

# Natural Language Processing

## Lecture 01

Dirk Hovy

[dirk.hovy@unibocconi.it](mailto: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

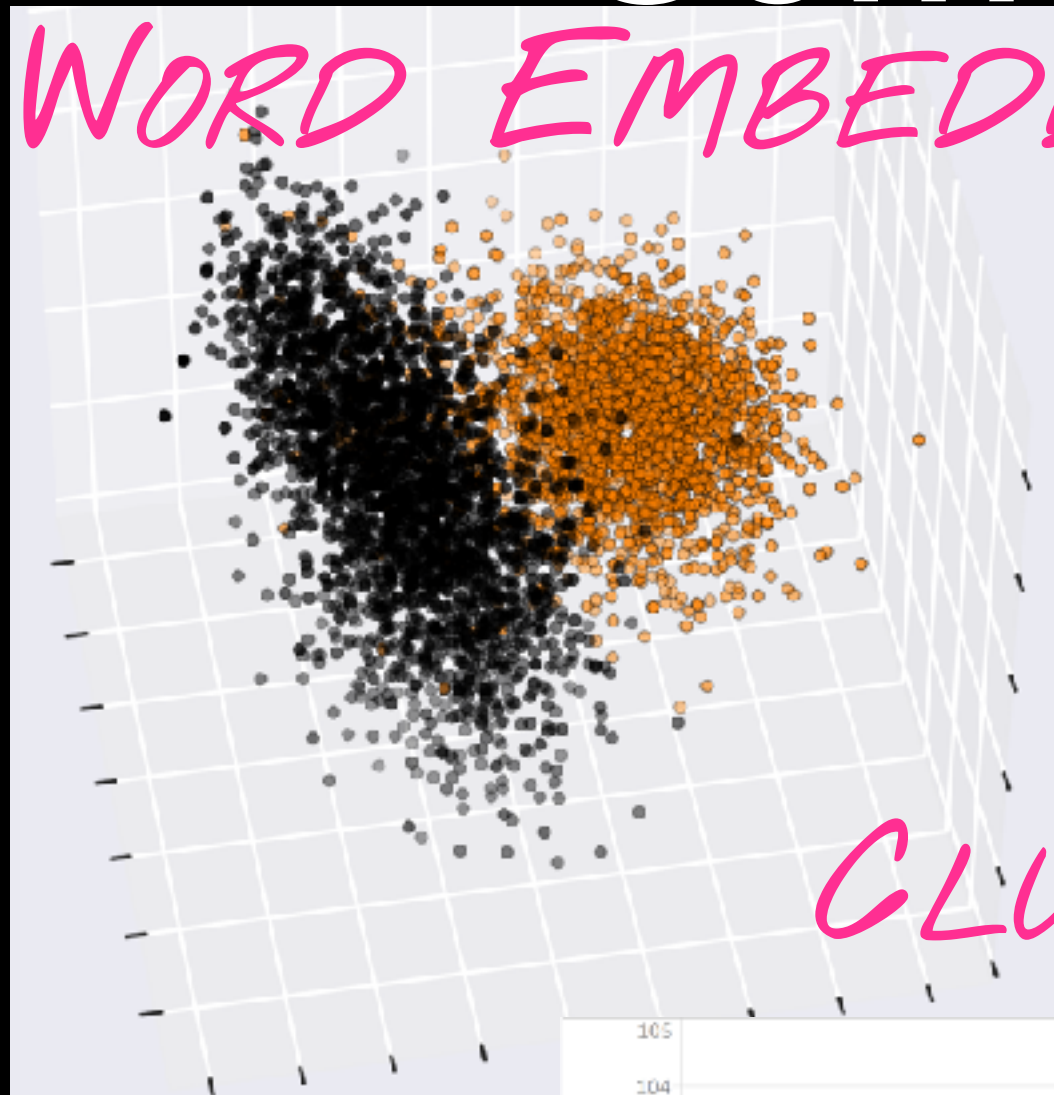
2023

2024

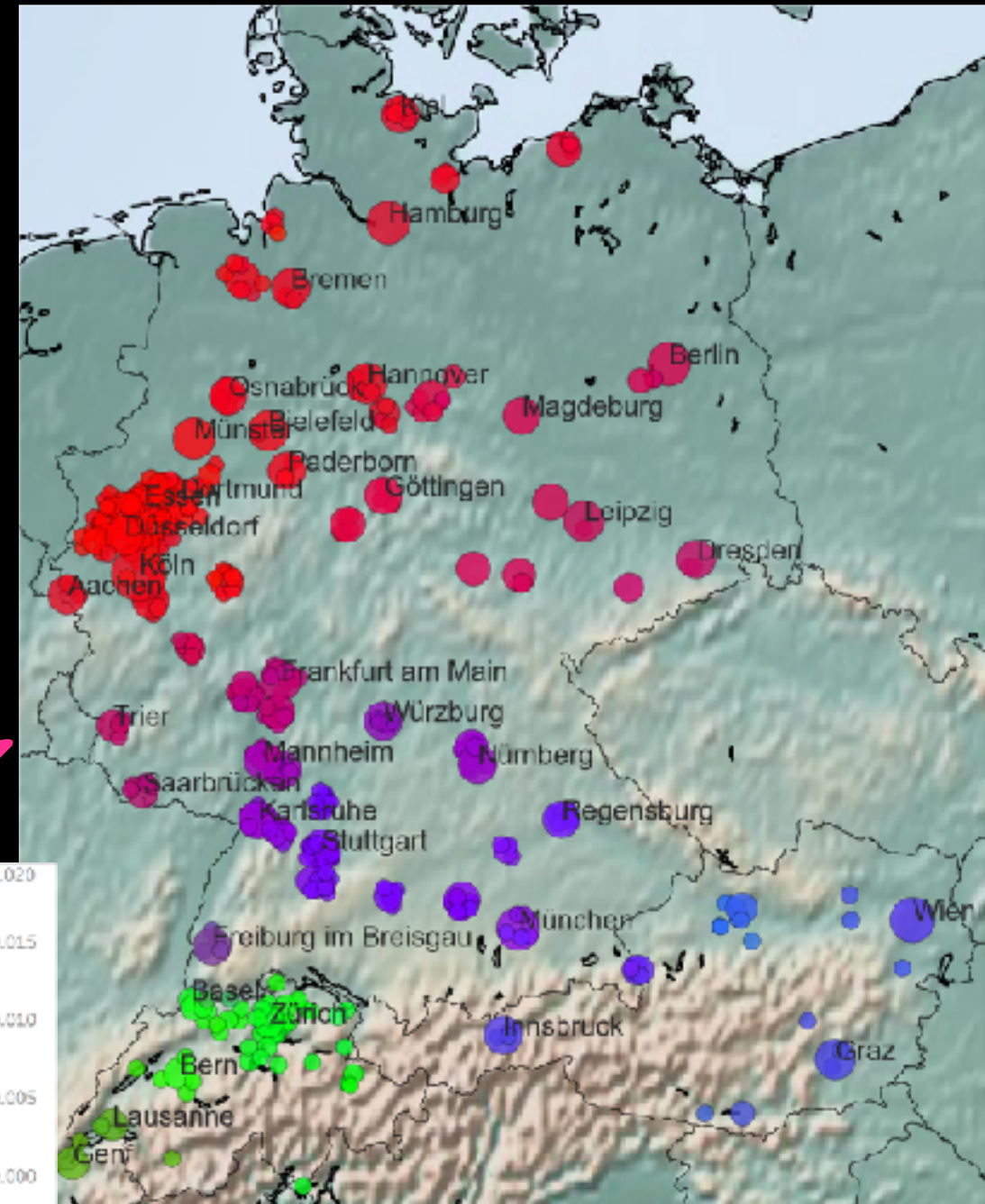
2025

# Some examples

WORD EMBEDDINGS



CLUSTERING



VISUALIZATION

SENTIMENT ANALYSIS



# Syllabus

Lesson	Topic
1	Intro 1
2	Intro 2
3	Representations
4	Embeddings
5	Information retrieval 1
6	Information retrieval 2
7	Language models 1
8	Language models 2
9	Topic models 1
10	Topic models 2
11	Dimensionality reduction and Clustering
12	Visualization
13	Retrofitting
14	Midterm Project Presentations
15	Text classification
16	Application: Sentiment Analysis
17	Improving classification performance
18	Application: Author attribute prediction
19	Neural networks basics
20	Multilayer Perceptron
21	Sequence models
22	RNNs and Bi-LSTMs
23	Final Project Presentations
24	Final Project Presentations

# Class structure

- Thursdays: intuition, theory, math (slides)
- Fridays: exercises and practice (Jupyter Notebooks)

# Grading

1. *Individual* midterm project (50%): Exploration and visualization
  2. Final *group* project (50%): Data annotation, exploration, visualization, and prediction
- Both projects are to be handed in as Jupyter Notebooks
  - Graded on data set size, annotation quality, correctness of implementations, performance of prediction
  - No point changes, only complete regrades (total can go down)!



# How do I succeed?

- Code well
- Pay attention
- Code some more



# What to do with a problem

1. Google it. [stackoverflow.com](https://stackoverflow.com) is your friend
2. Talk to your classmates
3. Ask the TA, Tommaso Fornaciari  
`fornaciari@unibocconi.it`
4. Make an appointment

