

# Deep Dive (3)

*Text Vectors, Word Embeddings,  
Word2vec, Glove, Implement with  
Gensim, Use-cases and overview of BERT*



# Agenda

1. Recap Week 3
2. Text Vectorization & Challenges
3. Word Embeddings
  - Word2vec
  - Stanford GloVe
4. Use-cases
  - Doc Similarity, Question Answering, Information Retrieval
5. State of NLP today - intro to BERT
6. Theory Wrap-up
7. Implement with Google Colab

# Recap - Week 3

---

TF-IDF

Bag-of-Words

Document Similarity - Use Cases

# Recap TF-IDF

$$w_{x,y} = \text{tf}_{x,y} \times \log \left( \frac{N}{\text{df}_x} \right)$$

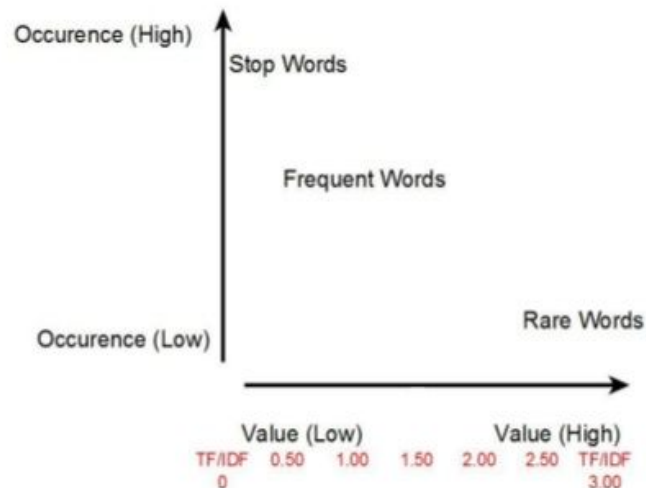
## TF-IDF

Term  $x$  within document  $y$

$\text{tf}_{x,y}$  = frequency of  $x$  in  $y$

$\text{df}_x$  = number of documents containing  $x$

$N$  = total number of documents



# Recap

## Bag-of-Words

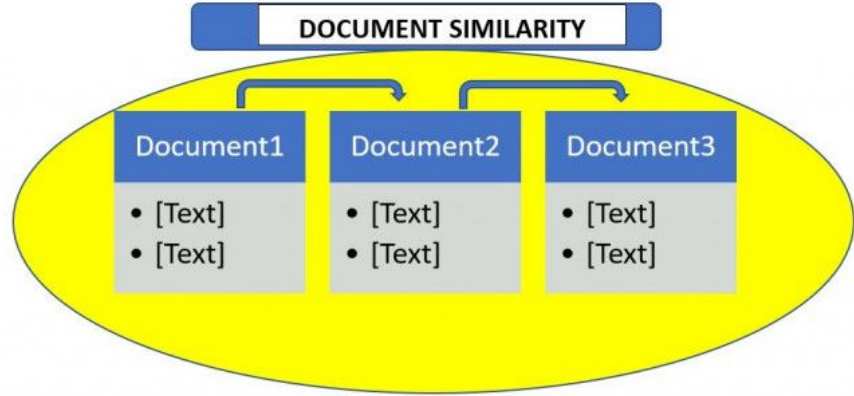
Raw Text

Bag-of-words  
vector

it is a puppy and it  
is extremely cute

it	2
they	0
puppy	1
and	1
cat	0
aardvark	0
cute	1
extremely	1
...	...

