

# Deep Learning for Vision & Language

Natural Language Processing I: Introduction



# Second Assignment

- To be released over the next day or two. Stay tuned. Deadline two weeks from release date. Tentative release date: Wednesday.
- Also updated syllabus (schedule) on our website to account for actual pace of the class.

# Today

- What is Natural Language Processing?
- Why is it hard?
- Common Tasks in NLP
- Language Modeling
- Word and Sentence representations for ML

# Natural Language Processing

The study of automatic reasoning over text / language

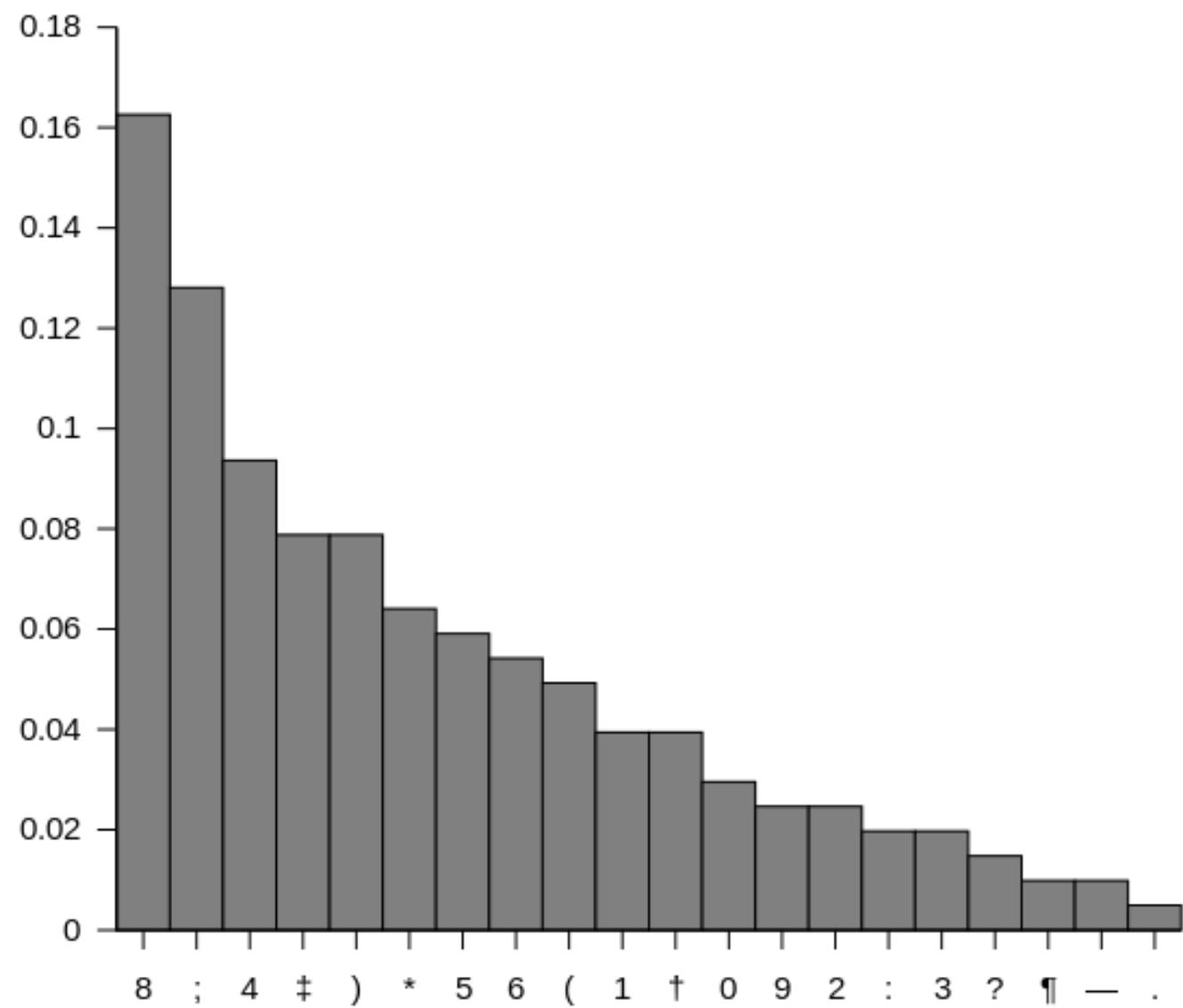


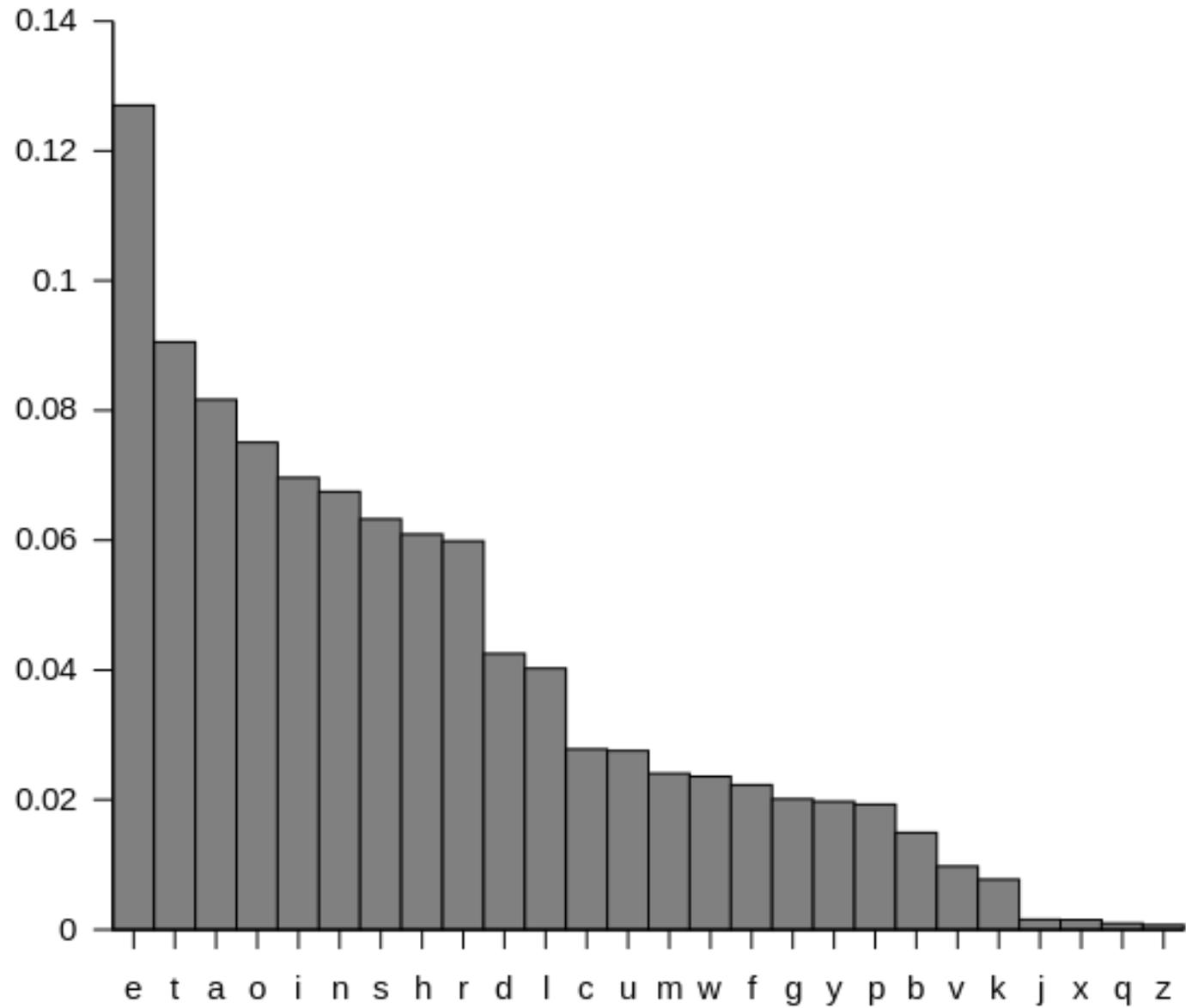
- **Fundamental goal: *deep* understand of *broad* language**
  - Not just string processing or keyword matching!
- **End systems that we want to build:**
  - Ambitious: speech recognition, machine translation, information extraction, dialog interfaces, question answering...
  - Modest: spelling correction, text categorization...

# Challenges in Natural Language Understanding:

53†††305))6\*;4826)4†.)4†);806\*;48†8  
