

Why NLP is Hard?

1. Ambiguity
2. Scale
3. Sparsity
4. Variation
5. Expressivity
6. Unmodeled Variables
7. Unknown representations



Why NLP is Hard?

1. Ambiguity
2. Scale
3. Sparsity
4. Variation
5. Expressivity
6. Unmodeled Variables
7. Unknown representations



Ambiguity

- Ambiguity at multiple levels
 - Word senses: **bank** (finance or river ?)
 - Part of speech: **chair** (noun or verb ?)
 - Syntactic structure: **I can see a man with a telescope**
 - Multiple: **I made her duck**





“One morning I shot
an elephant in my pajamas”

I made her duck

[SLP2 ch. 1]

- I cooked waterfowl for her
- I cooked waterfowl belonging to her
- I created the (plaster?) duck she owns
- I caused her to quickly lower her head or body
- ...

The Challenges of “Words”

- Segmenting text into words
- Morphological variation
- Words with multiple meanings: bank, mean
- Domain-specific meanings: latex
- Multiword expressions: make a decision, take out, make up

Part of Speech Tagging

ikr smh he asked fir yo last name

so he can add u on fb lololol

Part of Speech Tagging

I know, right shake my head for your
ikr smh he asked fir yo last name

so he can add you on Facebook laugh out loud
u fb lololol

Part of Speech Tagging

I know, right

ikr

!

interjection

shake my head

smh

G

acronym

he

O

pronoun

asked

V

verb

for

fir

P

prep.

your

yo

D

det.

last

A

adj.

name

N

noun

so

P

preposition

he

O

can

V

add

V

you

u

O

on

P

Facebook

fb

^

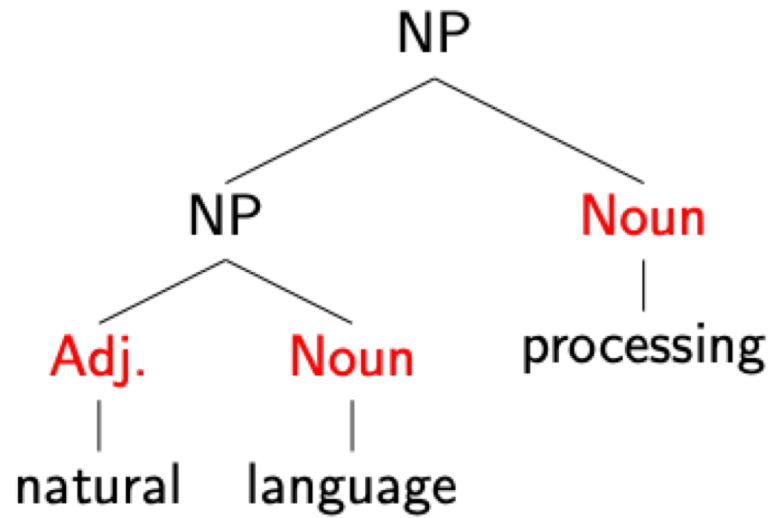
proper noun

laugh out loud

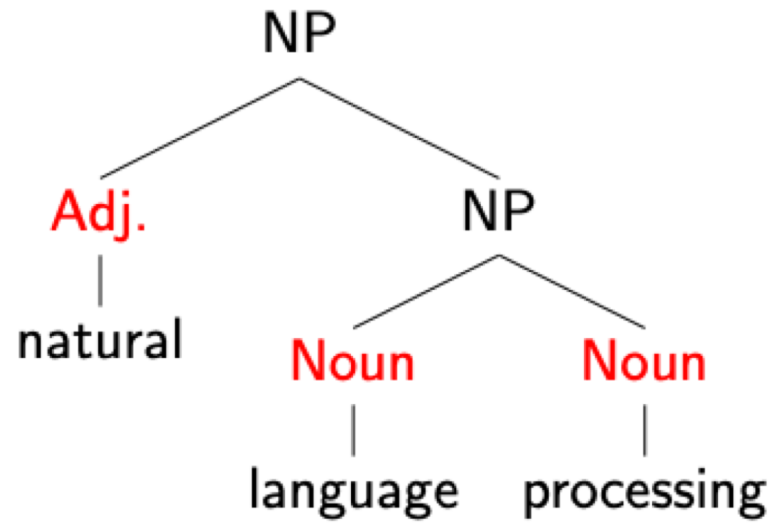
lololol

!