# Recap Document Similarity



# 1

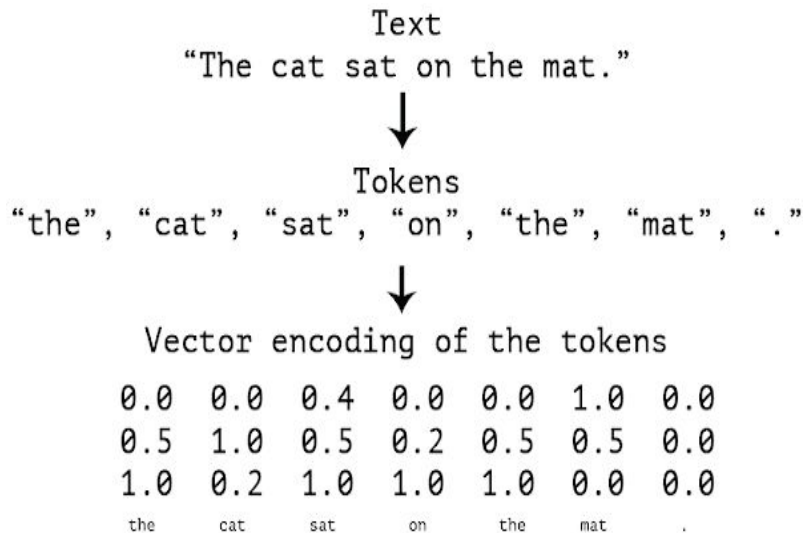
## What is Word Embedding?

We need embeddings!



# Word Vectorization Review

- Vectorization is the process of converting text into numerical representation
- Few techniques available, each one with its own pros and cons
  - simplest encoding techniques do not retain word order
  - fast and intuitive, but the size of document vectors grows quickly with the size of the dictionary
  - optimize dimension but lose in interpretability
- Not only words - sentences, documents etc.

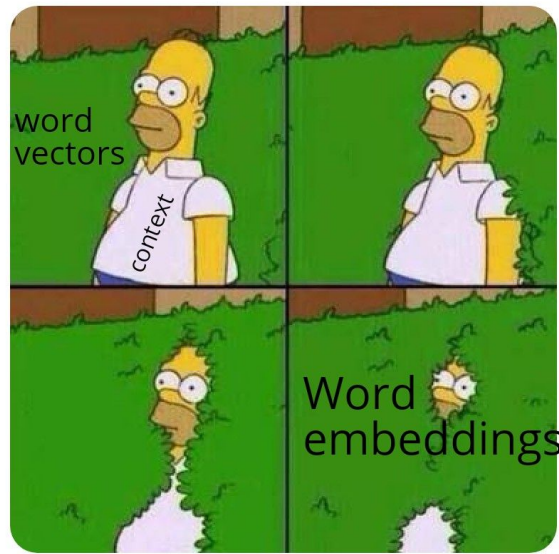




# Word Embeddings

- Embeddings are types of knowledge representation where each textual variable is represented with a vector
  - where words or phrases from the vocabulary are mapped to vectors of real numbers
  - mathematical embedding from a space with many dimensions per word to a continuous vector space with a lower dimension.
- Methods to generate this mapping include:
  - neural networks
  - dimensionality reduction on the word co-occurrence matrix
  - probabilistic models

