

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

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

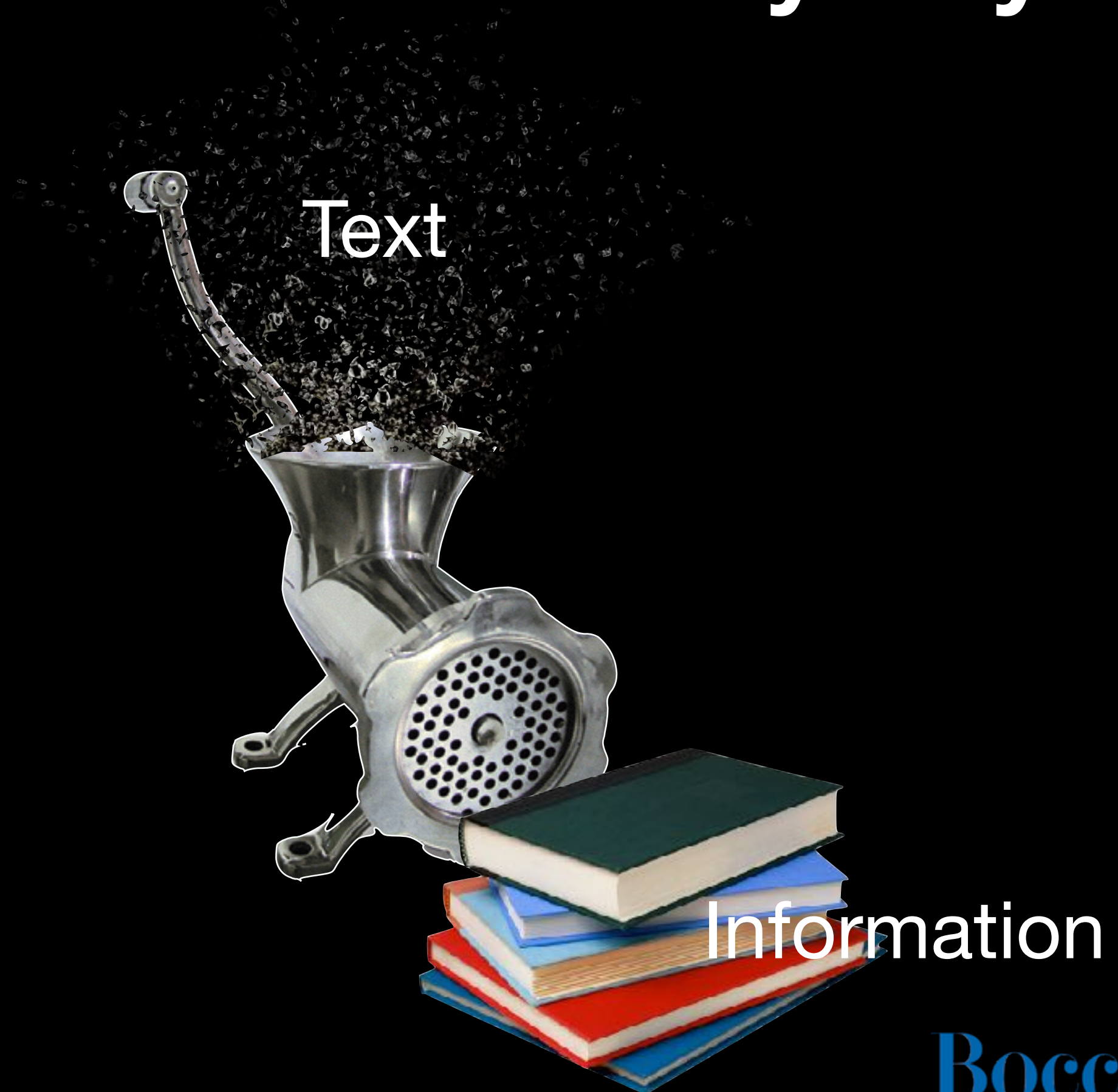
2022

2023

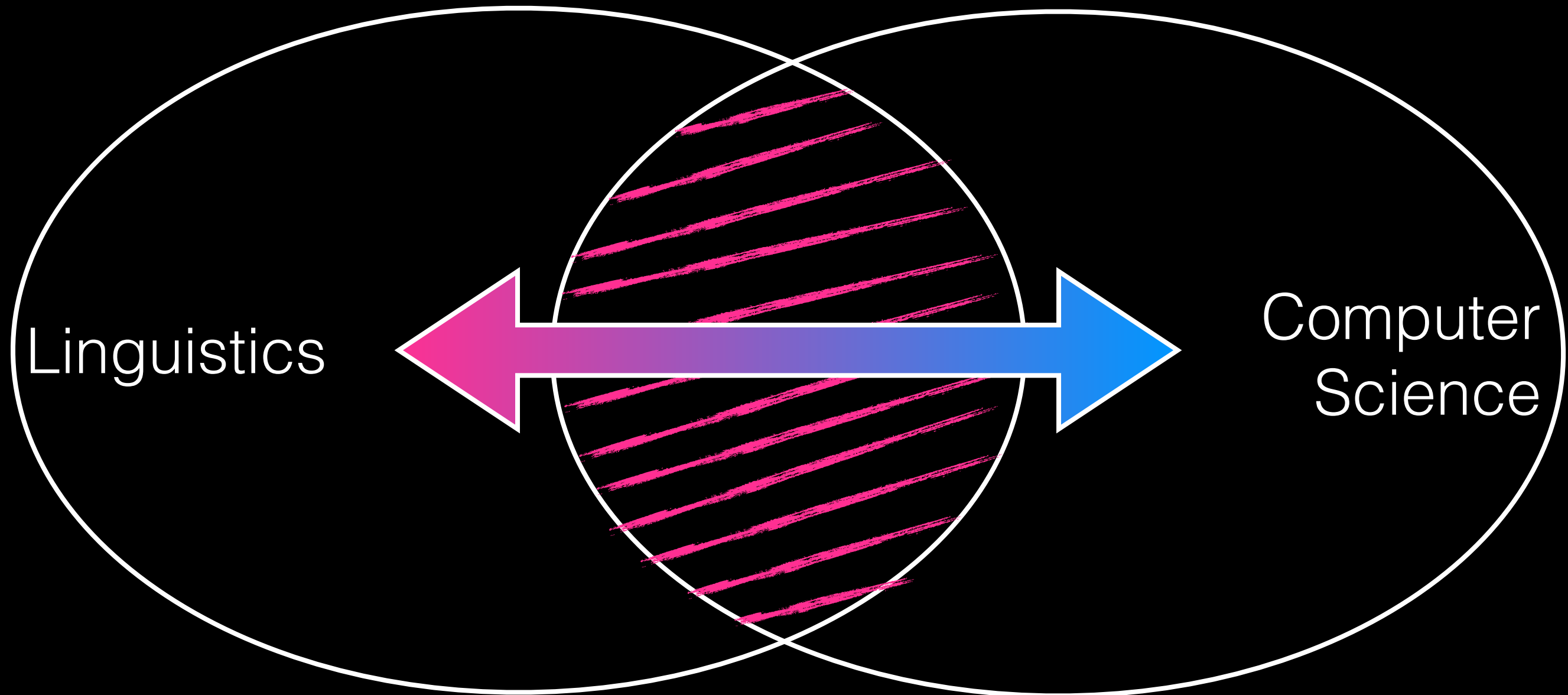
2024

2025

So, what's NLP anyway?