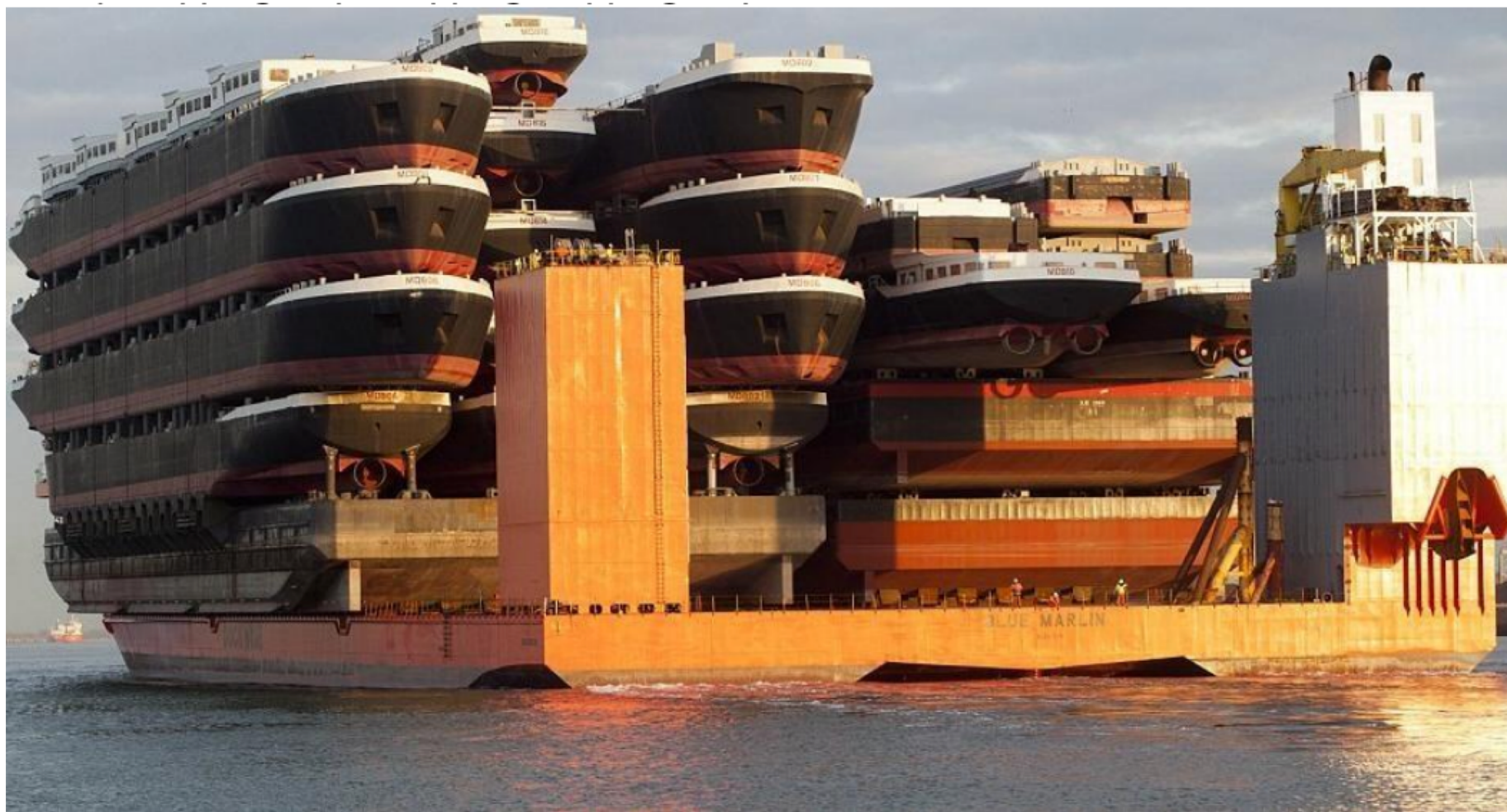
Syntax



vs.