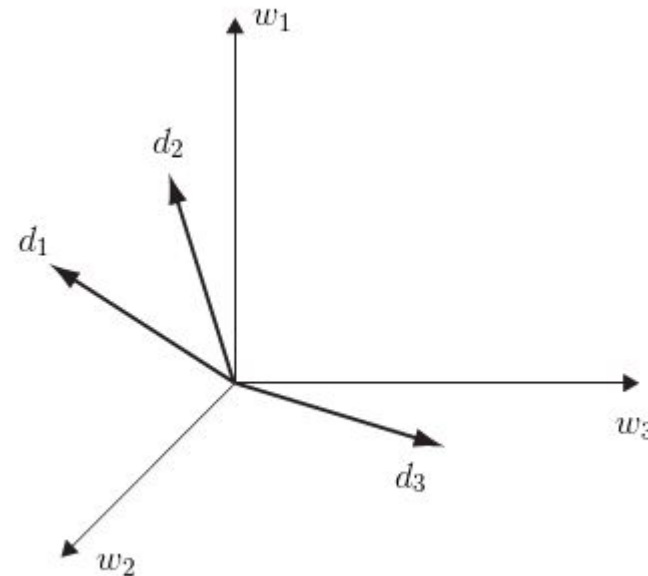
Word embeddings  
in a nutshell



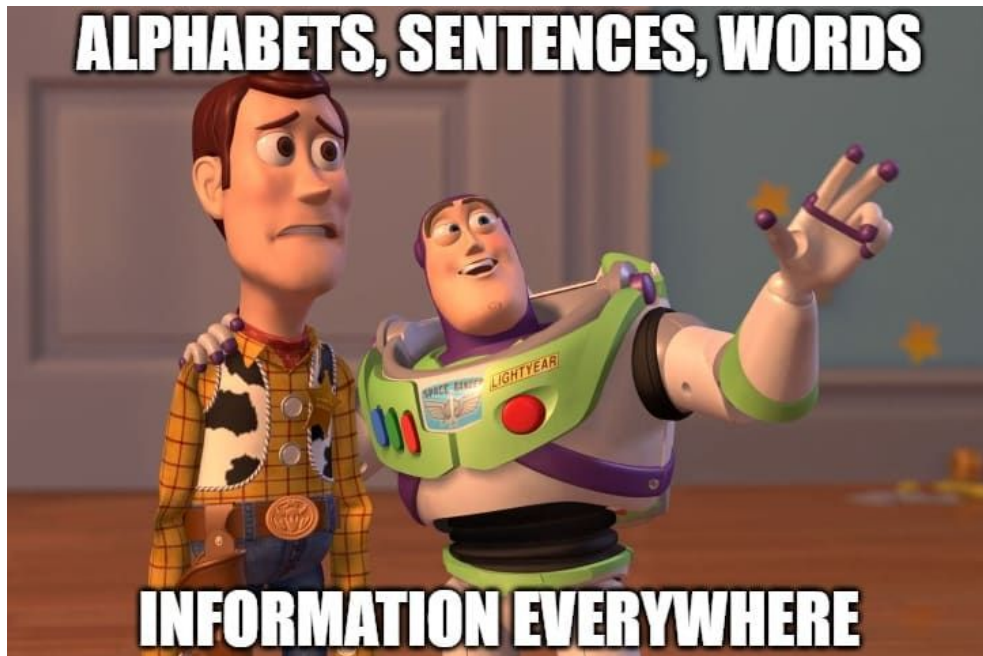
# What is a Vector Space

Vector space model is a statistical model for representing text information for -

- how many documents contain a term
- what are important terms each document has etc.
- dimensionality - whether vectors are sparse
- e.g. vocabulary size  $|V|=105$ , but documents may contain only 500 distinct words
- lexicon of document, word correlations etc.



# Google Word2Vec



# Word2Vec: Idea



Fundamental idea behind Word2vec?

**You shall know a word by the company it keeps — J.R. Firth (1957)**

It's all about **Context**

The baseline pre-trained word2vec model has -

- 300 dimensions
- 3 million 'unique' words from google news data in the training corpus

# Word2Vec: Context Matters



## **Type of Similarity 1:** Semantic Relatedness

**Def:** some relation between words

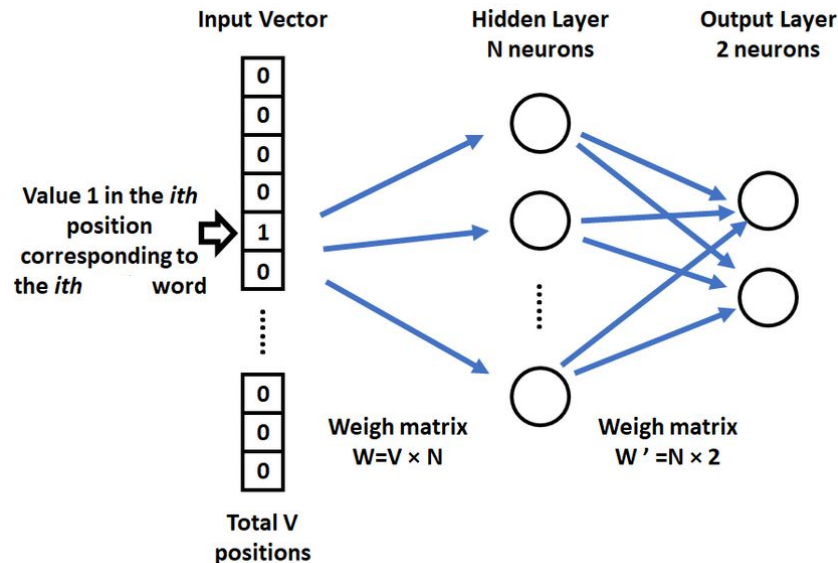
**Eg:** Truck -- Road, Bee -- Honey

## **Type of Similarity 2:** Semantic Similarity

**Def:** words used in a same way and interchangeable context

**Eg:** Car -- Auto, Doctor -- Surgeon

# Word2Vec: Definition



- A word2vec model is a **feed forward shallow neural network** model with a **single hidden layer**.
- Each word is represented by a vector (Array of numbers based on Embedding Size)
- Word2Vec finds relation (Semantic or Syntactic) between the words which was NOT possible by the Traditional TF-IDF or Frequency based approach
- Transforms the unlabeled raw corpus into labeled data (by mapping the target word to its context word), and learns the representation of words in a classification task