The two sides of NLP

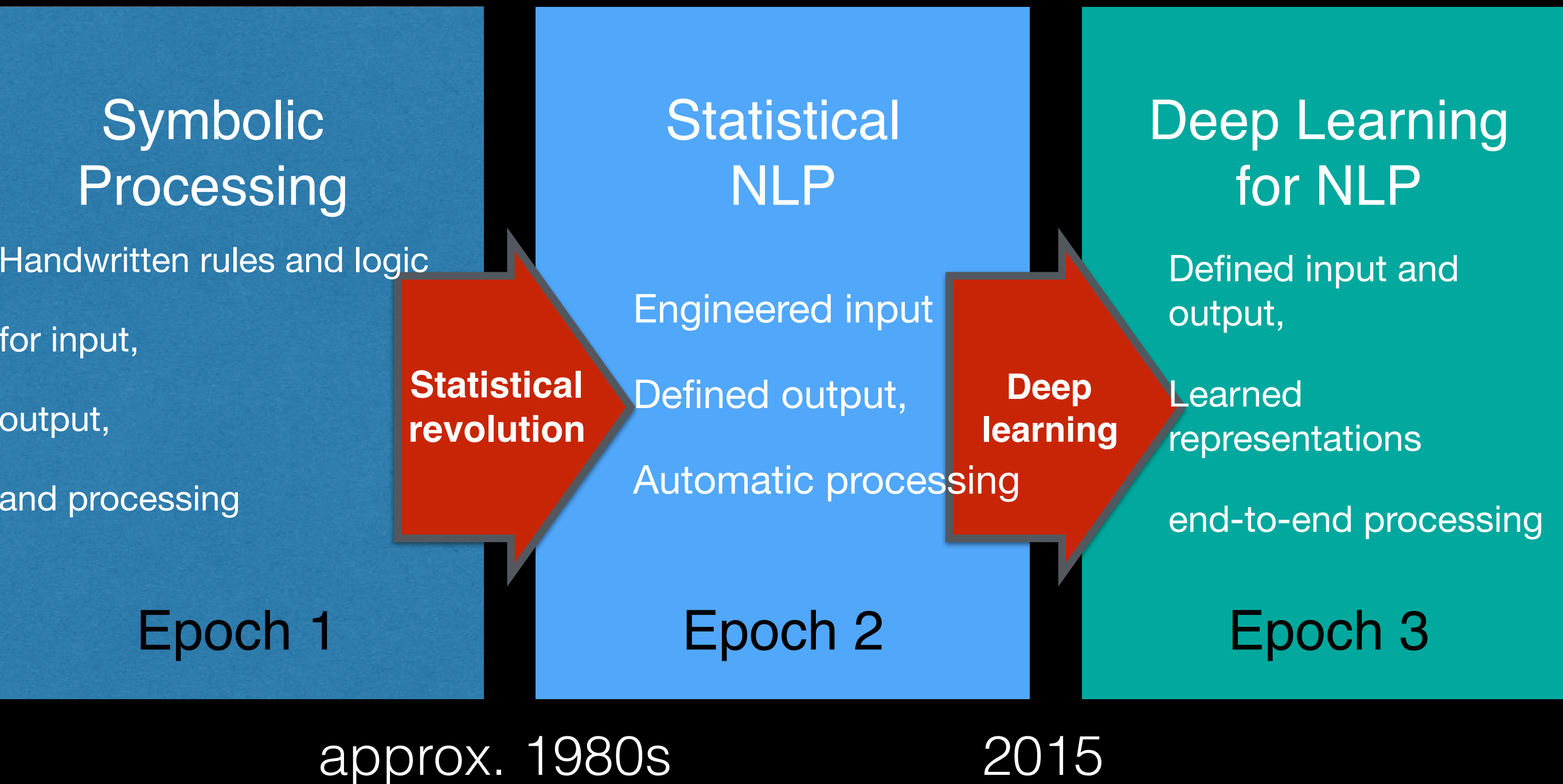


Linguistics

Computer
Science

