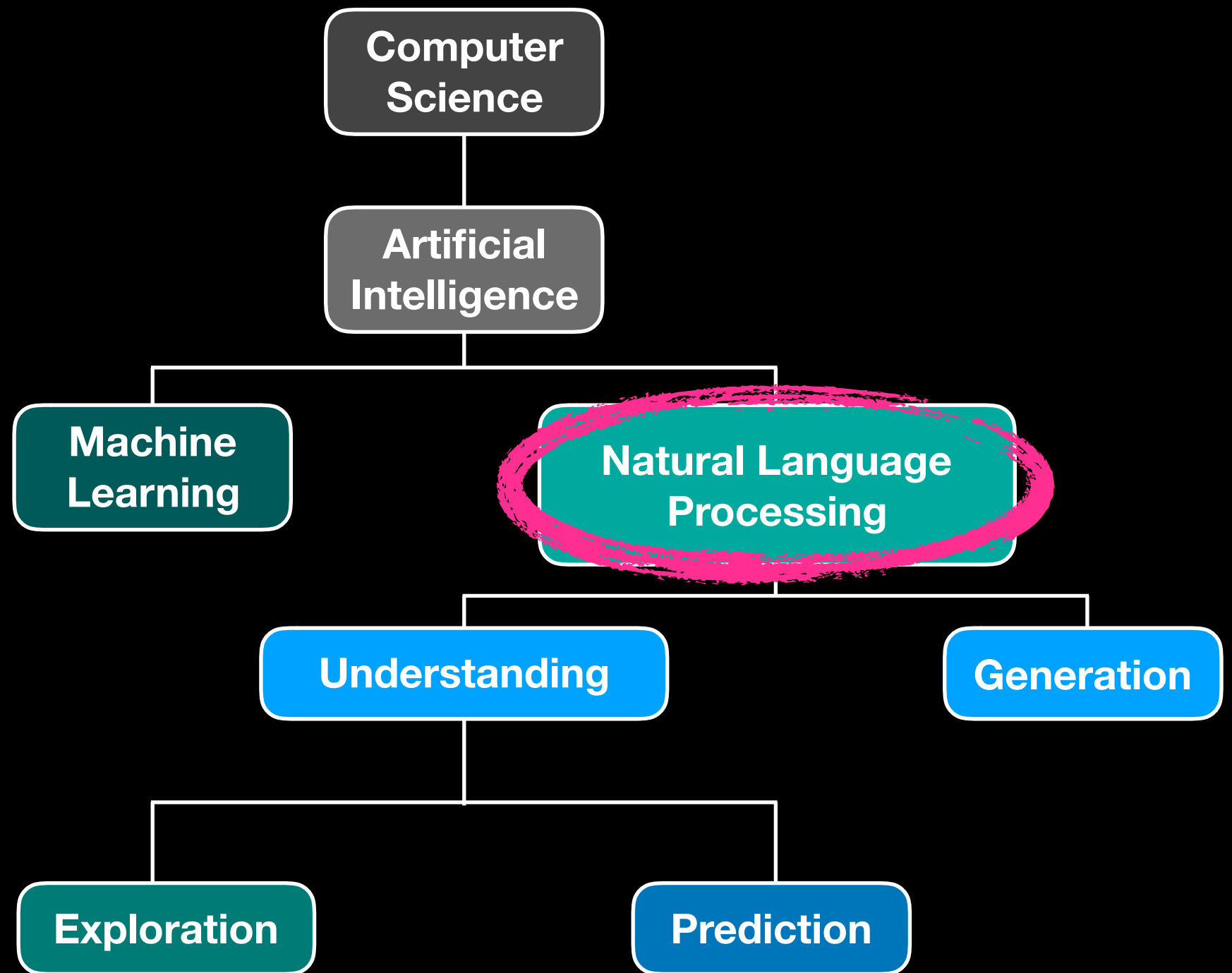
informed linguistic hypotheses large-scale statistical analysis

