# Deep Learning for Vision & Language

Natural Language Processing I: Introduction

RICE UNIVERSITY

# Natural Language Processing
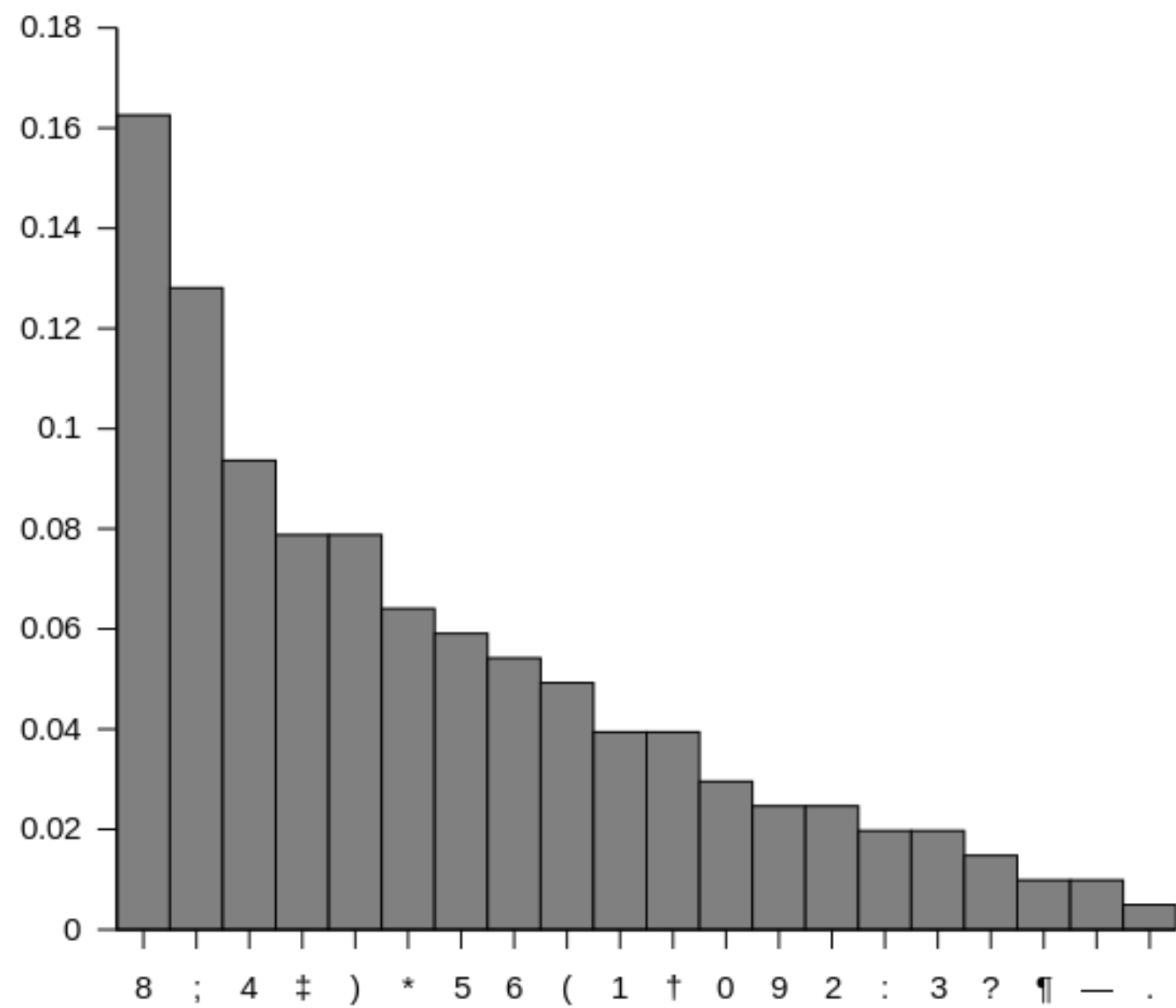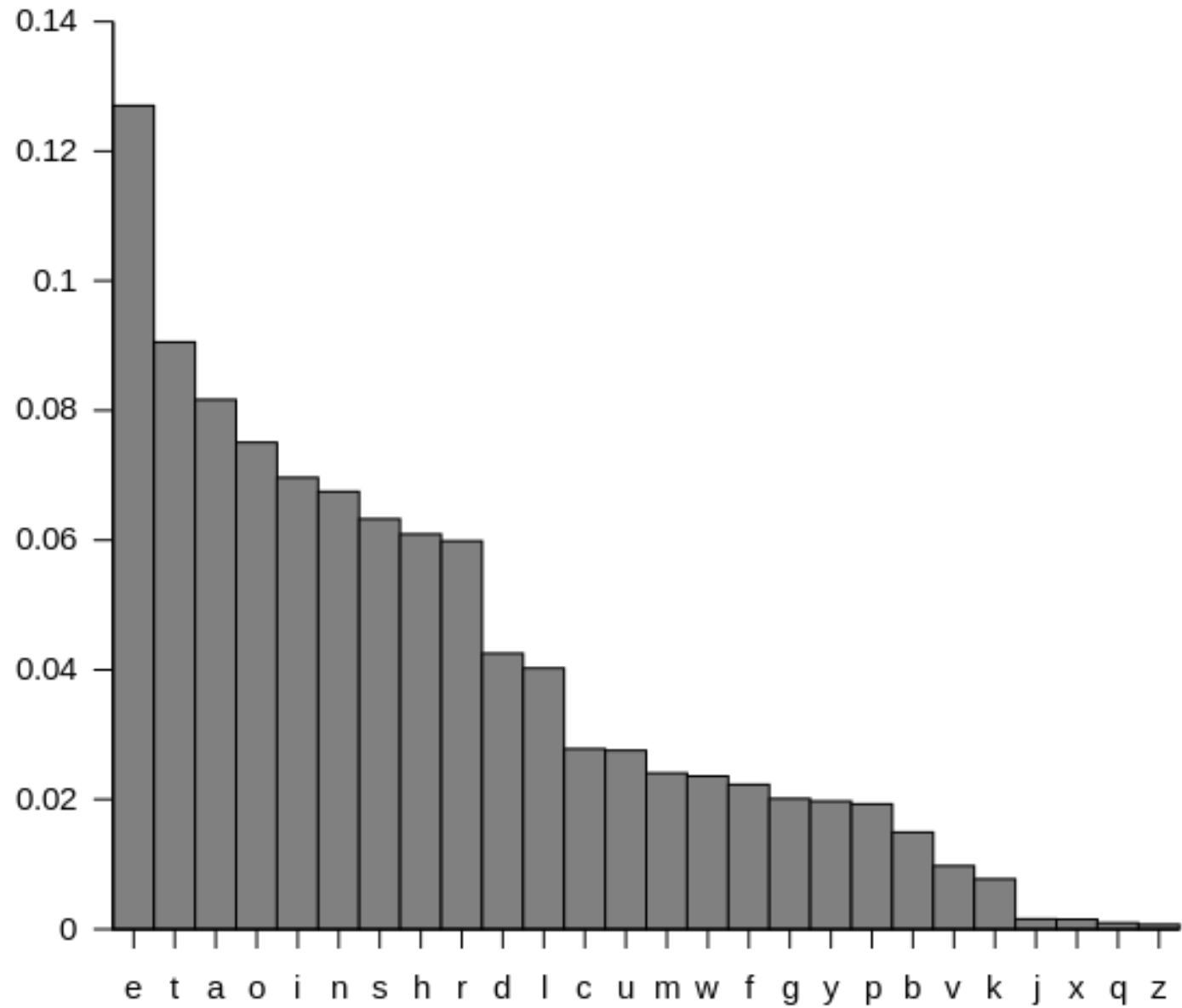
The study of automatic reasoning over text / language