# Word2Vec: General Process

- When we train a model, each one hot encoded word gets a point in a dimensional space where it learns and groups the words with similar meaning
- Create an Embedding Look up Layer

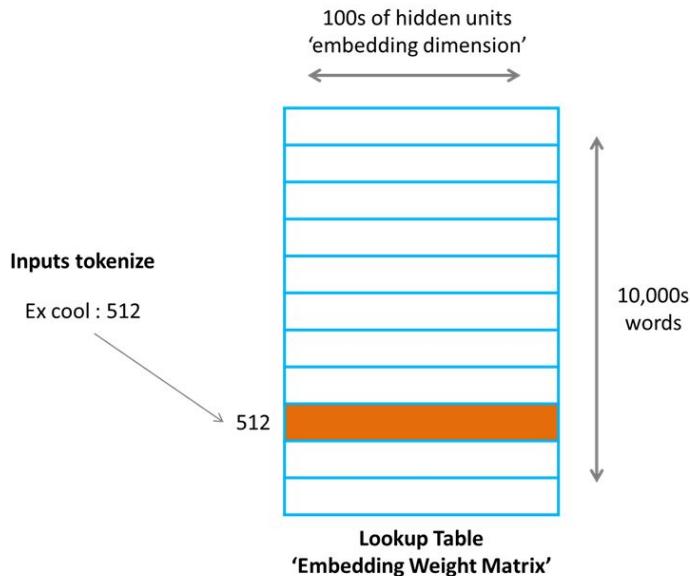
$$\begin{array}{c}
 [0 \quad 0 \quad 0 \quad 1 \quad 0] \\
 \text{One-hot vector}
 \end{array}
 \times
 \begin{array}{c}
 \begin{bmatrix}
 8 & 2 & 1 & 9 \\
 6 & 5 & 4 & 0 \\
 7 & 1 & 6 & 2 \\
 1 & 3 & 5 & 8 \\
 0 & 4 & 9 & 1
 \end{bmatrix} \\
 \text{Embedding Weight Matrix}
 \end{array}
 =
 \begin{array}{c}
 [1 \quad 3 \quad 5 \quad 8] \\
 \text{Hidden layer output}
 \end{array}$$

# Word2Vec: General Process

## Key Points:

- a) The embedding layer is just a hidden layer
- b) The lookup table is just a embedding weight matrix
- c) The lookup is just a shortcut for matrix multiplication
- d) The lookup table is trained just like any weight matrix

Word2vec falls under [prediction based embeddings](#) which tends to predict a word in a given context





# Word2Vec: Example (1)

Eg 1: Vector of the words similar

1. Cat
2. Dog

Eg 2: Vector of the words not similar

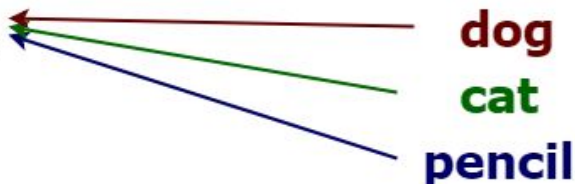
1. Cat
2. Pencil

**Similarity** is defined by the frequency of the two words in discussion

- [cat,dog]
- [cat,pencil]