A very Brief History of NLP

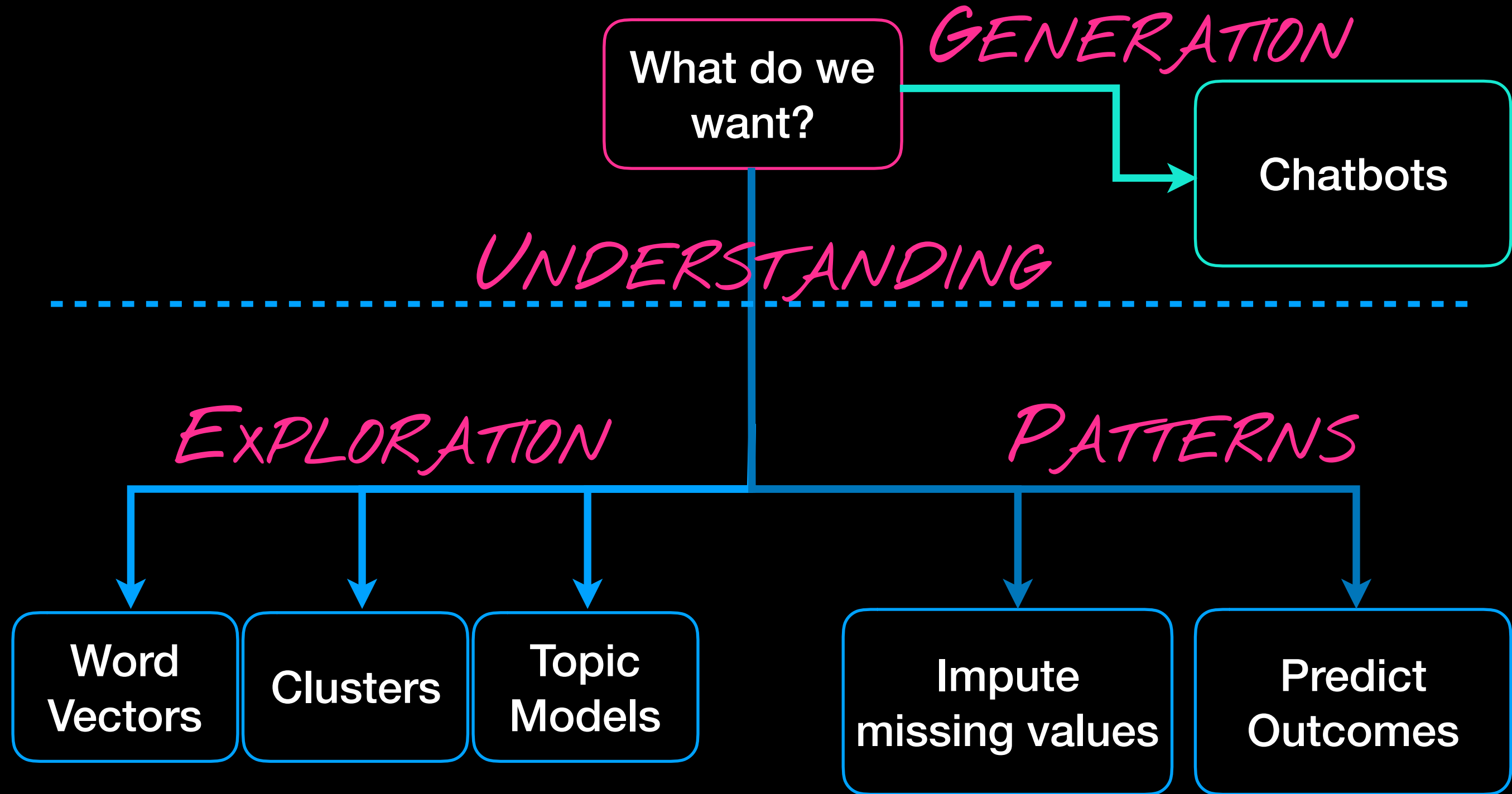


The NLP Family

- Extract information from text: topics, trends
- Classify text sentiment, content type, author profiles
- Generate text: translations, automated responses



Two Uses of NLP



Linguistic Analysis

Examples of Analysis

NER



PERSON

PERSON

PARSING

nsubj

~~dobj~~

punct

nn

POS

PRON

VERB

PROP N

PROP N

PUNCT

1

1

1

1

1

1

admire

Rosa Parks

☐

Pre-processing



Pre-processing steps

```
<div id="text">I've been in New York  
in 2011, but didn't like it. I  
preferred Los Angeles.</div>
```

GOAL: MINIMIZE VARIATION



Pre-processing steps

- Remove formatting (e.g. HTML)
- Segment sentences
- Tokenize words
- Normalize words
 - numbers
 - lemmas vs. stems
- Remove unwanted words
 - stopwords
 - content words (use POS tagging!)
- join collocations

I've been in New York in
2011, but didn't like
it. I preferred Los
Angeles.