*WARNING :*

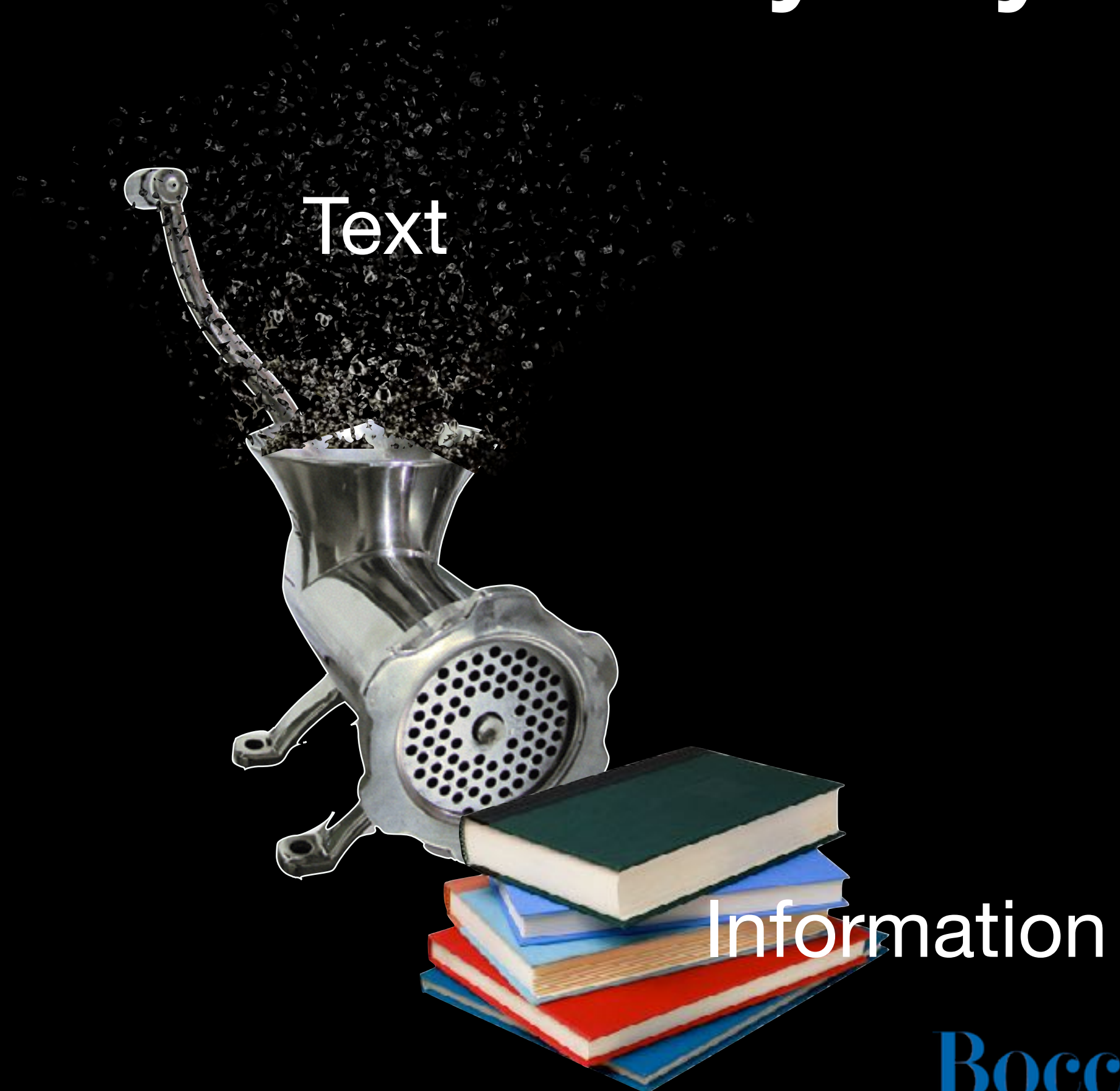
For any question we can solve with a Google search, we deduct points!