Morphology + Syntax



A ship-shipping
ship, shipping
shipping-ships

Semantics

- Every fifteen minutes a woman in this country gives birth. Our job is to find this woman, and stop her!

– Groucho Marx



Syntax + Semantics

- We saw the woman with the telescope wrapped in paper.
 - Who has the telescope?
 - Who or what is wrapped in paper?
 - An even of perception, or an assault?

Dealing with Ambiguity

- How can we model ambiguity?
 - Non-probabilistic methods (CKY parsers for syntax) return **all possible analyses**
 - Probabilistic models (HMMs for POS tagging, PCFGs for syntax) and algorithms (Viterbi, probabilistic CKY) return **the best possible analyses**, i.e., the most probable one
- But the “best” analysis is only good if our probabilities are accurate. Where do they come from?

Corpora

- A corpus is a collection of text
 - Often annotated in some way
 - Sometimes just lots of text
- Examples
 - Penn Treebank: 1M words of parsed WSJ
 - Canadian Hansards: 10M+ words of French/English sentences
 - Yelp reviews
 - The Web!



Rosetta Stone

Statistical NLP

- Like most other parts of AI, NLP is dominated by statistical methods
 - Typically more robust than rule-based methods
 - Relevant statistics/probabilities are **learned from data**
 - Normally requires lots of data about any particular phenomenon

Why NLP is Hard?