- **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…

Slide by Dan Klein

# 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‡
venthlimbeastsideshootfromthelefteyeo

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.

The Gold-Bug, Edgar Allan Poe

# Why is NLP Hard?

- Human Language is Ambiguous

Task: <u>Pronoun Resolution</u>
  - 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: <u>Co-Reference Resolution</u>
  - 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**.

Adapted from a slide by Yejin Choi

# 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

Slide by Yejin Choi

# 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 $\Rightarrow$ jump the shark .com
  - myspace.com/pluckerswingbar

    $\Rightarrow$ myspace .com pluckers wing bar

    $\Rightarrow$ 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 $\Longrightarrow$ carry + ed (past tense)
  - independently $\Longrightarrow$ in + (depend + ent) + ly
  - Googlers $\Longrightarrow$ (Google + er) + s (plural)
  - unlockable $\Longrightarrow$ un + (lock + able) ?
    $\Longrightarrow$ (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.

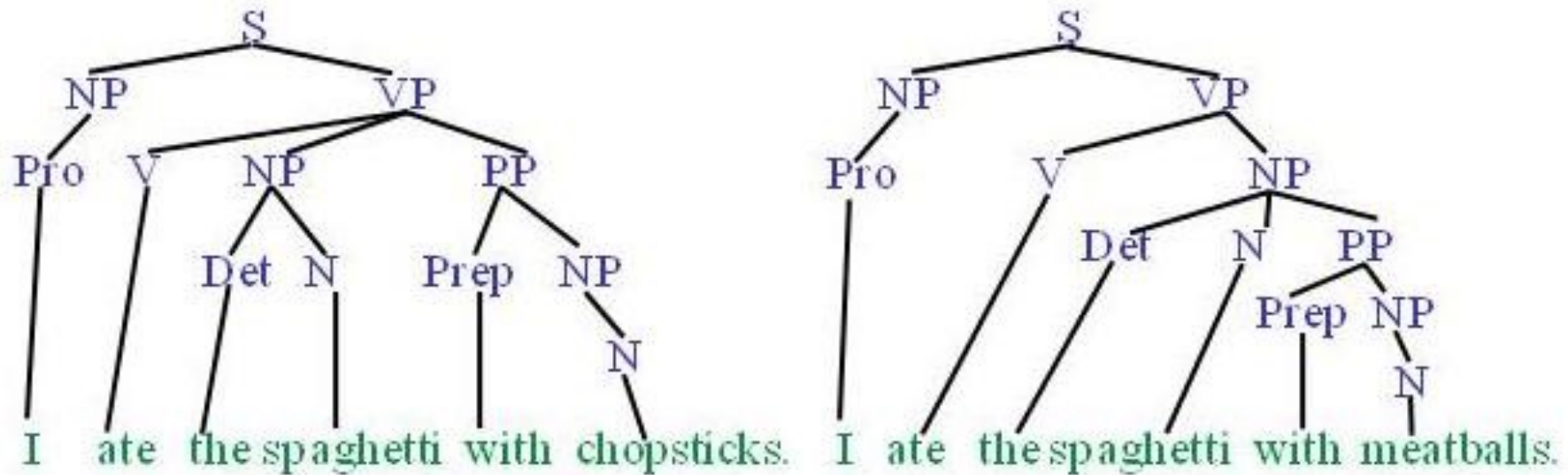- Useful for subsequent syntactic parsing and word sense disambiguation.