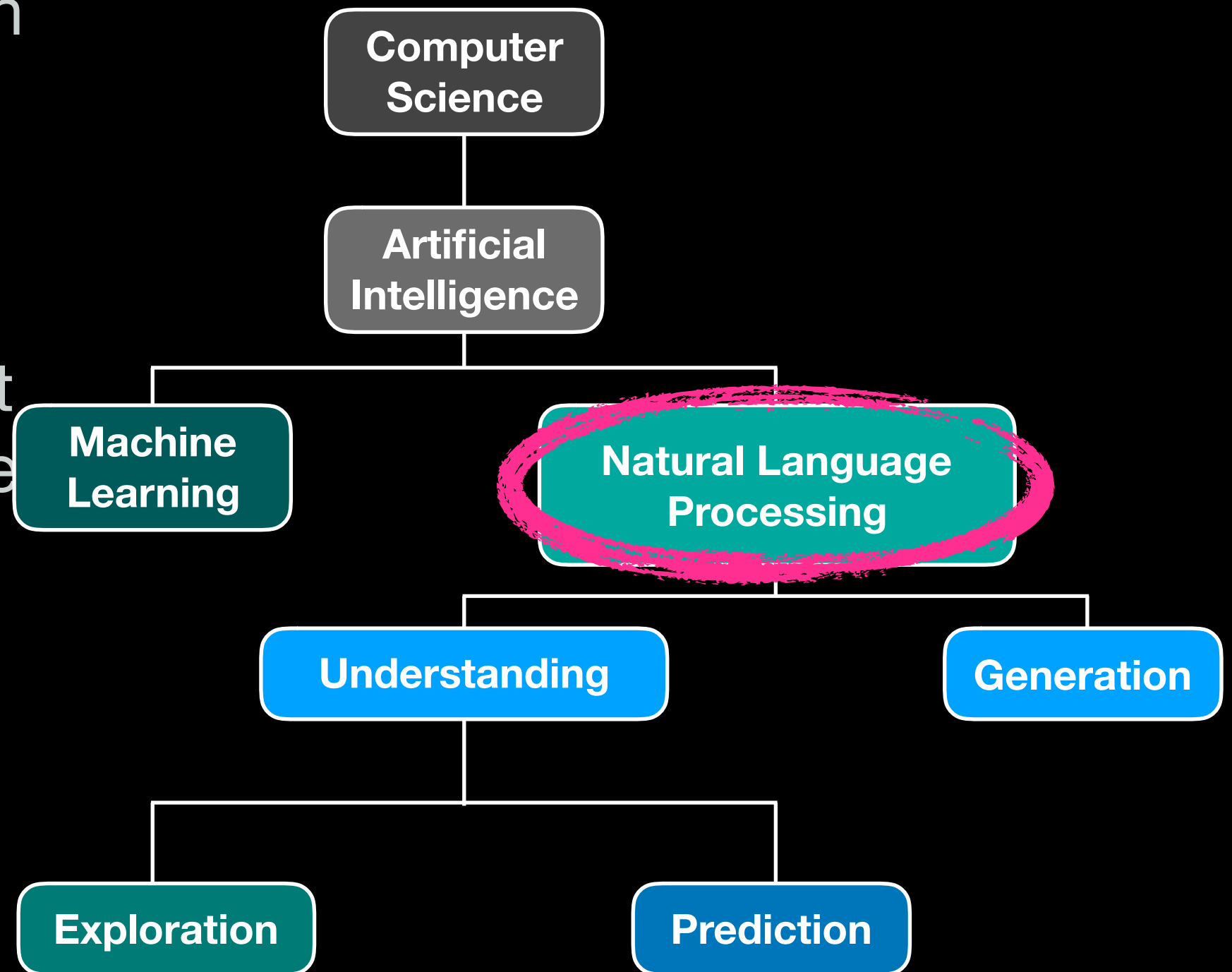
informed linguistic hypotheses large-scale statistical analysis