1. Ambiguity
2. Scale
3. **Sparsity**
4. Variation
5. Expressivity
6. Unmodeled Variables
7. Unknown representations



Sparsity

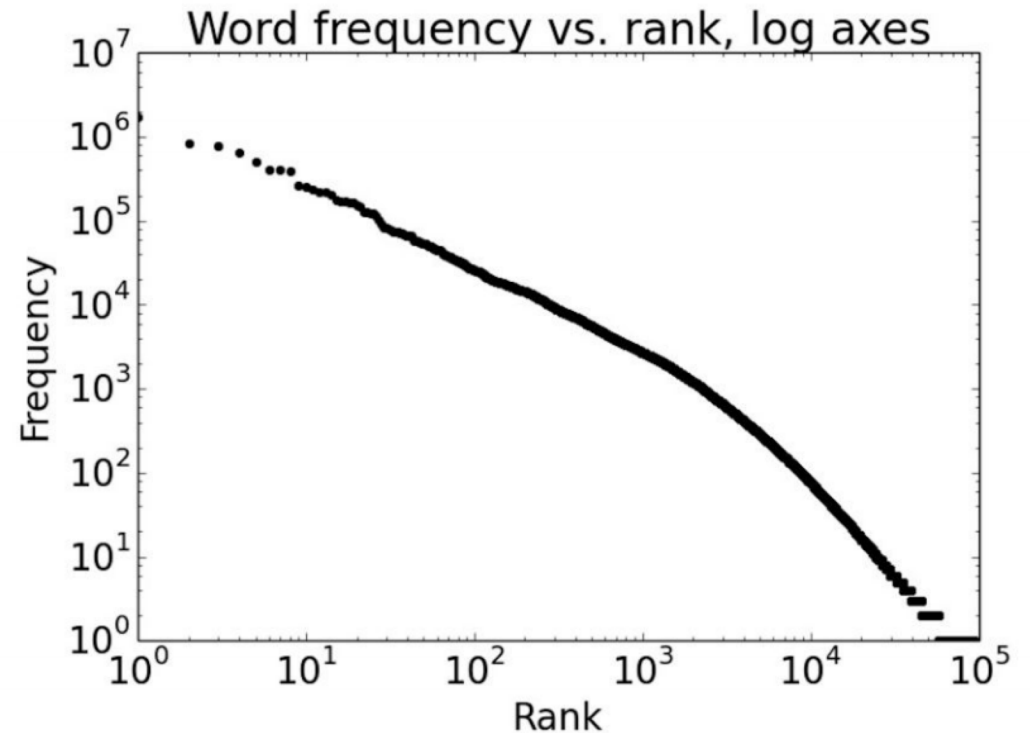
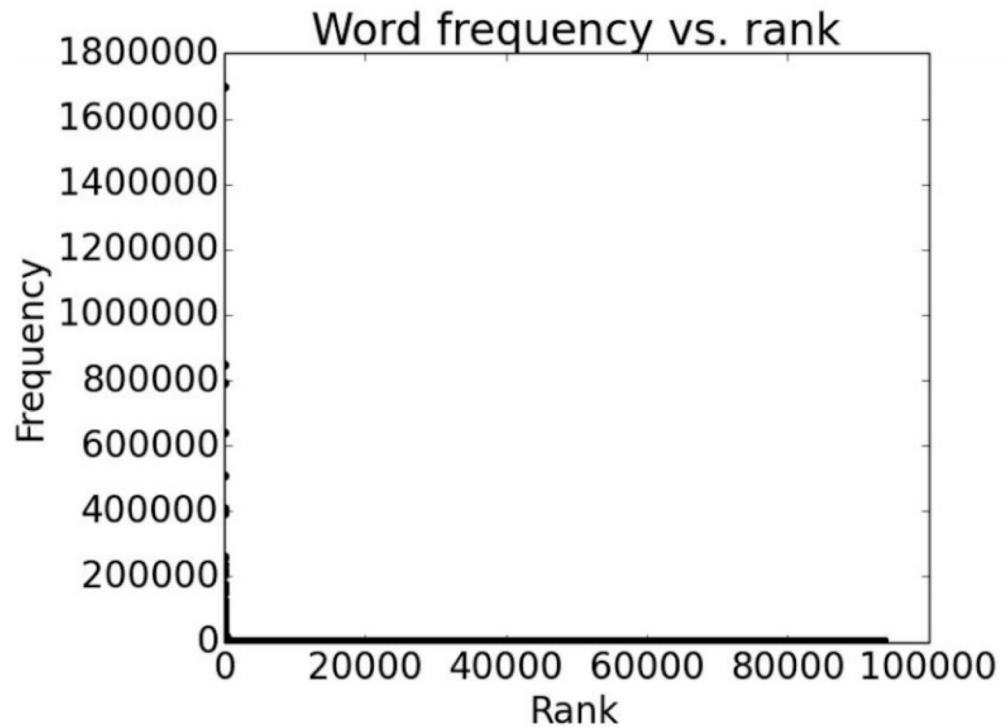
- Sparse data due to **Zipf's Law**
- Example: the frequency of different words in a large text corpus

any word	
Frequency	Token
1,698,599	the
849,256	of
793,731	to
640,257	and
508,560	in
407,638	that
400,467	is
394,778	a
263,040	I

nouns	
Frequency	Token
124,598	European
104,325	Mr
92,195	Commission
66,781	President
62,867	Parliament
57,804	Union
53,683	report
53,547	Council
45,842	States