¶60))85;;]8\*,:†\*8†83(88)5\*†;46(;88\*96  
\*?;8)\*†(;485);5\*†2:\*†(;4956\*2(5\*—4)8  
¶8\*;4069285);)6†8)4††;1(†9;48081;8:8†  
1;48†85;4)485†528806\*81(†9;48;(88;4  
(†?34;48)4†;161;:188;†?;

Any idea about what does it mean the text above?





# Challenges in Natural Language Understanding:

53††+305))6\*;4826)4†.)4†);806\*;48†8  
agoodglassinthebishopshostelinthede

¶60))85;;]8\*;;†\*8†83(88)5\*†;46(;88\*96  
vilsseattwentyonedegreesandthirteenmi

\*?;8)\*†(;485);5\*†2:\*†(;4956\*2(5\*—4)8  
nutesnortheastandbynorthmainbranchse

¶8\*;4069285);)6†8)4††;1(†9;48081;8:8†  
venthlimbeastsideshootfromthellefteyeo

1;48†85;4)485†528806\*81(†9;48;(88;4  
fthedeathsheadabeelinefromthetreeth

(†?34;48)4†;161;:188;†?;  
roughtheshotfiftyfeetout

# Challenges in Natural Language Understanding:

A good glass in the bishop's hostel in the  
devil's seat  
twenty-one degrees and  
thirteen minutes northeast and by north  
main branch seventh limb east side  
shoot from the left eye of the death's-head  
a bee line from the tree through  
the shot fifty feet out.