A very Brief History of NLP

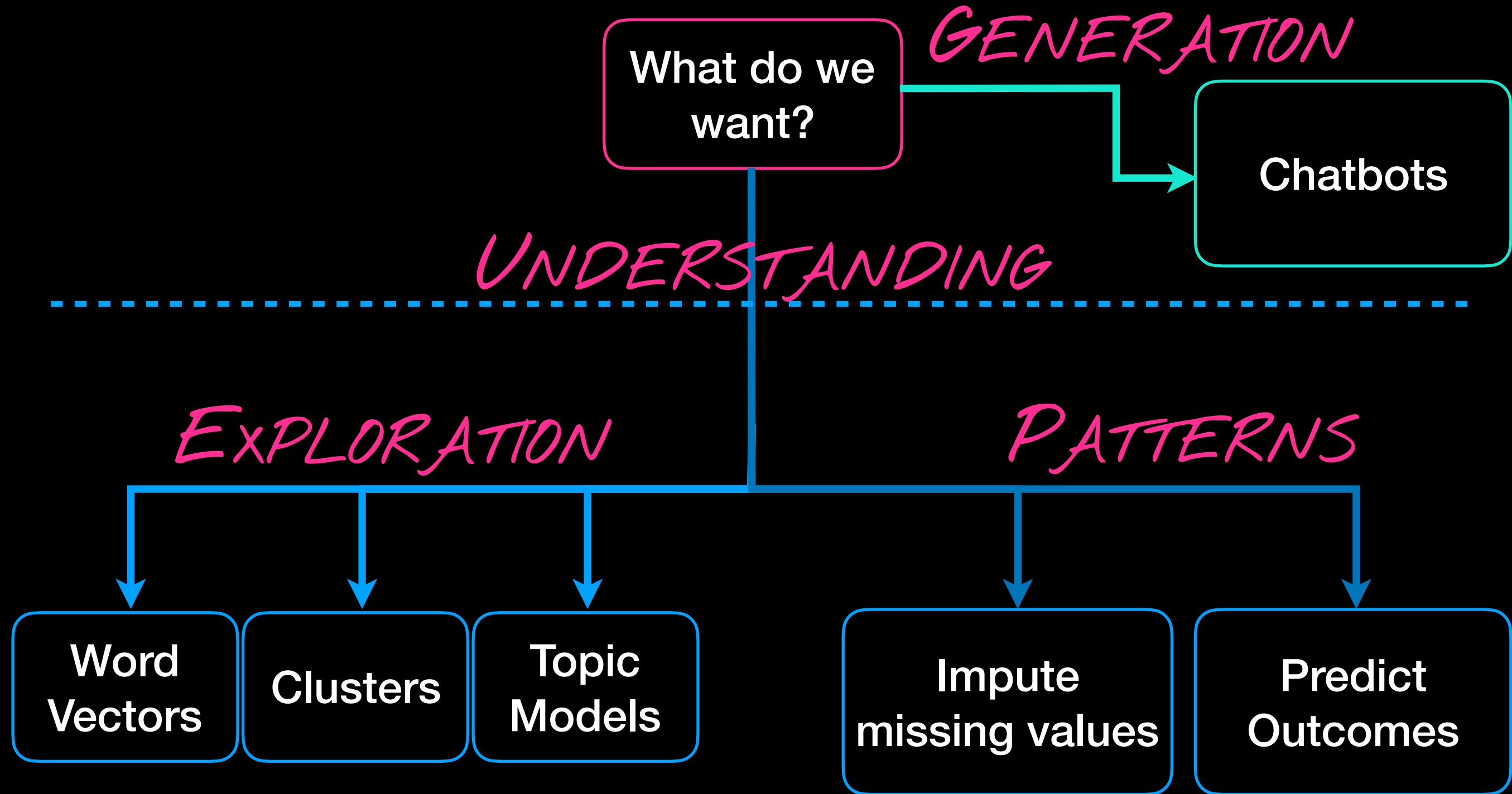


Structure of NLP

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

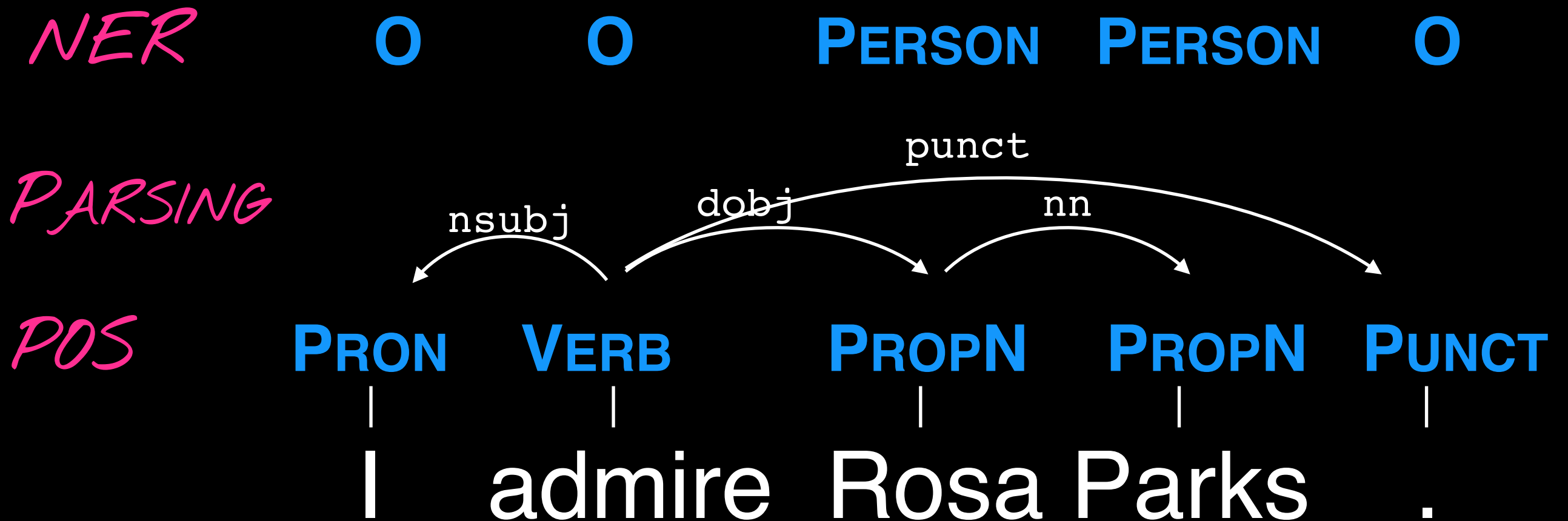


Two Uses of NLP



Linguistic Analysis

Examples of Analysis



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)

`new_york 0000 like`

- Segment sentences

- Tokenize words

`prefer los_angeles`

- Normalize words

- numbers

- lemmas vs. stems

- Remove unwanted words

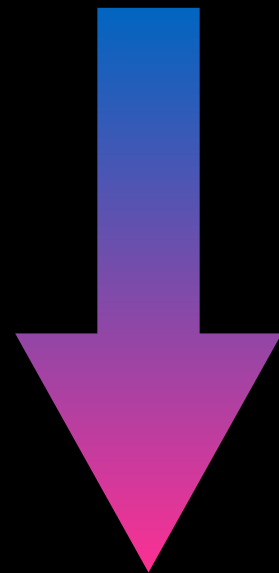
- stopwords

- content words (use POS tagging!)