# 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 ]

Slide from Ray Mooney

# Syntactic Parsing

- Produce the correct syntactic parse tree for a sentence.

# 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.

Slide from Ray Mooney

# Textual Entailment

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

# Textual Entailment Problems from PASCAL Challenge

| TEXT | HYPOTHESIS | ENTAILMENT |
|------|-----------|-------------|
| *Eyeing the huge market potential, currently led by Google, Yahoo took over search company Overture Services Inc last year.* | *Yahoo bought Overture.* | |
| *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.* | *Microsoft bought Star Office.* | |
| *The National Institute for Psychobiology in Israel was established in May 1971 as the Israel Center for Psychobiology by Prof. Joel.* | *Israel was established in May 1971.* | |
| *Since its formation in 1948, Israel fought many wars with neighboring Arab countries.* | *Israel was established in 1948.* | TRUE |

# 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 word?

# 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}

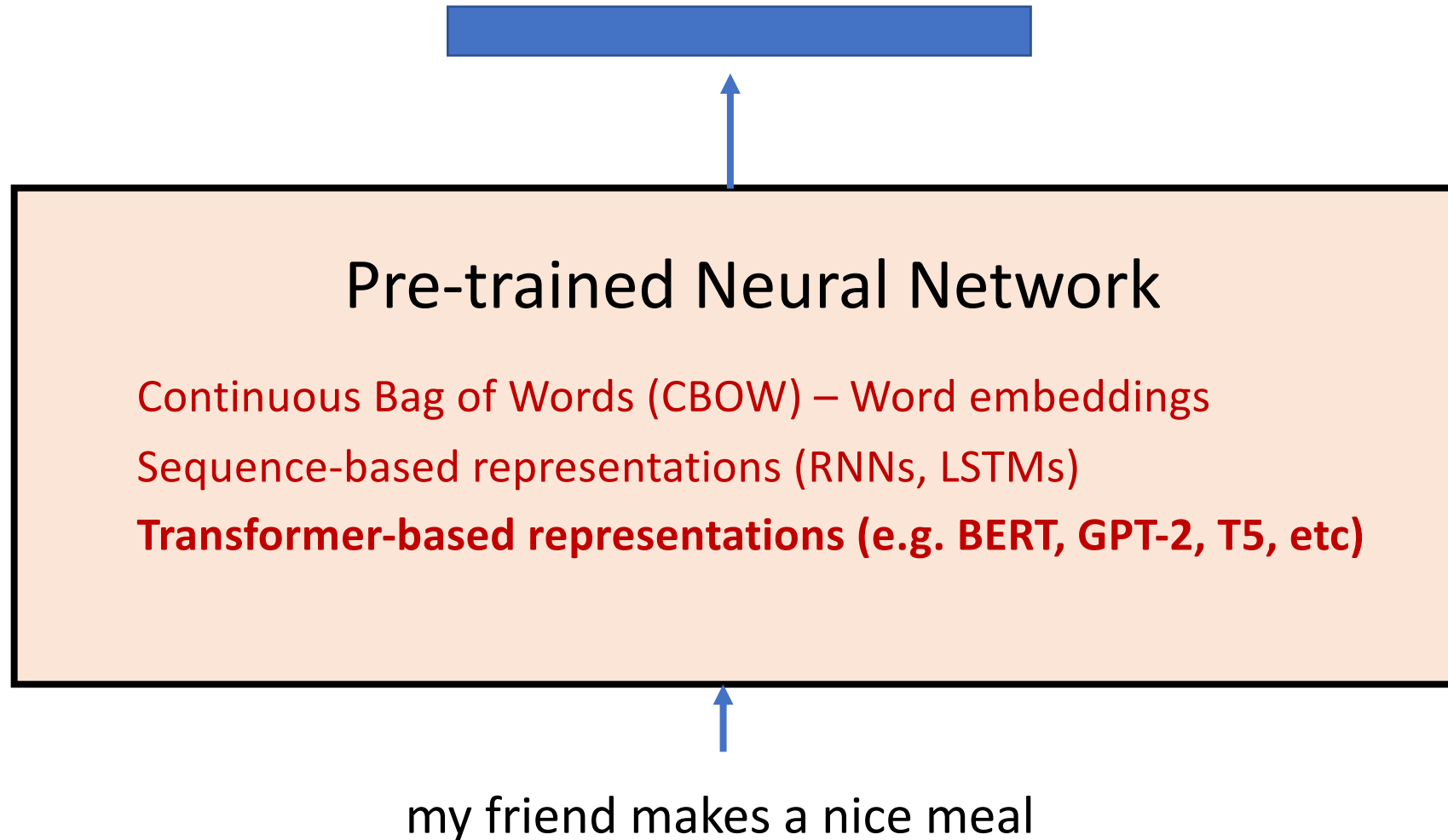A dense vector-representation would be very inefficient
Think about tri-grams and n-grams