# Why is NLP Hard?

- Human Language is Ambiguous

## Task: Pronoun Resolution

- Jack drank the wine on the table. **It** was red and round.
- Jack saw Sam at the party. **He** went back to the bar to get another drink.
- Jack saw Sam at the party. **He** clearly had drunk too much.

[Adapted from Wilks (1975)]

# Why is NLP Hard?

- Human Language Requires World Knowledge

## Task: Co-Reference Resolution

- The doctor hired a secretary because she needed help with new patients.
- The physician hired the secretary because he was highly recommended.

[From some of our group's work]

[Gender Bias in Coreference Resolution: Evaluation and Debiasing Methods](#)

Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, Kai-Wei Chang.

North American Chapter of the Association for Computational Linguistics. NAACL 2018.

# Why is NLP Hard?

- Human Language is Ambiguous

## Learning mother tongue (native language)

-- you might think it's easy, but...

- compare 5 year old V.S. 10 year old V.S. 20 year old
- Learning foreign languages
  - even harder

# Is NLP really that hard?

- In the back of your mind, if you're still thinking...
  - *“My native language is so easy. How hard can it be to type all the grammar rules, and idioms, etc into a software program? Sure it might take a while, but with enough people and money, it should be doable!”*
- You are not alone!

# Brief History of NLP

- Mid 1950's – mid 1960's: Birth of NLP and Linguistics
  - At first, people thought NLP is easy! Researchers predicted that “machine translation” can be solved in 3 years or so.
  - Mostly hand-coded rules / linguistics-oriented approaches
  - The 3 year project continued for 10 years, but still no good result, despite the significant amount of expenditure.
- Mid 1960's – Mid 1970's: A Dark Era
  - After the initial hype, a dark era follows -- people started believing that machine translation is impossible, and most abandoned research for NLP.

# Brief History of NLP

- 1970's and early 1980's – Slow Revival of NLP
  - Some research activities revived, but the emphasis is still on linguistically oriented, working on small toy problems with weak empirical evaluation
- Late 1980's and 1990's – Statistical Revolution!
  - By this time, the computing power increased substantially .
  - Data-driven, statistical approaches with simple representation win over complex hand-coded linguistic rules.  
→ *“Whenever I fire a linguist our machine translation performance improves.” (Jelinek, 1988)*
- 2000's – Statistics Powered by Linguistic Insights
  - With more sophistication with the statistical models, richer linguistic representation starts finding a new value.
- 2010's – Neural Networks – Word Embeddings – Neural Language Modeling.
- 2018's – 2020's --- Transformers – Large Scale Language Pretraining
- 2030s' -- ??

# Ambiguity is Explosive

- Ambiguities compound to generate enormous numbers of possible interpretations.
- In English, a sentence ending in  $n$  prepositional phrases has *over*  $2^n$  syntactic interpretations.
  - “I saw the man with the telescope”: **2 parses**
  - “I saw the man on the hill with the telescope.”: **5 parses**
  - “I saw the man on the hill in Texas with the telescope”:  
**14 parses**
  - “I saw the man on the hill in Texas with the telescope at noon.”: **42 parses**
  - “I saw the man on the hill in Texas with the telescope at noon on Monday” **132 parses**

# Humor and Ambiguity

- Many jokes rely on the ambiguity of language:
  - Groucho Marx: One morning I shot an elephant in my pajamas. How he got into my pajamas, I'll never know.
  - She criticized my apartment, so I knocked her flat.
  - Noah took all of the animals on the ark in pairs. Except the worms, they came in apples.
  - Policeman to little boy: "We are looking for a thief with a bicycle." Little boy: "Wouldn't you be better using your eyes."
  - Why is the teacher wearing sun-glasses. Because the class is so bright.

# Why is Language Ambiguous?

# Why is Language Ambiguous?

- Having a unique linguistic expression for every possible conceptualization that could be conveyed would make language overly complex and linguistic expressions unnecessarily long.
- Allowing resolvable ambiguity permits shorter linguistic expressions, i.e. data compression.
- Language relies on people's ability to use their knowledge and inference abilities to properly resolve ambiguities.
- Infrequently, disambiguation fails, i.e. the compression is lossy.