And how many time they are used in the **same context**

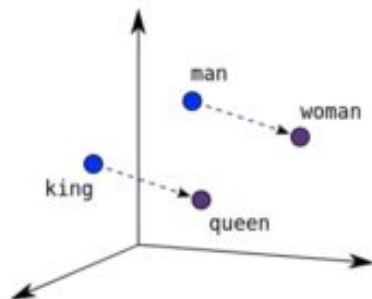
I like to pet my \_\_\_\_\_



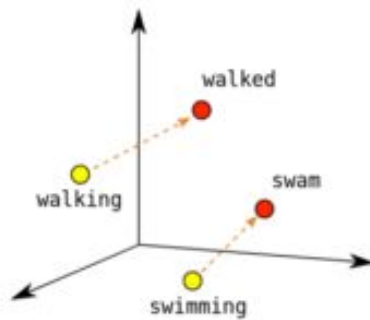
dog  
cat  
pencil

My \_\_\_\_\_ does not like the postman

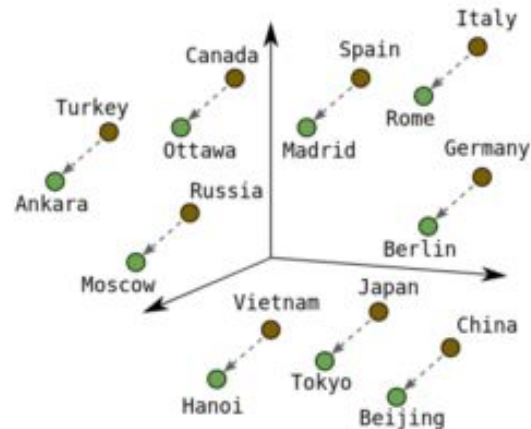
# Word2Vec: Example (2)