# Recommended reading: n-gram language models

Yejin Choi's course on Natural Language Processing

http://www3.cs.stonybrook.edu/~ychoi/cse628/lecture/02-ngram.pdf
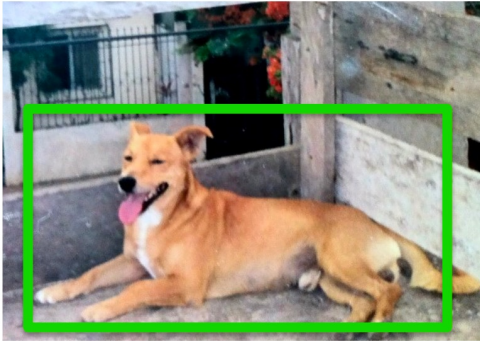
# Modern way of representing Phrases/Text

Pre-trained Neural Network

Continuous Bag of Words (CBOW) – Word embeddings

Sequence-based representations (RNNs, LSTMs)

**Transformer-based representations (e.g. BERT, GPT-2, T5, etc)**

my friend makes a nice meal

# 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

$$dog = \begin{bmatrix} 5 \\ 5 \\ 0 \\ 5 \\ 0 \\ 0 \\ 5 \\ 5 \\ 0 \\ 2 \\ ... \end{bmatrix} = w1 \begin{bmatrix} 0 \\ 1 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ ... \end{bmatrix} + w2 \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ ... \end{bmatrix} + w3 \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ ... \end{bmatrix} + ...$$

legs, running, walking          tail, fur, ears          mirror, window, door

# Toward more Compact Representations

dog = $\begin{bmatrix} w1 & w2 & w3 \end{bmatrix}$

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

This is known as Latent Semantic Analysis in NLP

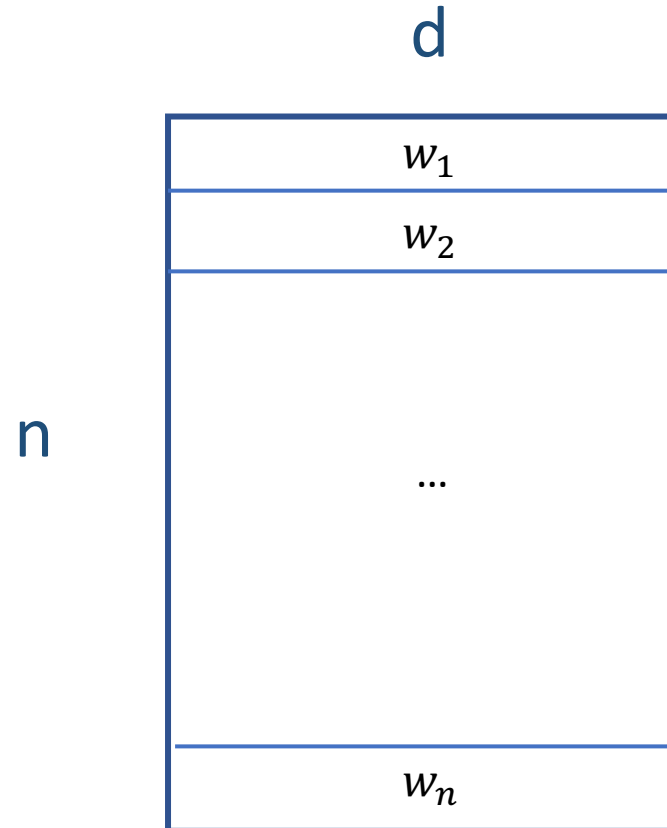# Toward more Compact Representations: Word Embeddings

dog =   $\begin{bmatrix} w1 & w2 & w3 \end{bmatrix}$

The weights w1, …, wn 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.