# Natural Languages vs. Computer Languages

- Ambiguity is the primary difference between natural and computer languages.
- Formal programming languages are designed to be unambiguous, i.e. they can be defined by a grammar that produces a unique parse for each sentence in the language.
- Programming languages are also designed for efficient (deterministic) parsing.

# Natural Language Tasks

- Processing natural language text involves many various syntactic, semantic and pragmatic tasks in addition to other problems.

# Syntactic Tasks

# Word Segmentation

- Breaking a string of characters into a sequence of words.
- In some written languages (e.g. Chinese) words are not separated by spaces.
- Even in English, characters other than white-space can be used to separate words [e.g. , ; . - : ( ) ]
- Examples from English URLs:
  - jumptheshark.com ⇒ jump the shark .com
  - myspace.com/pluckerswingbar
    - ⇒ myspace .com pluckers wing bar
    - ⇒ myspace .com plucker swing bar

# Morphological Analysis

- **Morphology** is the field of linguistics that studies the internal structure of words. (Wikipedia)
- A **morpheme** is the smallest linguistic unit that has semantic meaning (Wikipedia)
  - e.g. “carry”, “pre”, “ed”, “ly”, “s”
- Morphological analysis is the task of segmenting a word into its morphemes:
  - carried  $\Rightarrow$  carry + ed (past tense)
  - independently  $\Rightarrow$  in + (depend + ent) + ly
  - Googlers  $\Rightarrow$  (Google + er) + s (plural)
  - unlockable  $\Rightarrow$  un + (lock + able) ?  
 $\Rightarrow$  (un + lock) + able ?

- ***German***

555 --> fünfhundertfünfundfünfzig

7254 → Siebentausendzweihundertvierundfünfzig

# Part Of Speech (POS) Tagging

- Annotate each word in a sentence with a part-of-speech.

I ate the spaghetti with meatballs.

John saw the saw and decided to take it to the table.

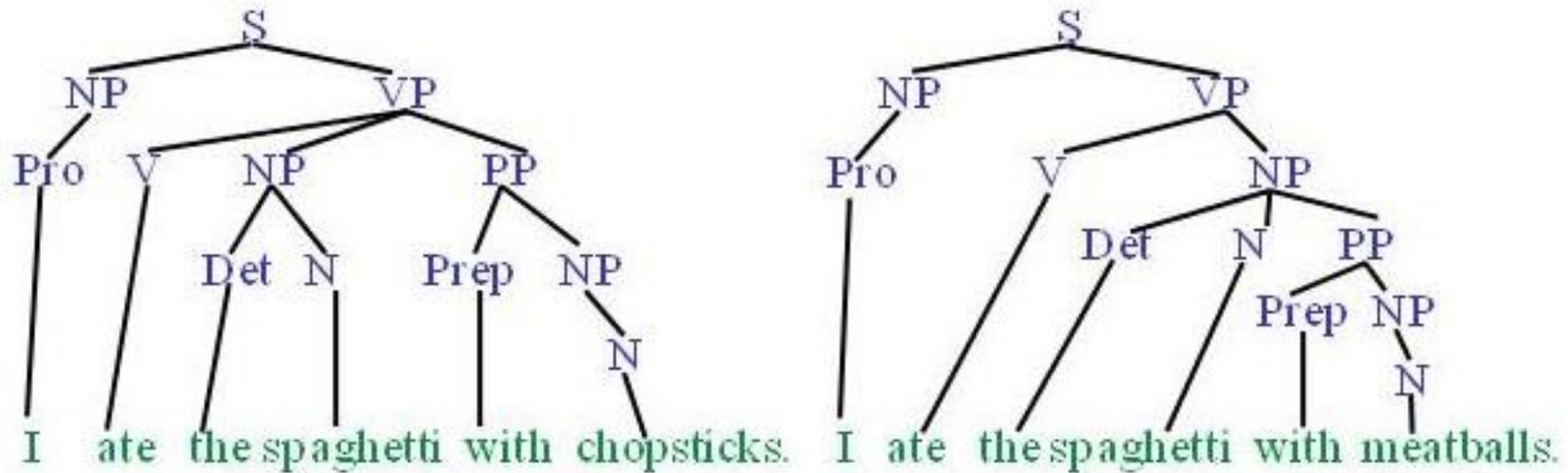
- Useful for subsequent syntactic parsing and word sense disambiguation.

# Phrase Chunking

- Find all noun phrases (NPs) and verb phrases (VPs) in a sentence.
  - [NP I] [VP ate] [NP the spaghetti] [PP with] [NP meatballs].
  - [NP He ] [VP reckons ] [NP the current account deficit ] [VP will narrow ] [PP to ] [NP only # 1.8 billion ] [PP in ] [NP September ]

# Syntactic Parsing

- Produce the correct syntactic parse tree for a sentence.