**Let's start!**

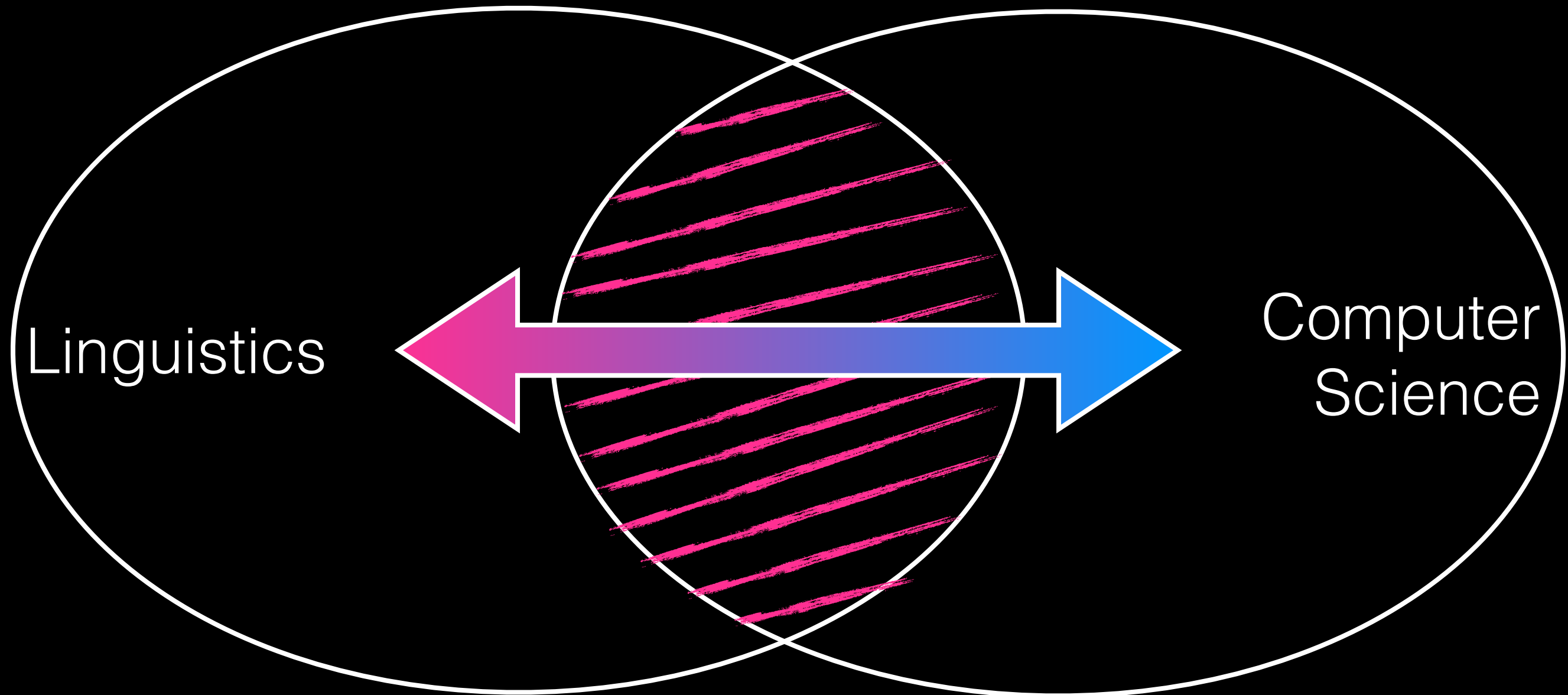
# Today's Goals

- 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?