Pre-processing steps

- Remove formatting (e.g. HTML)

- Segment sentences

- Tokenize words

- Normalize words

- numbers

- lemmas vs. stems

- Remove unwanted words

- stopwords

- content words (use POS tagging!)

- join collocations

I've been in New York in
2011, but didn't like
it.

I preferred Los Angeles.



Pre-processing steps

- Remove formatting (e.g. HTML)

- Segment sentences

- Tokenize words

- Normalize words

- numbers

- lemmas vs. stems

- Remove unwanted words

- stopwords

- content words (use POS tagging!)

- join collocations

I 've been in New York
in 2011 , but did n't
like it .

I preferred Los
Angeles .



Pre-processing steps

- Remove formatting (e.g. HTML)

- Segment sentences

- Tokenize words

- Normalize words

- numbers

- lemmas vs. stems

- Remove unwanted words

- stopwords

- content words (use POS tagging!)

- join collocations

i 've been in new york
in 0000 , but did n't
like it .

i preferred los
angeles .



Pre-processing steps

- Remove formatting (e.g. HTML)

- Segment sentences

- Tokenize words

- Normalize words

- numbers

- lemmas vs. stems

- Remove unwanted words

- stopwords

- content words (use POS tagging!)

- join collocations

i have be in new york in
0000 , but do not like
it .

i prefer los angeles .



Pre-processing steps

- Remove formatting (e.g. HTML)
- Segment sentences
- Tokenize words
- Normalize words
 - numbers
 - lemmas vs. stems
- Remove unwanted words
 - stopwords
 - content words (use POS tagging!)
- join collocations

i new york 0000 , like .

i prefer los angeles .



Pre-processing steps

- Remove formatting (e.g. HTML)

new york 0000 like

- Segment sentences

- Tokenize words

prefer los angeles

- Normalize words

- numbers

- lemmas vs. stems

CONTENT = (NOUN, VERB, NUM)

- Remove unwanted words

- stopwords

- content words (use POS tagging!)

- join collocations



Pre-processing steps

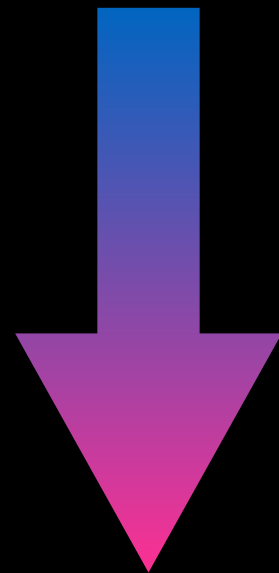
- Remove formatting (e.g. HTML)
- Segment sentences
- Tokenize words
- Normalize words
 - numbers
 - lemmas vs. stems
- Remove unwanted words
 - stopwords
 - content words (use POS tagging!)
- join collocations

`new_york 0000 like`

`prefer los_angeles`

Pre-processing steps

```
<div id="text">I've been in New York  
in 2011, but didn't like it. I  
preferred Los Angeles.</div>
```



*MINIMAL
VARIATION*

"BAG OF WORDS"

new_york 0000 like

prefer los_angeles

Parts of Speech

POS tagging

Grassfed highland Chianina beef with handcut fries and seasonal micro greens 29,—

Rich, tender, golden-brown beef with **crisp** fries and **tender** greens 18,—

Savory beef with **delicious** fries and **tasty** salad 12,—

ADJs = price?

POS tagging

POS

PRON

VERB

PROPN

PROPN

PUNCT

|

|

|

|

|

I

admire

Rosa Parks

.

POS tagging

Open class words	Closed class words	Other
ADJ adjectives: <i>awesome, red</i> ADV adverbs: <i>quietly, where, never</i> INTJ interjections: <i>ouch, shhh</i> NOUN nouns: <i>book, war</i> PROPN proper nouns: <i>Rosa, Twitter</i> VERB full verbs: <i>(she) codes, (they) submitted</i>	ADP adpositions: <i>over, before</i> AUX auxiliary/modal verbs: <i>have (been), could (do), will (change)</i> CCONJ coordinating conjunctions: <i>and, or, but</i> DET determiners: <i>a, they, which</i> NUM numbers. Exactly what you would think it is... PART particles: <i>'s</i> PRON pronouns: <i>you, her, myself</i> SCONJ subordinating conjunctions: <i>since, if, that</i>	PUNCT punctuation marks: <i>!, ?, –</i> SYM symbols: <i>%, \$, :)</i> x other: <i>pfffrt</i>