Male-Female



Verb Tense

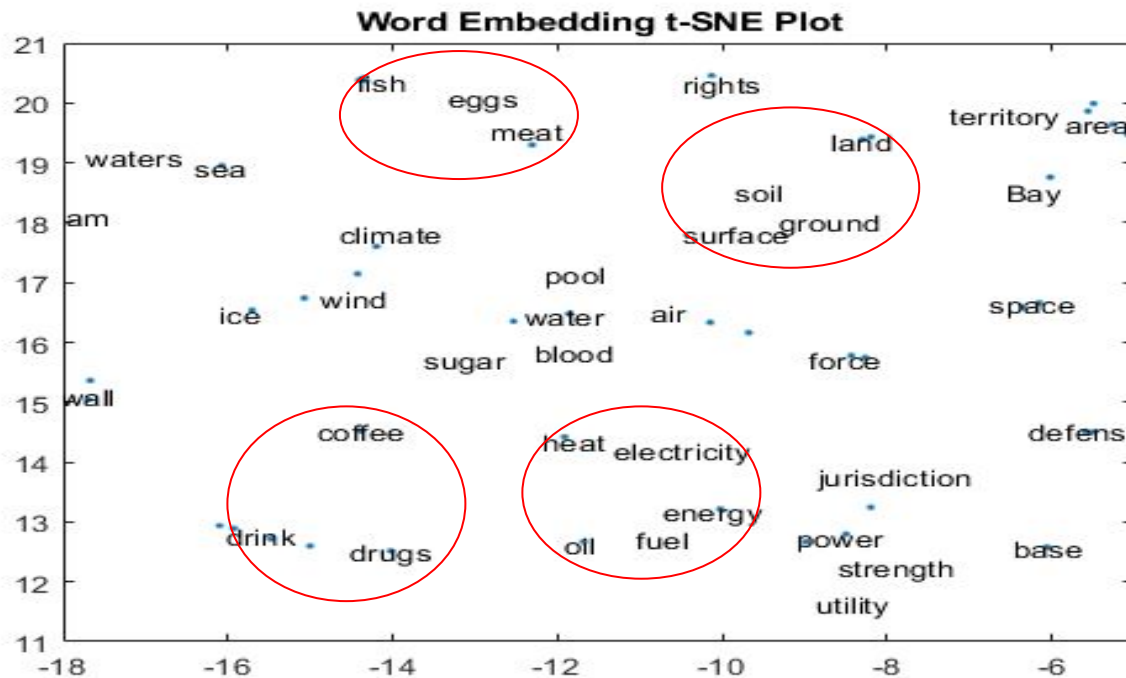


Country-Capital

Word2Vec allows some mathematical operations on vectors

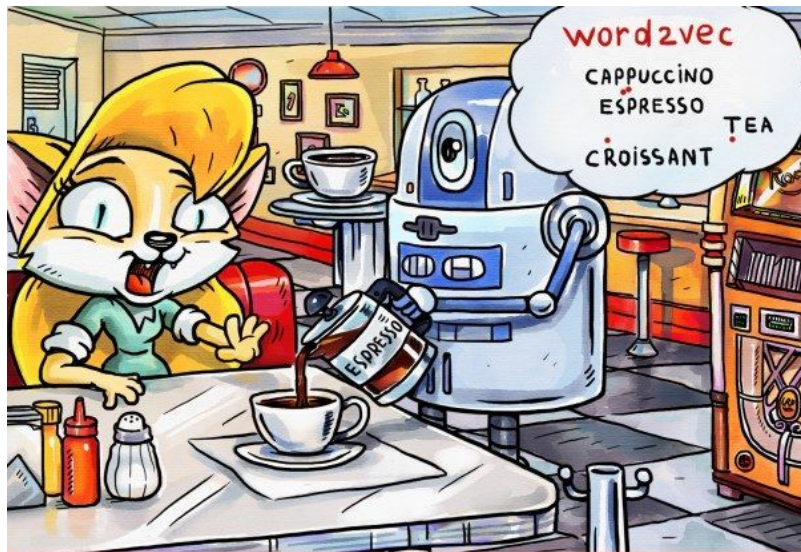
$$\text{king} - \text{man} + \text{woman} = \text{queen}$$

# Word2Vec: Example (3)



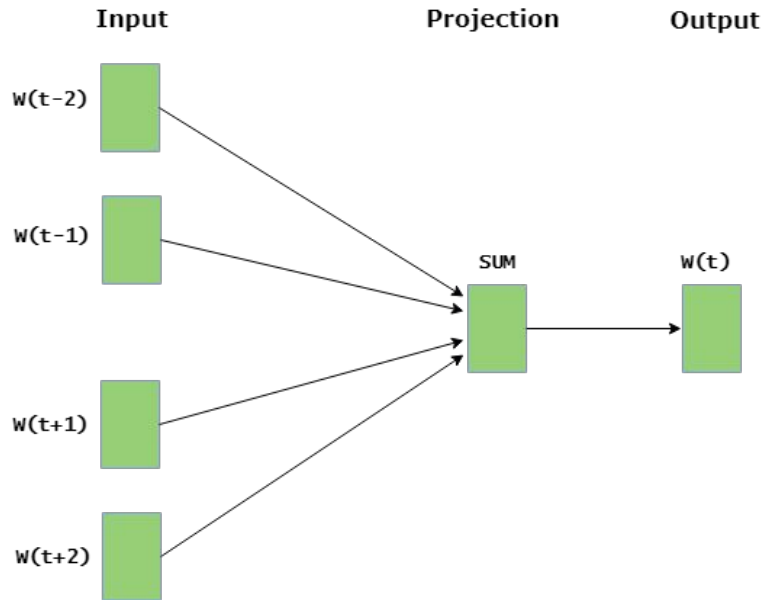
# Word2Vec: SOTA Applications

1. Topic Modeling
2. Document Similarity
3. Speech Recognition
4. Chatbots
5. Information Retrieval
6. Machine Translation
7. Question Answering
8. Recommendation Engines



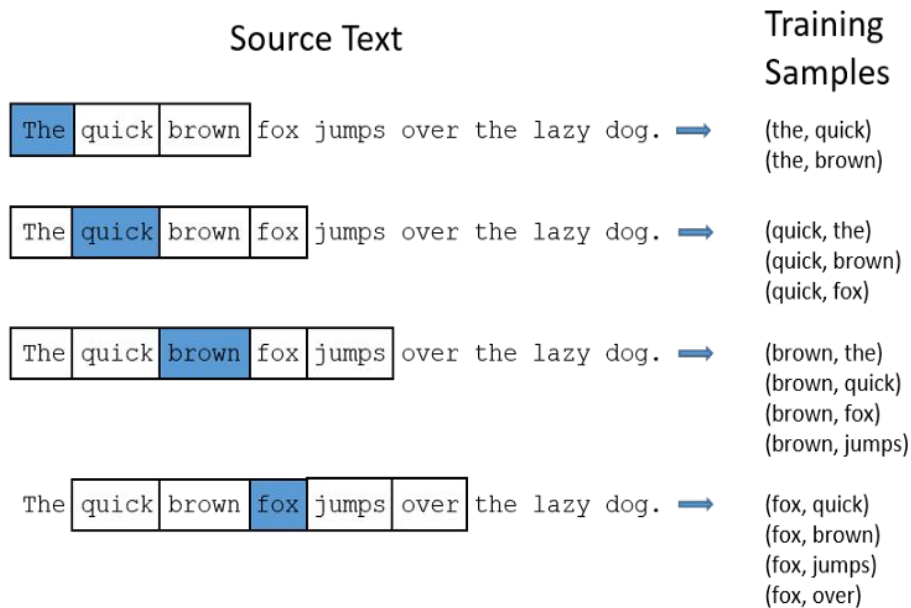
- Espresso? But I ordered a cappuccino!  
- Don't worry, the cosine distance between them is so small that they are almost the same thing.