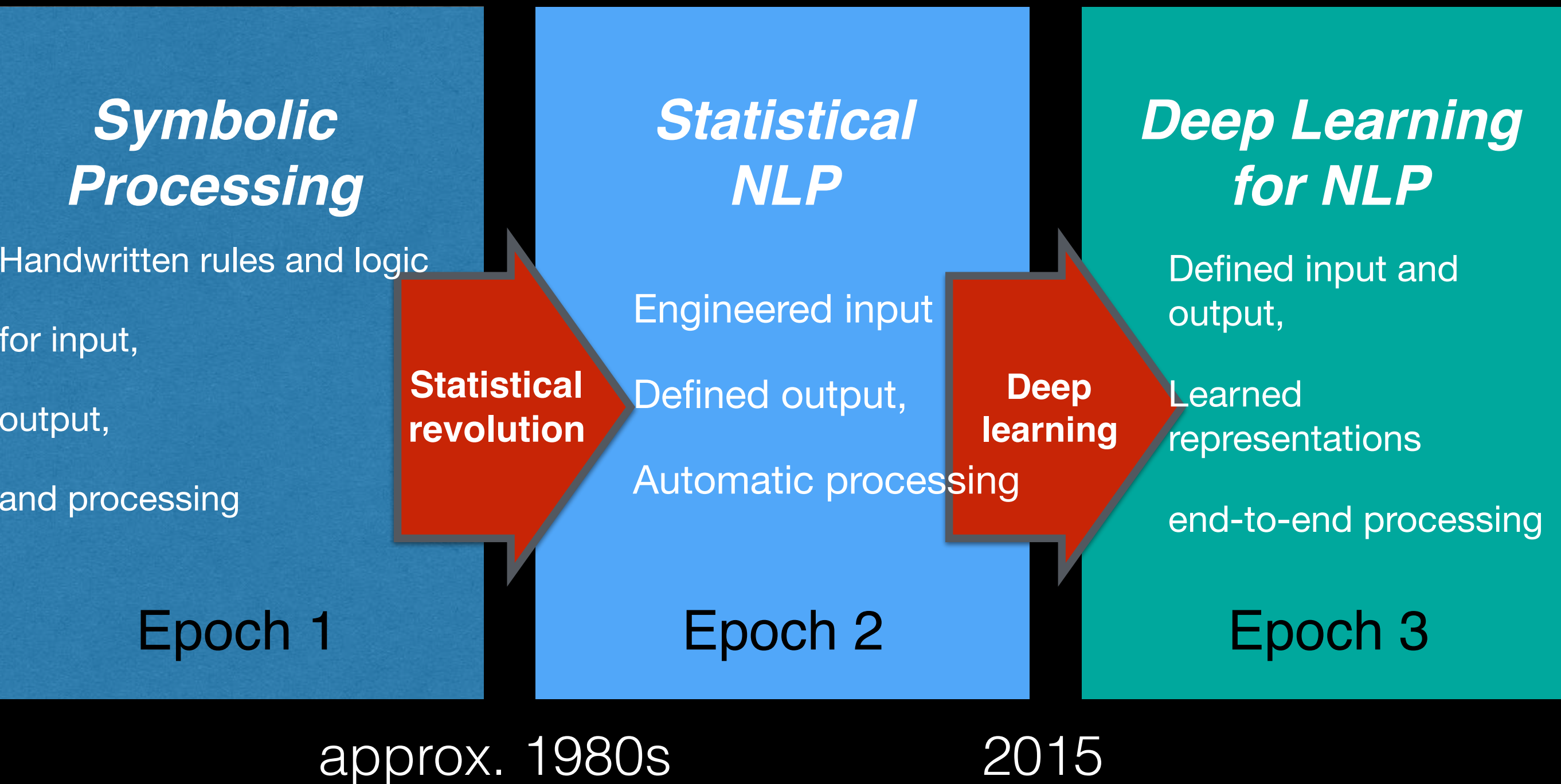


# The two sides of NLP



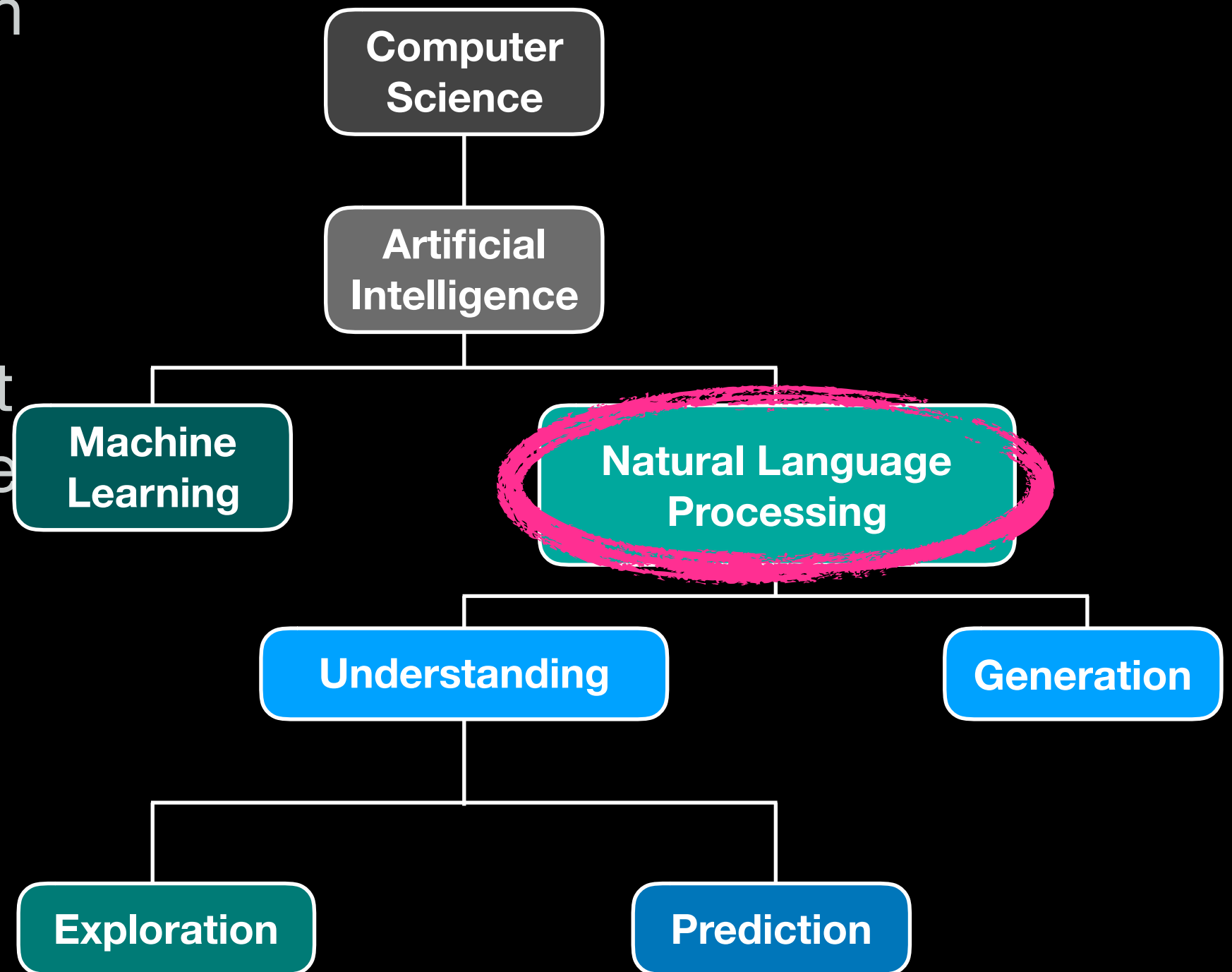
informed linguistic hypotheses    large-scale statistical analysis