# Semantic Tasks

# Word Sense Disambiguation (WSD)

- Words in natural language usually have a fair number of different possible meanings.
  - Ellen has a strong **interest** in computational linguistics.
  - Ellen pays a large amount of **interest** on her credit card.
- For many tasks (question answering, translation), the proper sense of each ambiguous word in a sentence must be determined.

# Semantic Role Labeling (SRL)

- For each clause, determine the semantic role played by each noun phrase that is an argument to the verb.
  - John drove Mary from Austin to Dallas in his Toyota Prius.
  - The hammer broke the window.
- Also referred to a “case role analysis,” “thematic analysis,” and “shallow semantic parsing”

# Textual Entailment

- Determine whether one natural language sentence entails (implies) another under an ordinary interpretation.

# Textual Entailment Problems from PASCAL Challenge

## TEXT

*Eyeing the huge market potential, currently led by Google, Yahoo took over search company Overture Services Inc last year.*

*Microsoft's rival Sun Microsystems Inc. bought Star Office last month and plans to boost its development as a Web-based device running over the Net on personal computers and Internet appliances.*

*The National Institute for Psychobiology in Israel was established in May 1971 as the Israel Center for Psychobiology by Prof. Joel.*

*Since its formation in 1948, Israel fought many wars with neighboring Arab countries.*

## HYPOTHESIS

*Yahoo bought Overture.*

*Microsoft bought Star Office.*

*Israel was established in May 1971.*

*Israel was established in 1948.*

## ENTAILMENT

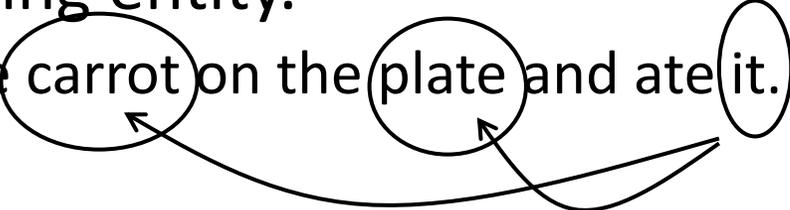
TRUE

# Pragmatics/Discourse Tasks

# Anaphora Resolution/ Co-Reference

- Determine which phrases in a document refer to the same underlying entity.

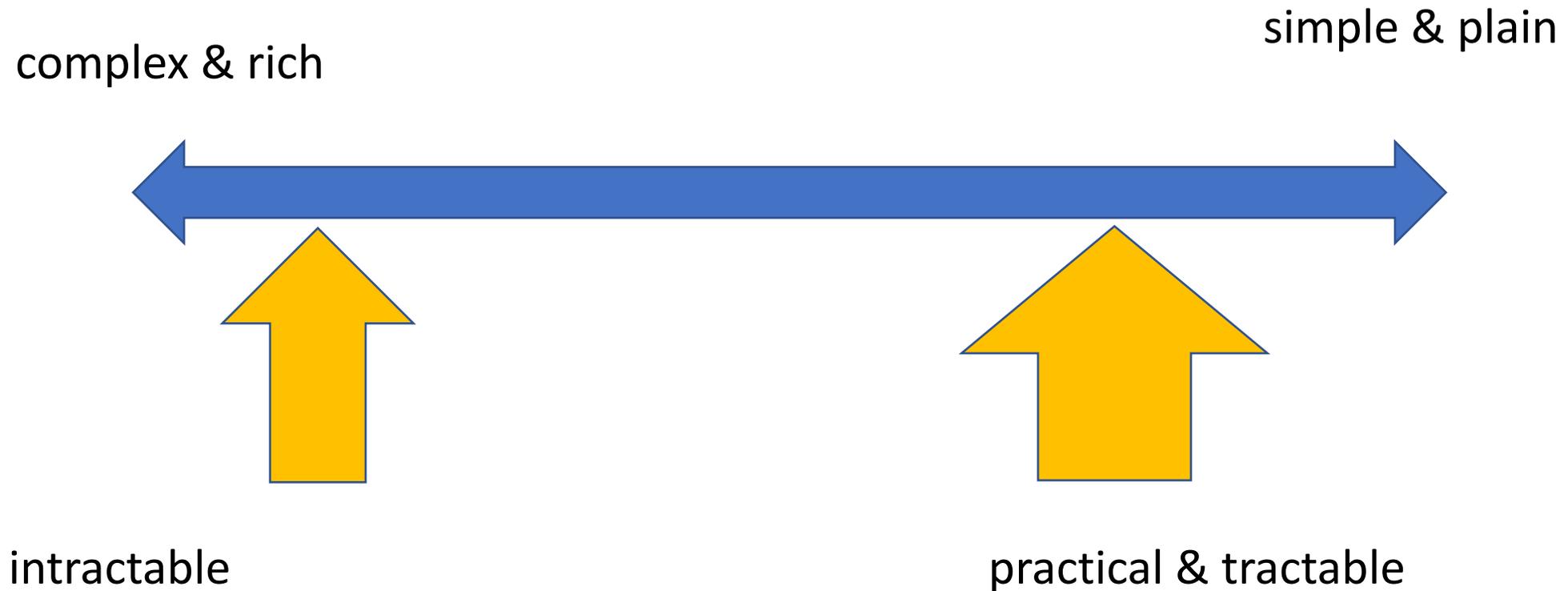
- John put the carrot on the plate and ate it.