# Word2Vec: Architecture (CBOW)



- CBOW - **predict** the current word based on context. Here the *input* will be the **context neighboring** words and *output* will be the **target word**.
- Based on Window Method, where we have to assign a Window size
- If window size is set to 1. So, 1 word from both the sides of target are being considered. Similarly, in each iteration, window will slide by single stride and our neighbors will keep changing

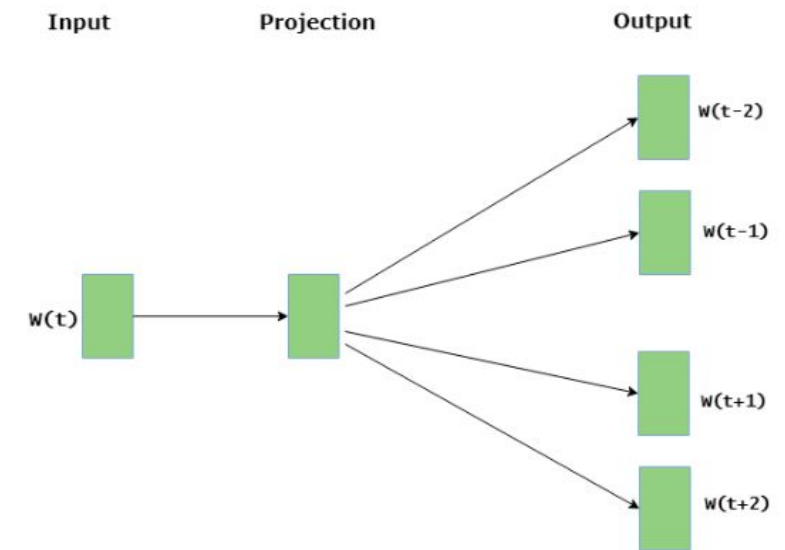
# Word2Vec: Architecture (CBOW)



- Model: CBOW
- INPUT Layer: White box content
- TARGET Layer: blue box word
- Window Size: 5

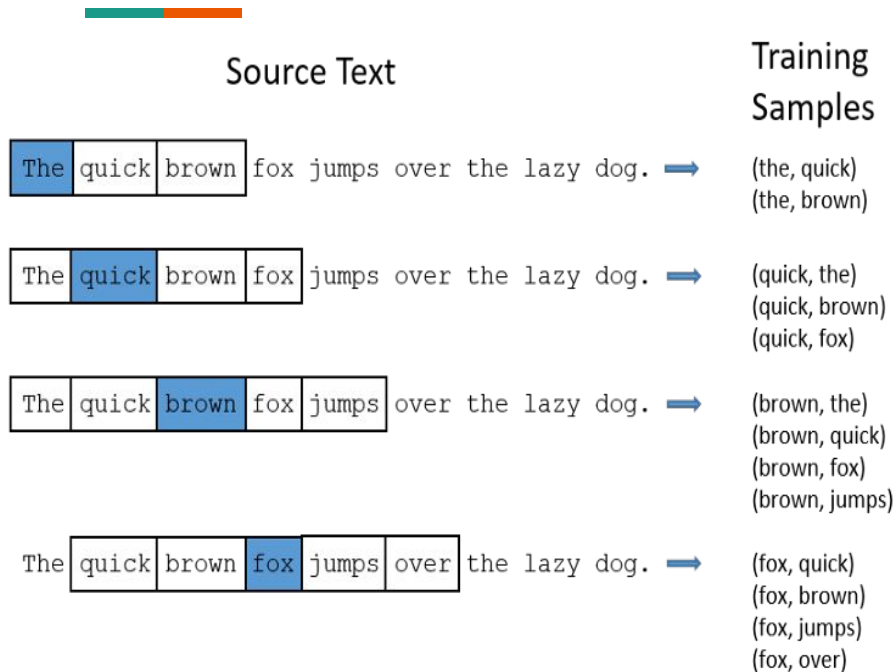
*CBOW architecture predicts the current word based on the context*