- join collocations

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

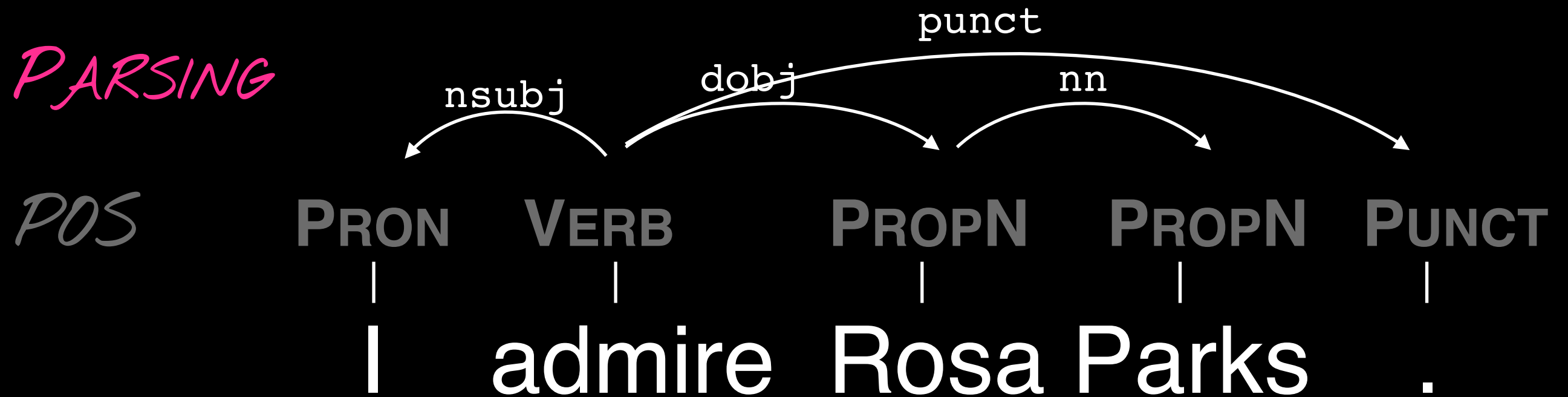
Dependency Parsing

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

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

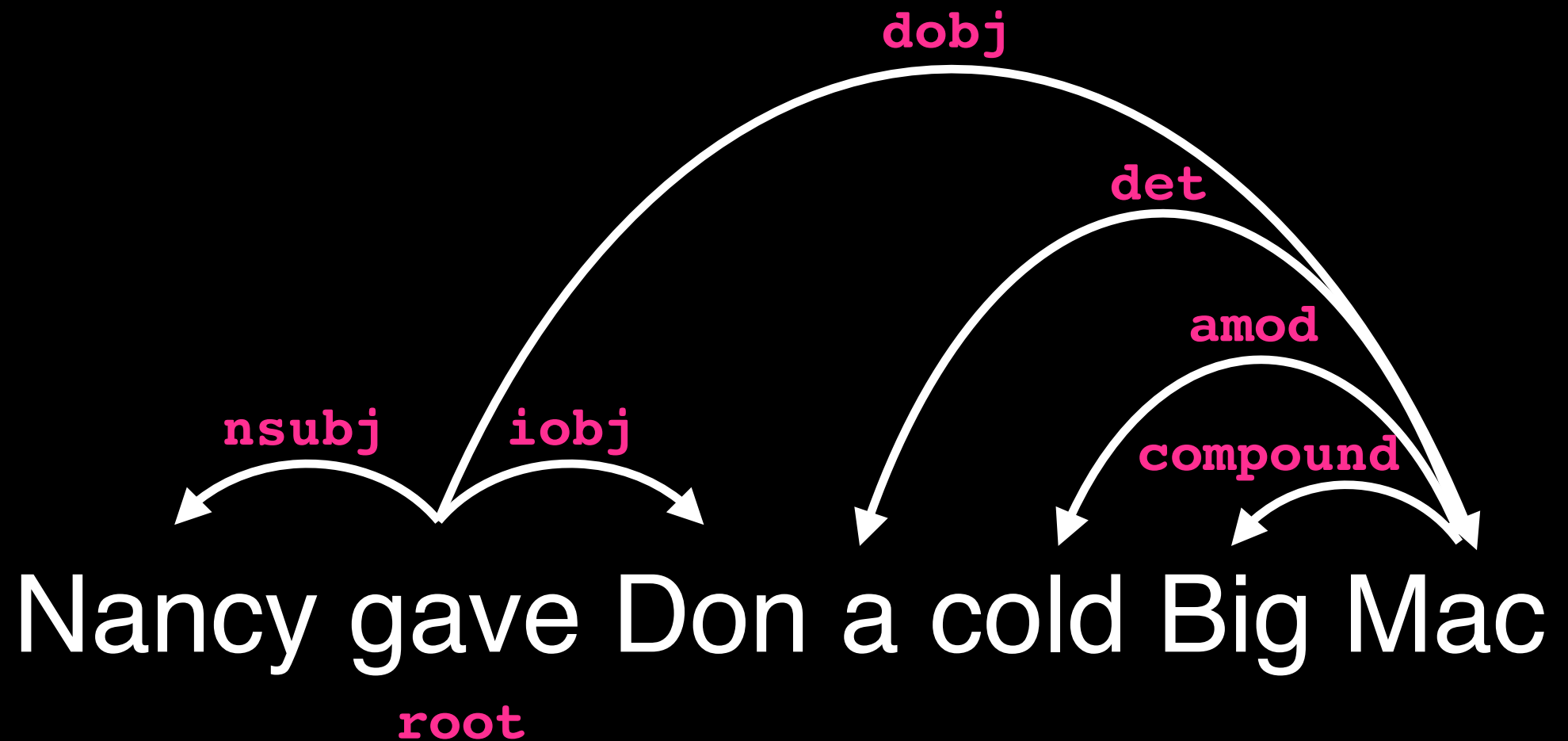
Facebook subsidiary  acquire(WhatsApp, look)
WhatsApp to acquire new
look.

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