# Efficient Estimation of Word Representations in Vector Space

**Tomas Mikolov**

Google Inc., Mountain View, CA

tmikolov@google.com

**Greg Corrado**

Google Inc., Mountain View, CA

gcorrado@google.com

**Kai Chen**

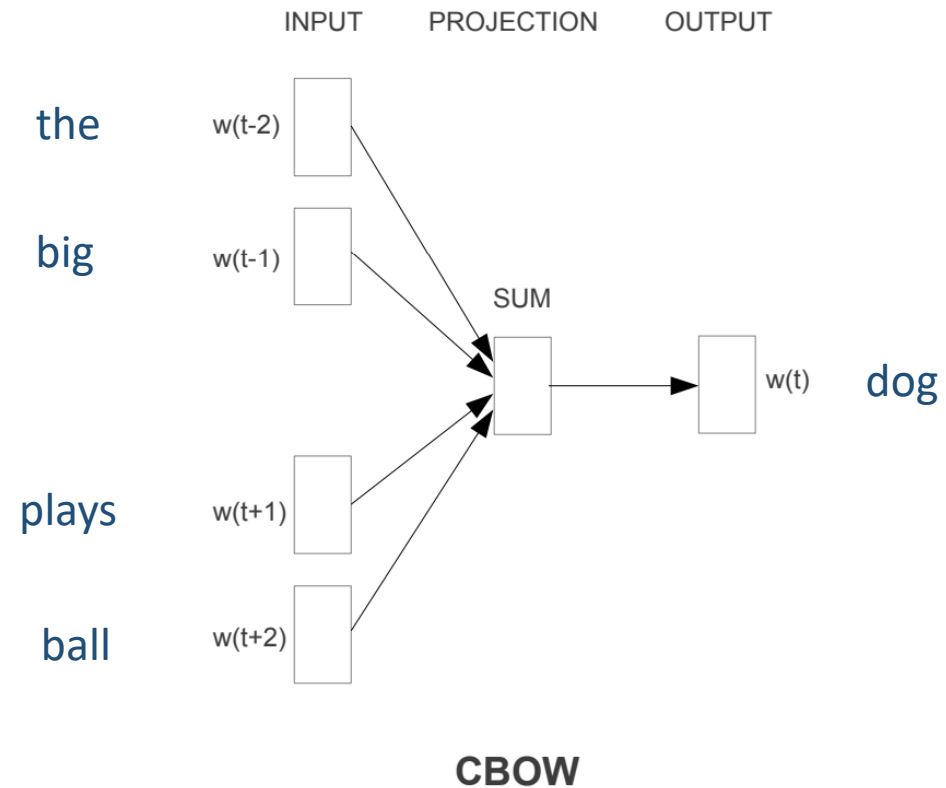Google Inc., Mountain View, CA

kaichen@google.com

**Jeffrey Dean**

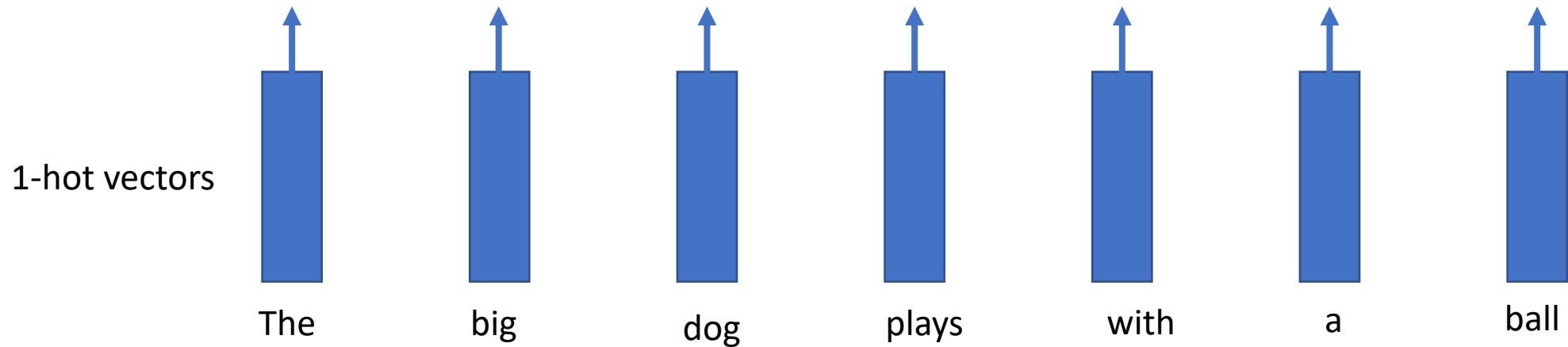Google Inc., Mountain View, CA

jeff@google.com

# Word2Vec – CBOW Version

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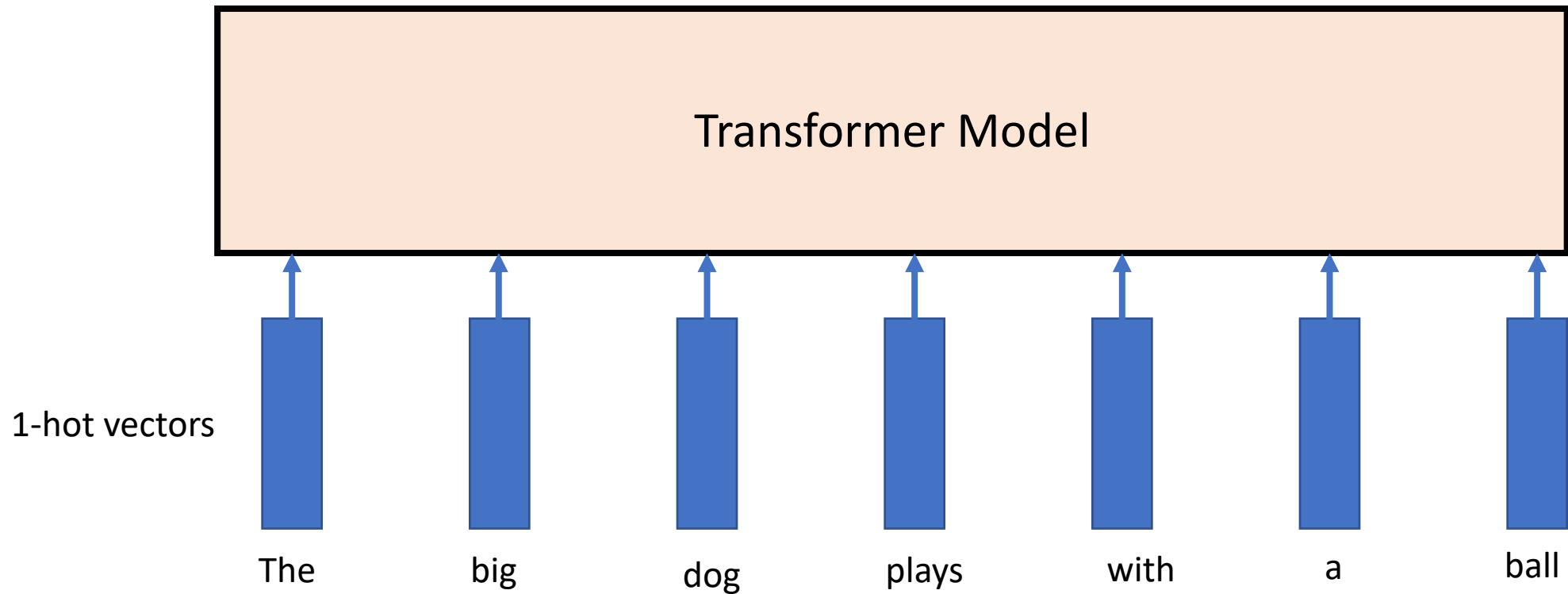