# A very Brief History of NLP



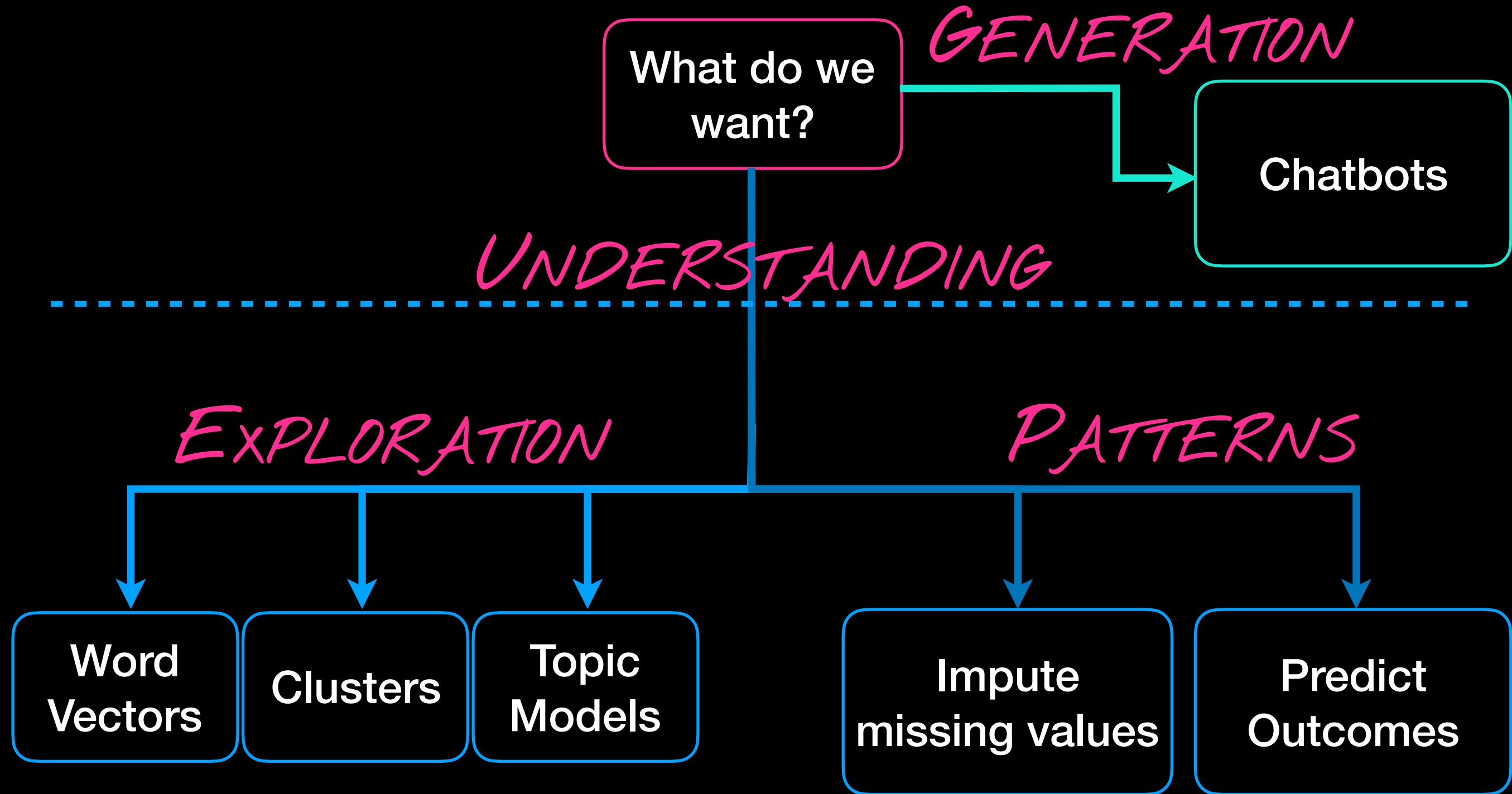
# Structure of NLP

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





# Two Uses of NLP



# Linguistic Analysis

# Examples of Analysis

# NER



# PERSON

# PERSON



# PARSING

nsubj

~~dob-j~~

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)

i new york 0000 , like .

- Segment sentences

- Tokenize words

i prefer los angeles .

- Normalize words

- numbers

- lemmas vs. stems

- Remove unwanted words

- stopwords

- content words (use POS tagging!)

- join collocations



# 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

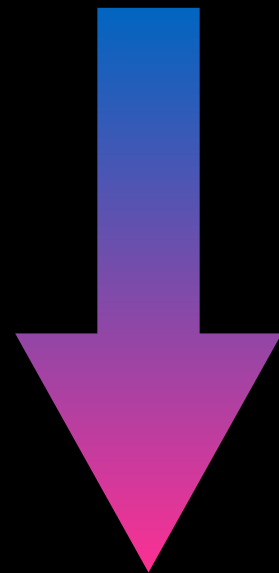
- Remove formatting (e.g. HTML)
- Segment sentences
- Tokenize words
- Normalize words
  - numbers
  - lemmas vs. stems
- Remove unwanted words
  - stopwords
  - content words (use POS tagging!)
- join collocations

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**

|

I

**VERB**

|

admire

**PROPN**

|

Rosa

**PROPN**

|

Parks

**PUNCT**

|

.