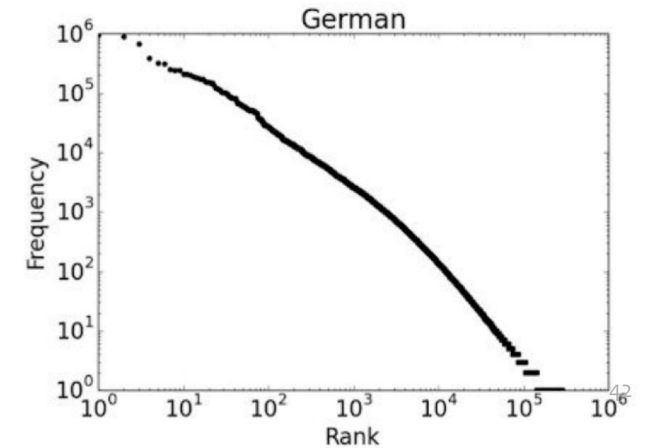
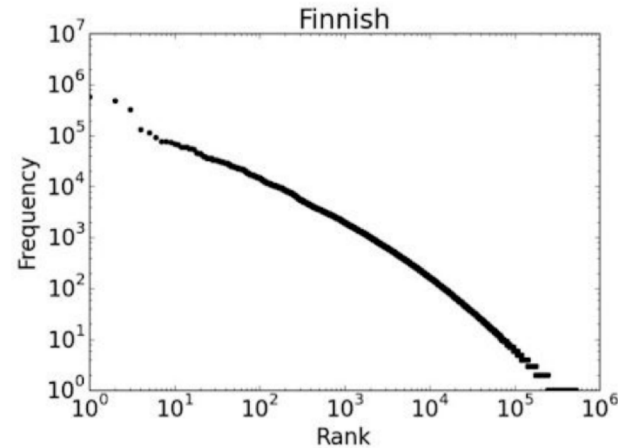
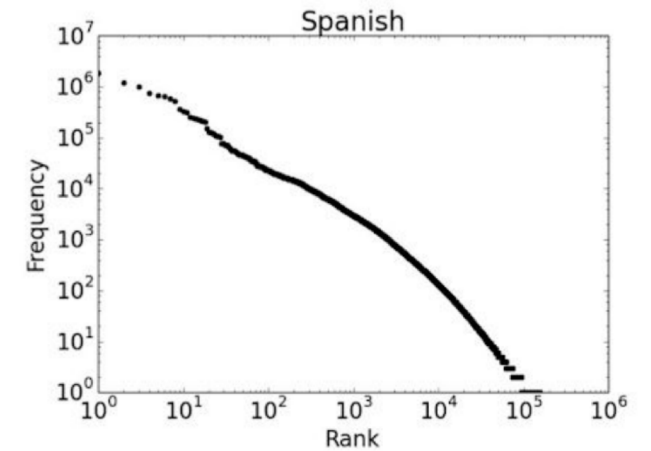
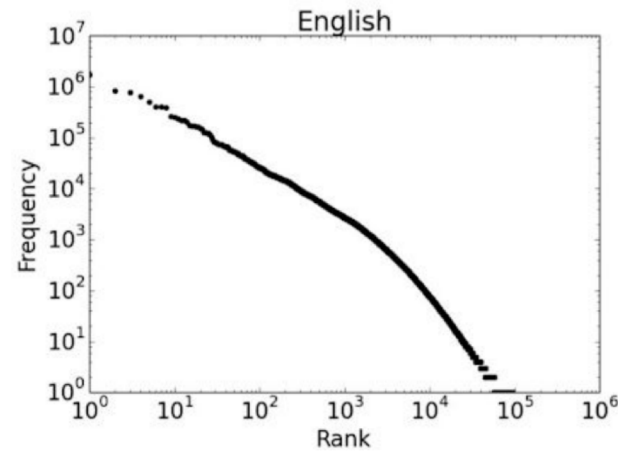
Sparsity

- Order words by frequency. What is the frequency of n th ranked word?