- Bush started the war in Iraq. But the president needed the consent of Congress.



# Representation vs Computability



# How to represent a word?

one-hot encodings

dog	1	[1 0 0 0 0 0 0 0 0 0]
cat	2	[0 1 0 0 0 0 0 0 0 0]
person	3	[0 0 1 0 0 0 0 0 0 0]
holding	4	[0 0 0 1 0 0 0 0 0 0]
tree	5	[0 0 0 0 1 0 0 0 0 0]
computer	6	[0 0 0 0 0 1 0 0 0 0]
using	7	[0 0 0 0 0 0 1 0 0 0]

# How to represent a phrase/sentence?

bag-of-words representation

person holding dog	{1, 3, 4}	[1	0	1	1	0	0	0	0	0	0]
person holding cat	{2, 3, 4}	[0	1	1	1	0	0	0	0	0	0]
person using computer	{3, 7, 6}	[0	0	1	0	0	1	1	0	0	0]
		dog	cat	person	holding	tree	computer	using			
person using computer person holding cat	{3, 3, 7, 6, 2}	[0	1	2	1	0	1	1	0	0	0]

What if vocabulary is very large?

# Sparse Representation

bag-of-words representation

person holding dog	{1, 3, 4}	indices = [1, 3, 4]	values = [1, 1, 1]
person holding cat	{2, 3, 4}	indices = [2, 3, 4]	values = [1, 1, 1]
person using computer	{3, 7, 6}	indices = [3, 7, 6]	values = [1, 1, 1]
person using computer person holding cat	{3, 3, 7, 6, 2}	indices = [3, 7, 6, 2]	values = [2, 1, 1, 1]

# Recap

- Bag-of-words encodings for text (e.g. sentences, paragraphs, captions, etc)

You can take a set of sentences/documents and classify them, cluster them, or compute distances between them using this representation.

# Problem with this bag-of-words representation

my friend makes a nice meal

These would be the same using bag-of-words

my nice friend makes a meal

# Bag of Bi-grams

indices = [10132, 21342, 43233, 53123, 64233]

values = [1, 1, 1, 1, 1]

my friend makes a nice meal

{my friend, friend makes, makes a,  
a nice, nice meal}

indices = [10232, 43133, 21342, 43233, 54233]

values = [1, 1, 1, 1, 1]

my nice friend makes a meal

{my nice, nice friend, friend makes,  
makes a, a meal}

A dense vector-representation would be very inefficient

Think about tri-grams and n-grams

# Recommended reading: n-gram language models

Kai-Wei's course on Natural Language Processing

<http://www.cs.virginia.edu/~kc2wc/teaching/NLP16/slides/02-ngram.pdf>

<http://www.cs.virginia.edu/~kc2wc/teaching/NLP16/slides/03-smooth.pdf>

Yejin Choi's course on Natural Language Processing

<http://www3.cs.stonybrook.edu/~ychoi/cse628/lecture/02-ngram.pdf>

# Problem with this bag-of-words representation

my friend makes a nice meal

chicken makes a nice meal

Alternatives:

Continuous Bag of Words (CBOW) – Word embeddings

Sequence-based representations (RNNs, LSTMs)

Transformer-based representations (e.g. BERT)

# Back to how to represent a word?

Problem: distance between words using one-hot encodings always the same

dog	1	[1 0 0 0 0 0 0 0 0 0]
cat	2	[0 1 0 0 0 0 0 0 0 0]
person	3	[0 0 1 0 0 0 0 0 0 0]

Idea: Instead of one-hot-encoding use a histogram of commonly co-occurring words.

# Distributional Semantics



Dogs are man's best friend.

I saw a dog on a leash walking in the park.

His dog is his best companion.

He walks his dog in the late afternoon

...

	friend	leash	park	walking	walks	food	legs	runs	sleeps	sits	...
dog	[3	2	3	4	2	4	3	5	6	7	...]

# Distributional Semantics

dog	[5	5	0	5	0	0	5	5	0	2	...]
cat	[5	4	1	4	2	0	3	4	0	3	...]
person	[5	5	1	5	0	2	5	5	0	0	...]
	food	walks	window	runs	mouse	invented	legs	sleeps	mirror	tail	...



This vocabulary can be extremely large

# Toward more Compact Representations

dog	[5	5	0	5	0	0	5	5	0	2	...
cat	[5	4	1	4	2	0	3	4	0	3	...
person	[5	5	1	5	0	2	5	5	0	0	...

food   walks   window   runs   mouse   invented   legs   sleeps   mirror   tail   ...