POS tagging

show {VERB, NOUN}

PART **show**
show
PRON **show**
show

DET **show**
show
show
ADJ **show**
show

Structured prediction: depends on the POS of a previous word

Parsing

Dependency Parsing

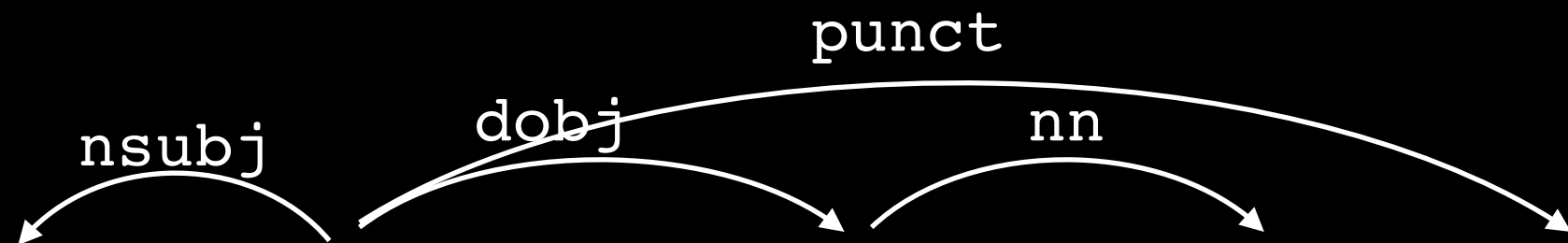
Facebook eventually  acquire(Facebook, WhatsApp) after hard negotiations.

WhatsApp was acquired  acquire(Facebook, WhatsApp) by Facebook.

Facebook subsidiary  acquire(WhatsApp, look) WhatsApp to acquire new look.