Sparsity

- Regardless of how large our corpus is, there will be a lot of infrequent words
- This means we need to find clever ways to estimate probabilities for things we have rarely or never seen



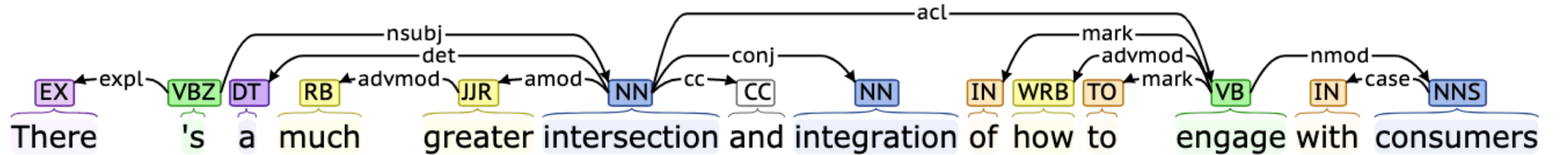
Why NLP is Hard?

1. Ambiguity
2. Scale
3. Sparsity
4. **Variation**
5. Expressivity
6. Unmodeled Variables
7. Unknown representations



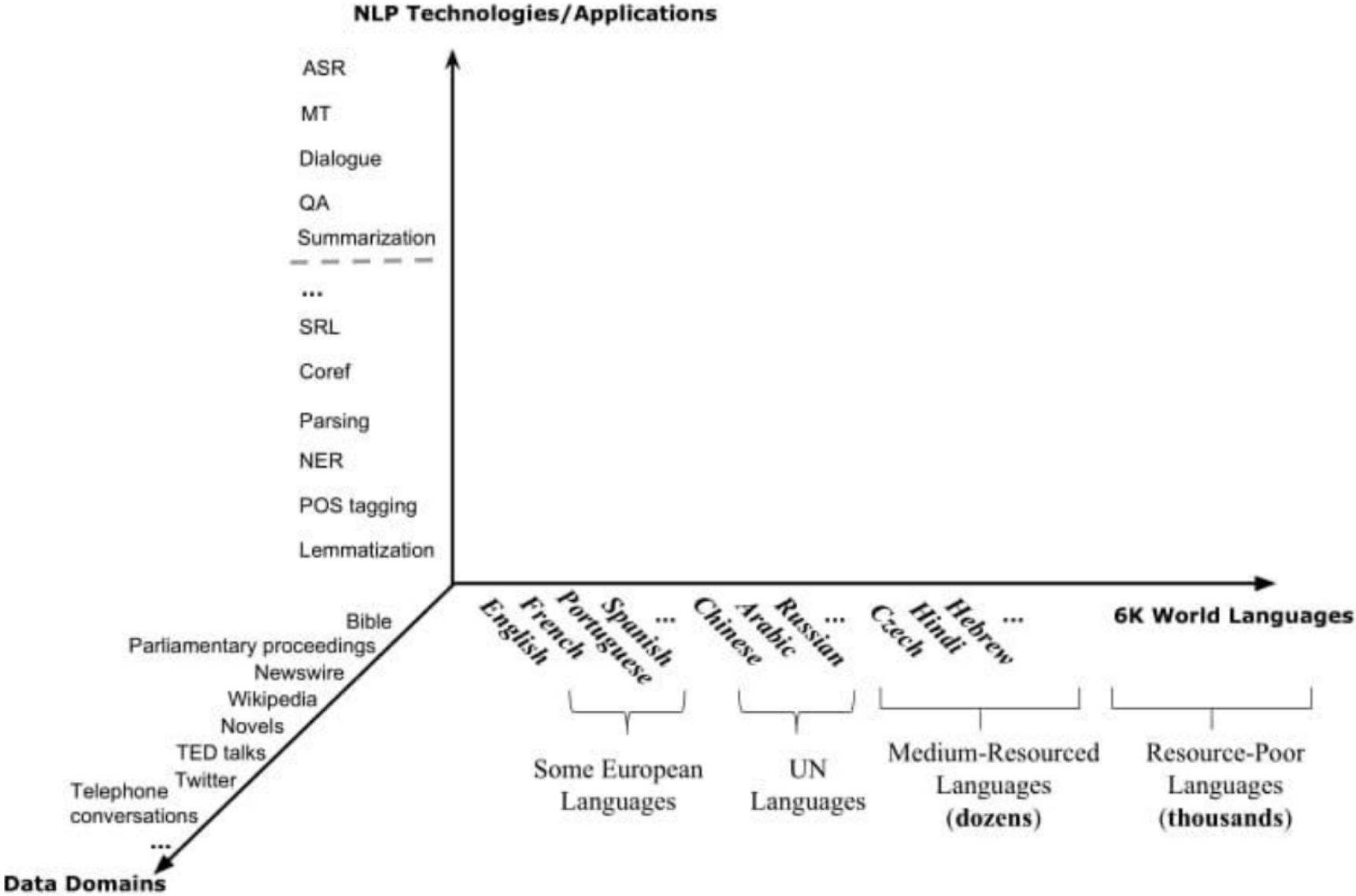
Variation

- Suppose we train a part of speech tagger or a parser on the **Wall Street Journal**



- What will happen if we try to use this tagger/parser for **social media**?
 - *“ikr smh he asked fir yo last name so he can add u on fb lololol”*

Variation