This vocabulary can be extremely large

# Toward more Compact Representations

$$\text{dog} = \begin{bmatrix} 5 \\ 5 \\ 0 \\ 5 \\ 0 \\ 0 \\ 5 \\ 5 \\ 0 \\ 2 \\ \dots \end{bmatrix} = w_1 \begin{bmatrix} 0 \\ 1 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ \dots \end{bmatrix} + w_2 \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ \dots \end{bmatrix} + w_3 \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ \dots \end{bmatrix} + \dots$$

legs, running,  
walking

tail, fur,  
ears

mirror, window,  
door

# Toward more Compact Representations

dog =  $[ w_1 \quad w_2 \quad w_3 ]$

The basis vectors can be found using Principal Component Analysis (PCA)

This is known as Latent Semantic Analysis in NLP

# Toward more Compact Representations: Word Embeddings

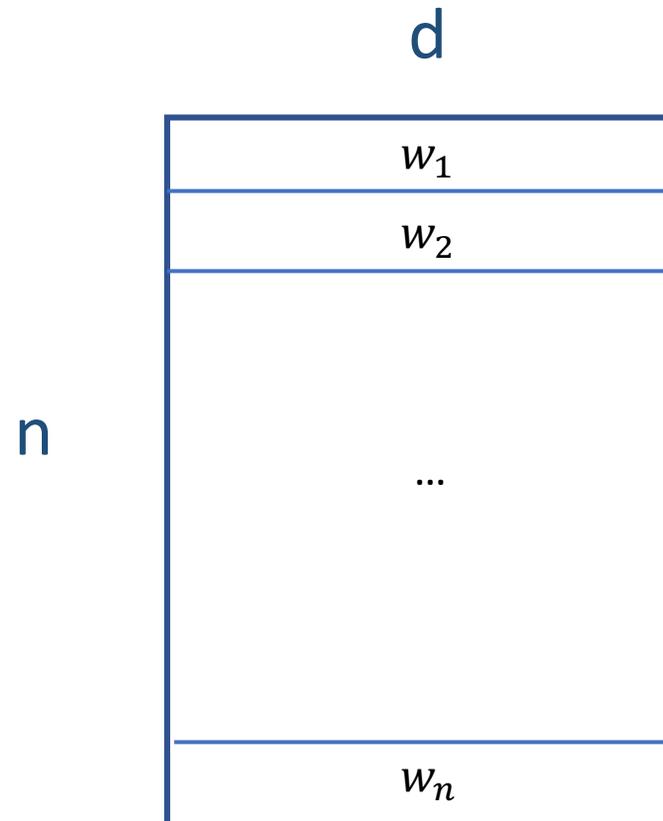
dog =  $[ w_1 \quad w_2 \quad w_3 ]$

The weights  $w_1, \dots, w_n$  are found using a neural network

Word2Vec: <https://arxiv.org/abs/1301.3781>

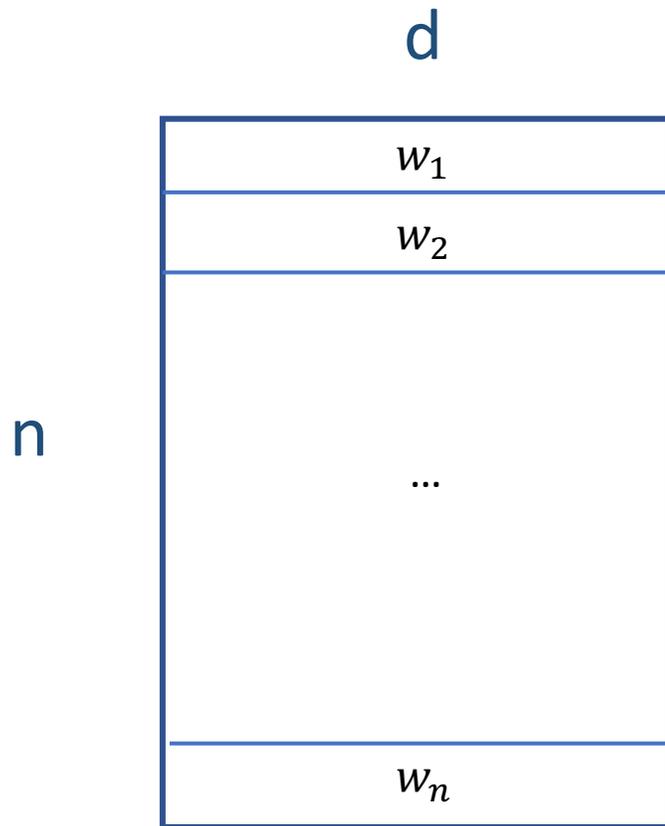
# Word2Vec – CBOW Version

- First, create a huge matrix of word embeddings initialized with random values – where each row is a vector for a different word in the vocabulary.