Dependency Parsing

PARSING

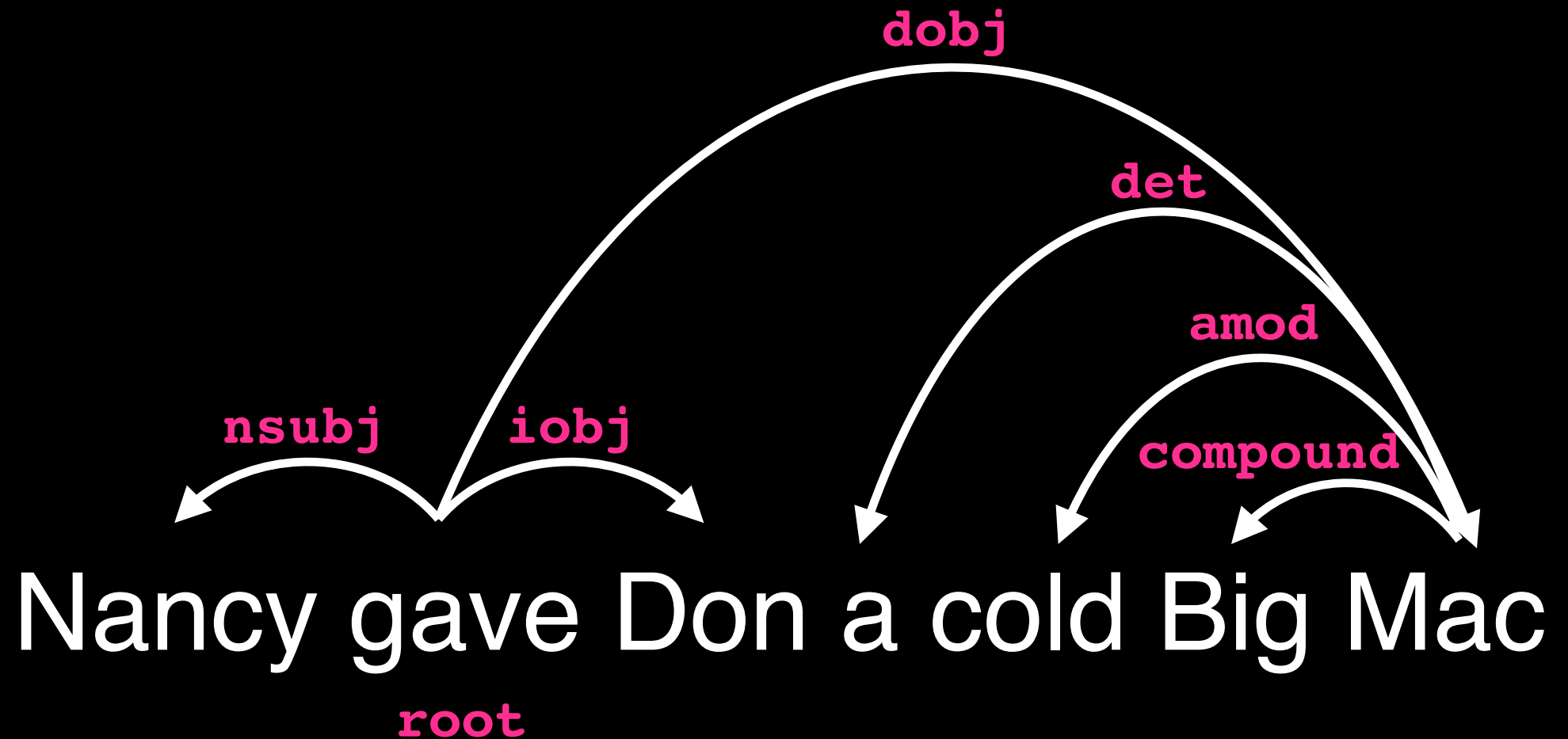


POS

PRON	VERB	PROPN	PROPN	PUNCT
I	admire	Rosa	Parks	.

Dependency Parsing

acl: adjectival clause
advcl: adverbial clause modifier
advmod: adverbial modifier
amod: adjectival modifier
appos: appositional modifier
aux: auxiliary
case: case marking
cc: coordinating conjunction
ccomp: clausal complement
clf: classifier
compound: compound
conj: conjunct
cop: copula
csbj: clausal subject
dep: unspecified dependency
det: determiner
discourse: discourse element
dislocated: dislocated elements
dobj: direct object
expl: expletive
fixed: fixed multiword expression
flat: flat multiword expression
goeswith: goes with
iobj: indirect object
list: list
mark: marker
nmod: nominal modifier
nsbj: nominal subject
nummod: numeric modifier
obl: oblique nominal
orphan: orphan
parataxis: parataxis
punct: punctuation
reparandum: overridden disfluency
root: root
vocative: vocative
xcomp: open clausal complement



Named Entities

Named Entities

Support The Guardian | Search jobs | Sign in | Search | International edition

Contribute → Subscribe →

The Guardian

News | Opinion | Sport | Culture | **Lifestyle** | More

Travel ► UK Europe US

Observer spring breaks
City breaks

Jane Dunford, Chris Moss, Mary Novakovich, Cella Topping

Mon 4 Feb 2019
11.00 GMT

1043

Spring breaks: 5 of the best cities in Europe



→ Places:

```
{ 'Ada',  
  'Antigone',  
  'Belgrade',  
  'Berlin',  
  'Constitución',  
  'Danube',  
  'Florence',  
  'France',  
  'Mikser',  
  'Rome',  
  'Santa Cruz',  
  'Savamala',  
  'Schlachtensee',  
  'Serbia',  
  'Spain',  
  'Tezga',  
  'Ville',  
  'Wannsee' }
```

Named Entities

NER

O

O

B-PERSON I-PERSON O

POS

PRON

VERB

PROPN

PROPN

PUNCT

|

|

|

|

|

I

admire

Rosa Parks

.

Named Entities

NE	Example
PERSON	
NORP (Nationality OR Religious or Political group)	
FAC (facility)	
ORG (organization)	
GPE (GeoPolitical Entity)	
LOC (locations, such as seas or mountains)	
PRODUCT	
EVENT (in sports, politics, history, etc.)	
WORK_OF_ART	
LAW	
LANGUAGE	
DATE	
TIME	
PERCENT	
MONEY	
QUANTITY	
ORDINAL	
CARDINAL (numbers)	

Wrapping up

Take Home Points

- NLP is a subfield of AI, using ML on linguistic problems to **explore, predict, and generate** text
- **Preprocessing** removes noise and unwanted variation
- Parts of speech (**POS**) denote a word's grammatical *category*
- **Parsing** denotes a word's grammatical *function*
- **Named entities** categorize a noun's semantic type