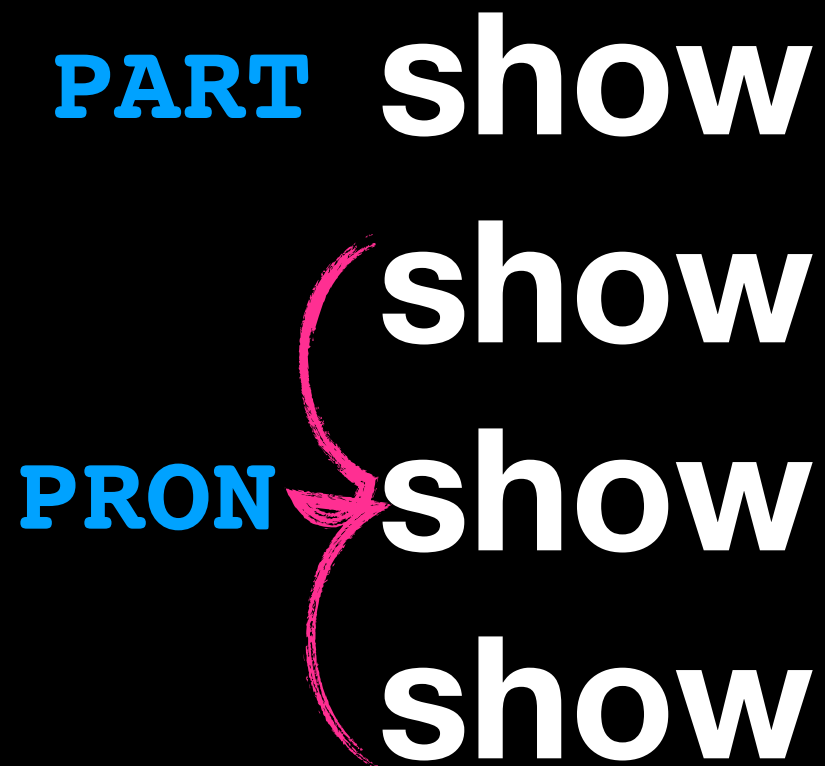
# POS tagging

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

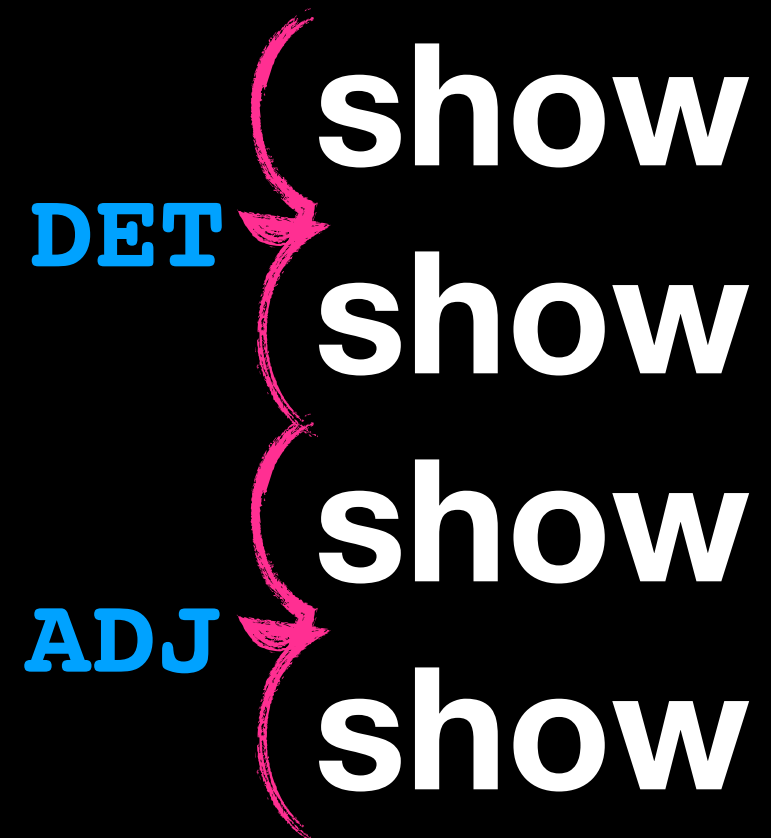
# POS tagging

**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

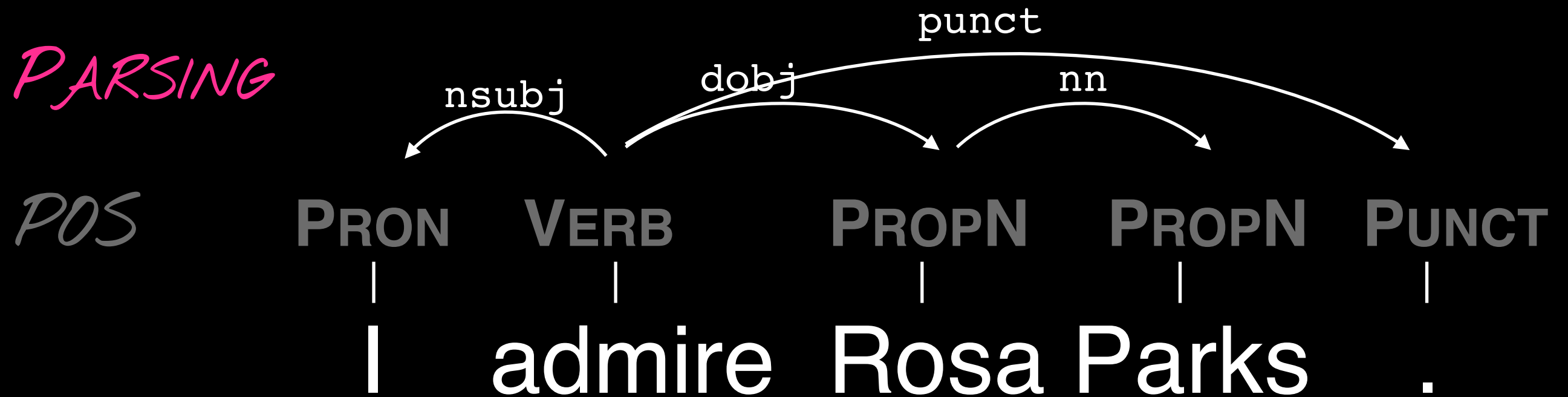
Facebook eventually  acquire(Facebook, WhatsApp)  