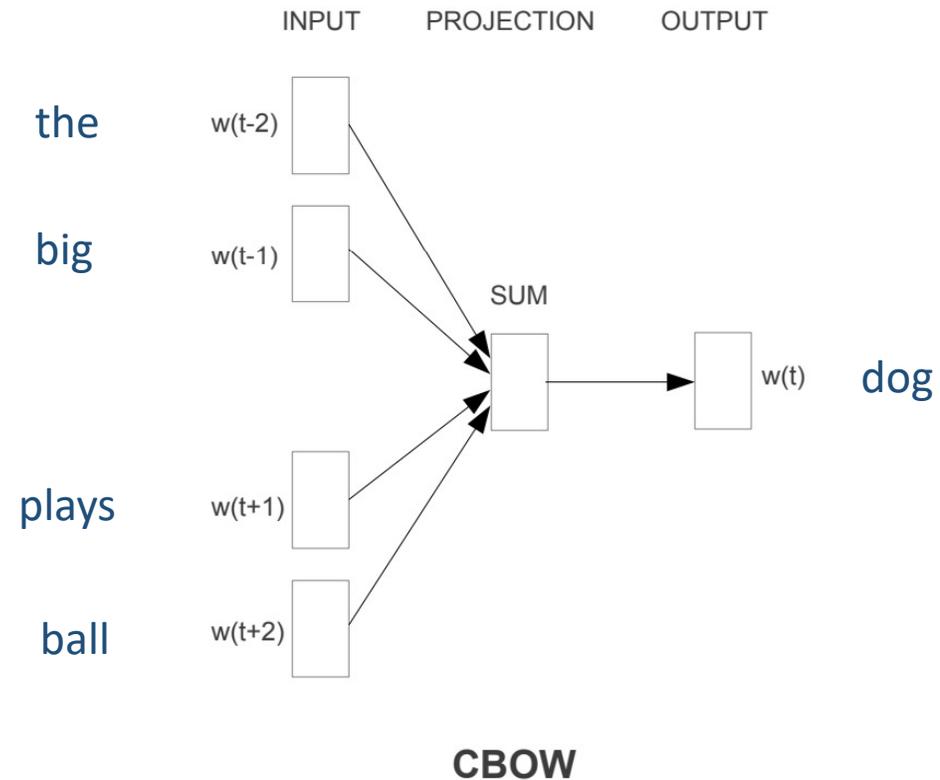


# Word2Vec – CBOW Version

- Then, collect a lot of text, and solve the following regression problem for a large corpus of text:



“the big dog plays ball”



# Practical Issues - Tokenization

- For each text representation we usually need to separate a sentence into tokens – we have assumed words in this lecture (or pairs of words) – but tokens could also be characters and anything in-between.
- Word segmentation can be used as tokenization.
  - In the assignment I was lazy I just did “my sentence”.split(“ ”) and called it a day.
  - However, even English is more difficult than that because of punctuation, double spaces, quotes, etc. For English I would recommend you too look up the great word tokenization tools in libraries such as Python’s NLTK and Spacy before you try to come up with your own word tokenizer.

# Issues with Word based Tokenization

- We already mentioned that tokenization can be hard even when word-based for other languages that don't use spaces in-between words.
- Word tokenization can also be bad for languages where the words can be “glued” together like German or Turkish.
  - Remember fünfhundertfünfundfünfzig? It wouldn't be feasible to have a word embedding for every number in the German language.
- It is problematic to handle words that are not in the vocabulary e.g. a common practice is to use a special <OOV> (out of vocabulary) token for those words that don't show up in the vocabulary.

# Solution: Sub-word Tokenization

- Byte-pair Encoding Tokenization (BPE)
  - Start from small strings and based on substring counts iteratively use larger sequences until you define a vocabulary that maximizes informative subtokens. That way most will correspond to words at the end.
- Byte-level BPE Tokenizer
  - Do the same but at the byte representation level not at the substring representation level.

We will discuss these more as we discuss Transformer Models

## Tokenizers

 Rust  license  downloads/week 

Provides an implementation of today's most used tokenizers, with a focus on performance and versatility.

### Main features:

- Train new vocabularies and tokenize, using today's most used tokenizers.
- Extremely fast (both training and tokenization), thanks to the Rust implementation. Takes less than 20 seconds to tokenize a GB of text on a server's CPU.
- Easy to use, but also extremely versatile.
- Designed for research and production.
- Normalization comes with alignments tracking. It's always possible to get the part of the original sentence that corresponds to a given token.
- Does all the pre-processing: Truncate, Pad, add the special tokens your model needs.

[huggingface/tokenizers](https://huggingface.co/tokenizers)

# Questions?