"the big dog plays ball"

d

n

$w_1$

$w_2$

…

$w_n$
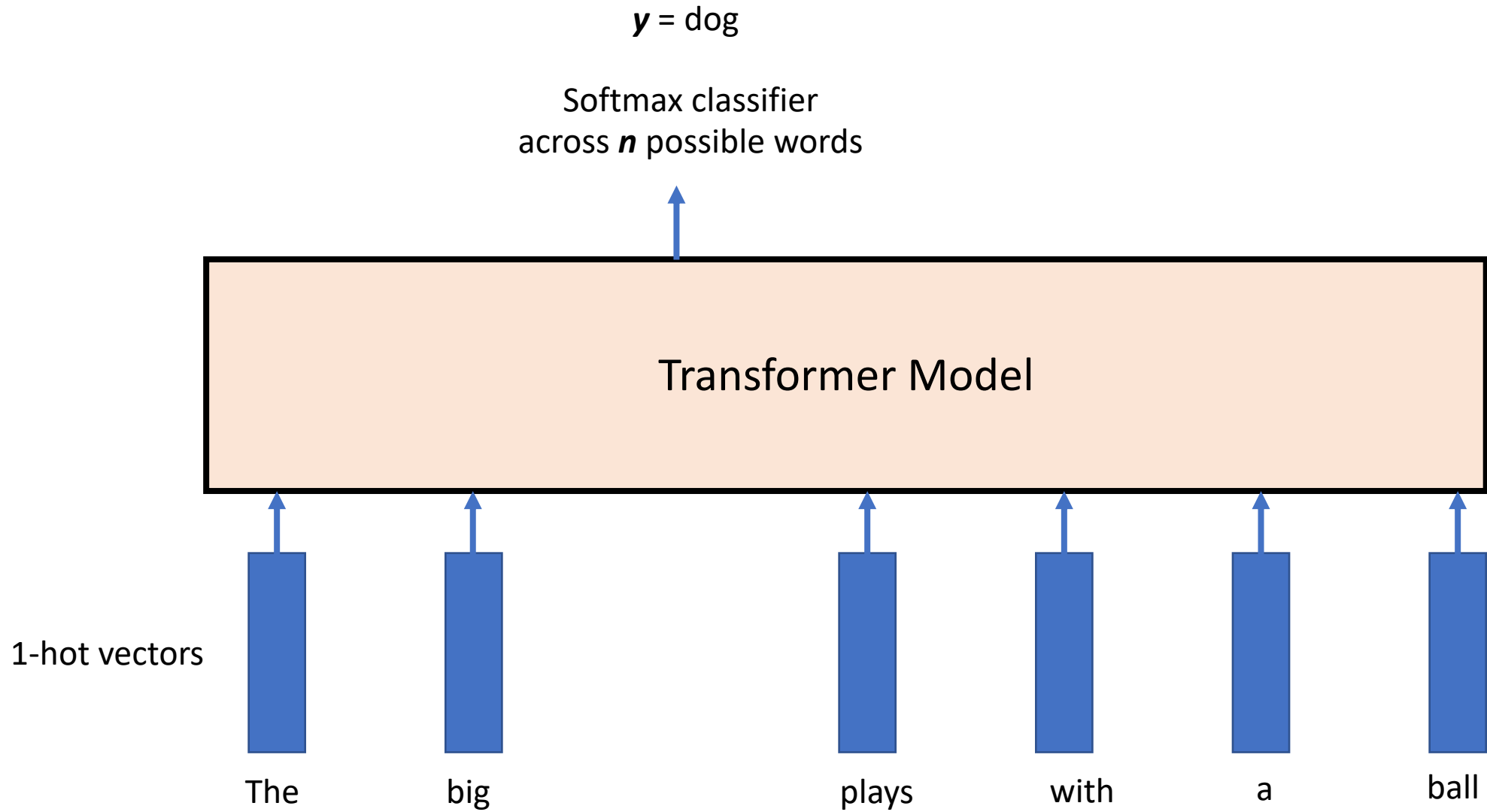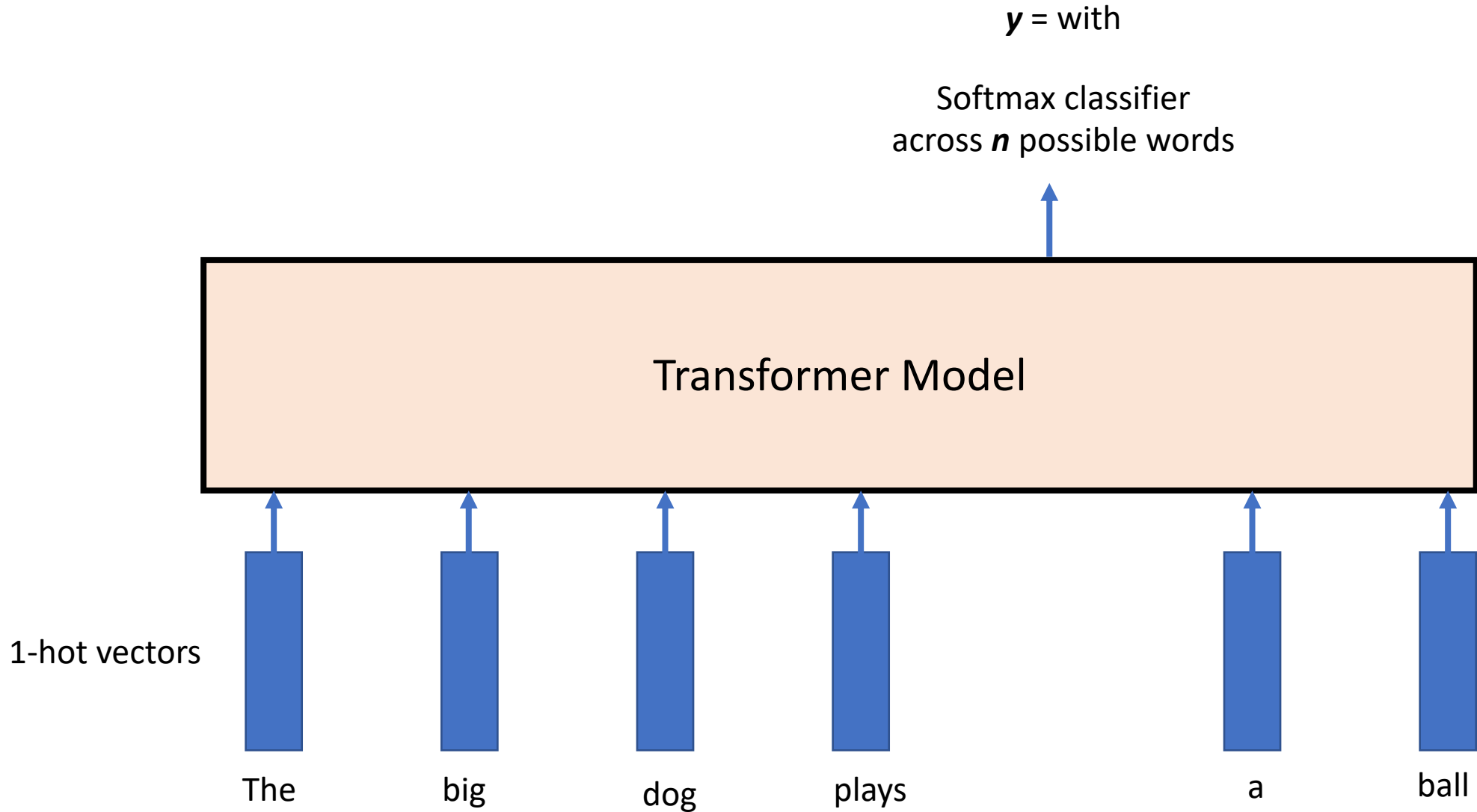
INPUT    PROJECTION    OUTPUT

the    w(t-2)

big    w(t-1)

SUM

plays    w(t+1)

ball    w(t+2)