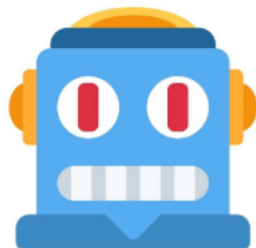
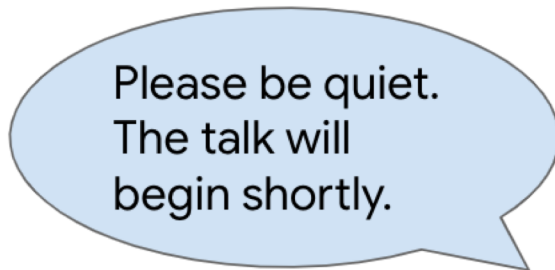
Why NLP is Hard?

1. Ambiguity
2. Scale
3. Sparsity
4. Variation
5. Expressivity
6. Unmodeled Variables
7. Unknown representations



Expressivity

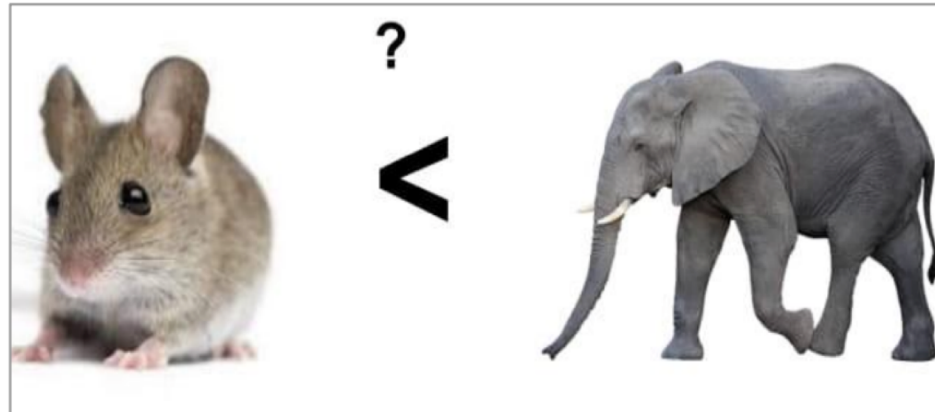
- Not only can one form have different meanings (ambiguity) but the same meaning can be expressed with different forms:
 - *She gave the book to Tom* vs. *She gave Tom the book*
 - *Some kids popped by* vs. *A few children visited*
 - *Is that window still open?* vs. *Please close the window*



Unmodeled Variables



“Drink this milk”



World knowledge

I dropped the glass on the floor and it broke

I dropped the hammer on the glass and it broke