acquired WhatsApp after  
hard negotiations.

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

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

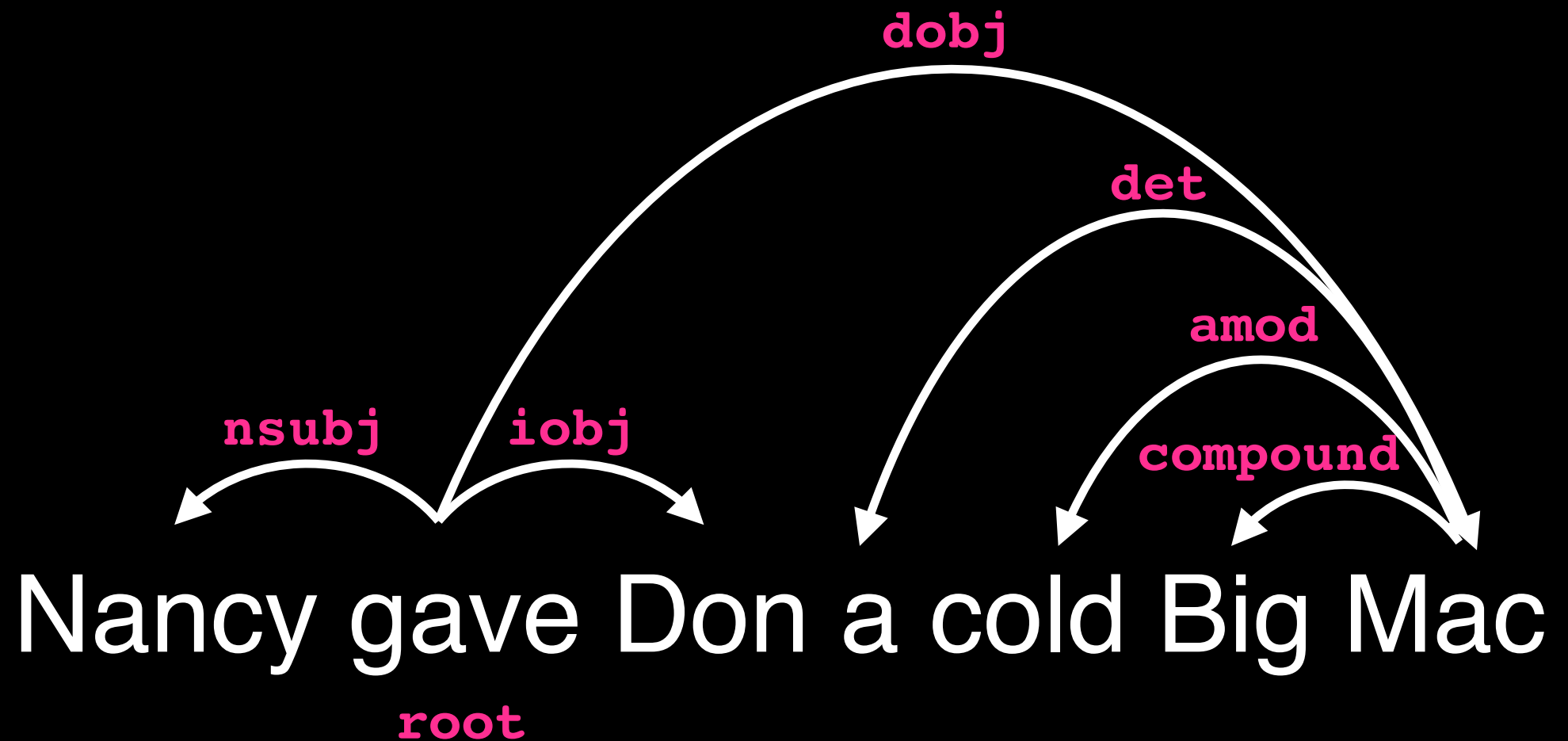


# Dependency Parsing



# 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