w(t)    dog

CBOW

# Pre-trained Language Models

1-hot vectors



The    big    dog    plays    with    a    ball

# Pre-trained Language Models



Transformer Model

1-hot vectors

The    big    dog    plays    with    a    ball

# Pre-trained Language Models

$y$ = dog

Softmax classifier
across $n$ possible words

Transformer Model

1-hot vectors

The        big        plays        with        a        ball

# Pre-trained Language Models

$y$ = with

Softmax classifier
across $n$ possible words

Transformer Model

1-hot vectors

The      big      dog      plays                    a      ball
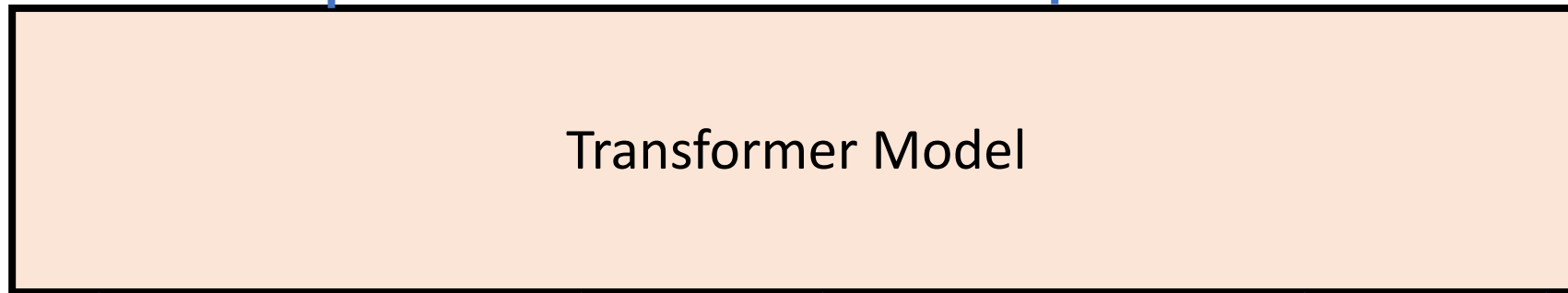
# Pre-trained Language Models

$y_1$ = big

$y_2$ = with

Softmax classifier
across *n* possible words

Softmax classifier
across *n* possible words

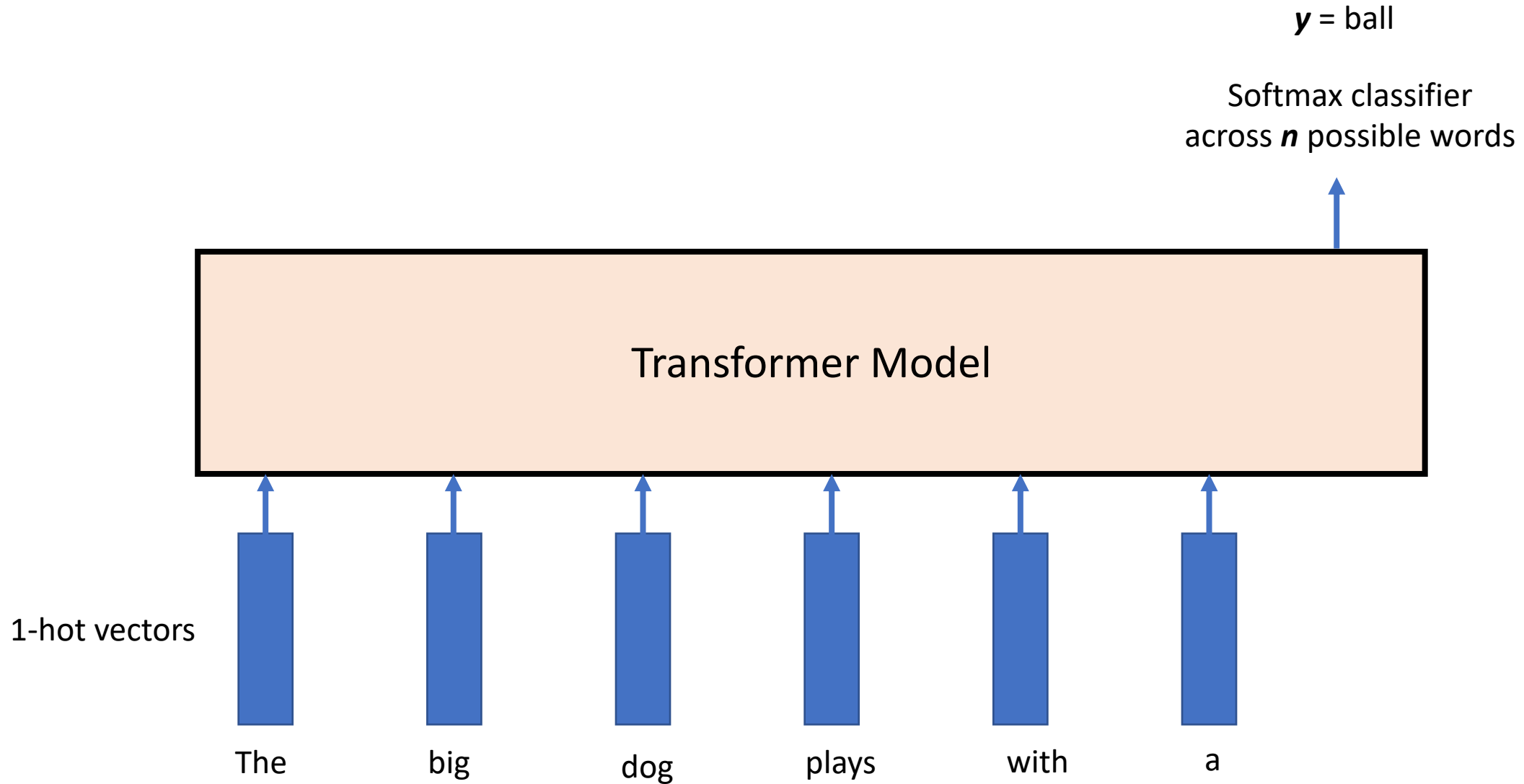Transformer Model

1-hot vectors

The          dog     plays                    a        ball

# Generative Language Models

$y$ = ball

Softmax classifier
across $n$ possible words

Transformer Model

1-hot vectors

The    big    dog    plays    with    a

# 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 passing | license Apache-2.0 | downloads/week 169k

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

# Questions?