# Word2Vec: Architecture (skip-grams)



- Skip gram predicts the surrounding context words within specific window given current word
- The input layer contains the current word
- The output layer contains the context words
- The hidden layer contains the number of dimensions in which we want to represent current word present at the input layer

# Word2Vec: Architecture (skip-grams)



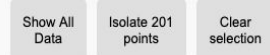
- Model: Skip Gram
- INPUT Layer: Blue box word
- TARGET Layer: White box content
- Window Size: 5

*Skip-gram architecture predicts surrounding words given the current word*



WOMEN WHO  
**CODE**  
DATA SCIENCE

Points: 820 | Dimension: 200 | Selected 201 points



Search by  
politics

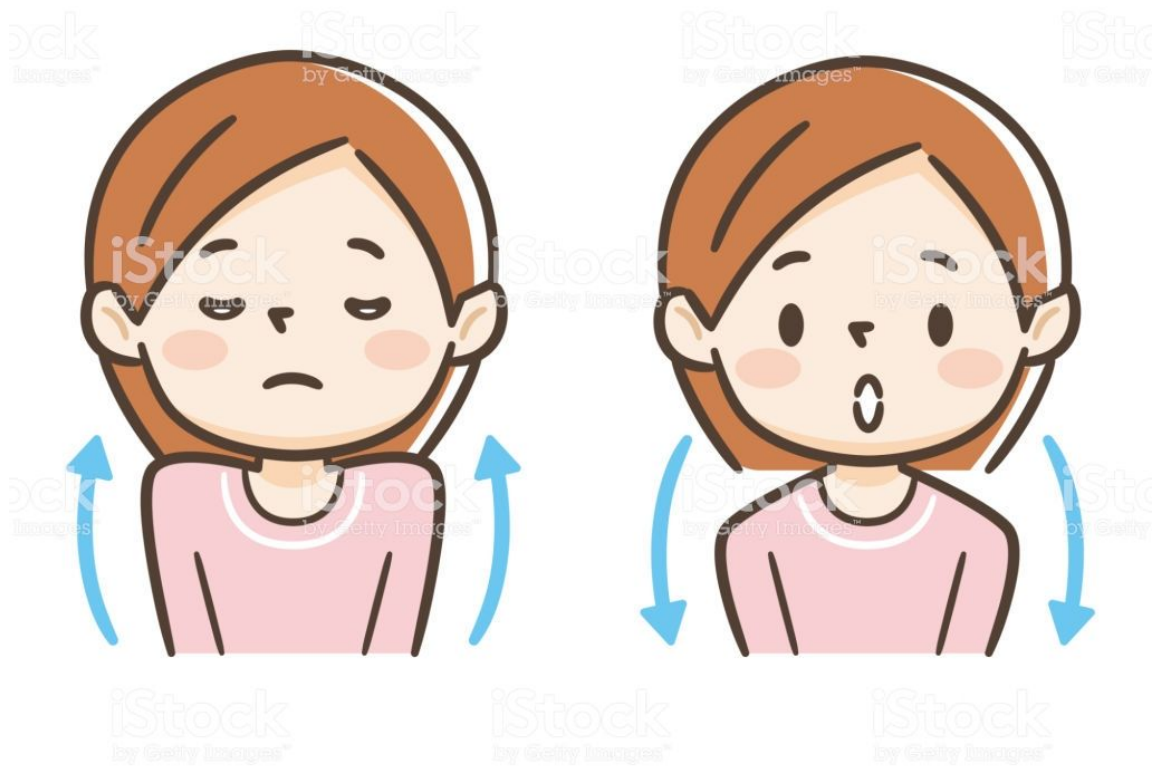
neighbors ?  200

distance COSINE EUCLIDEAN

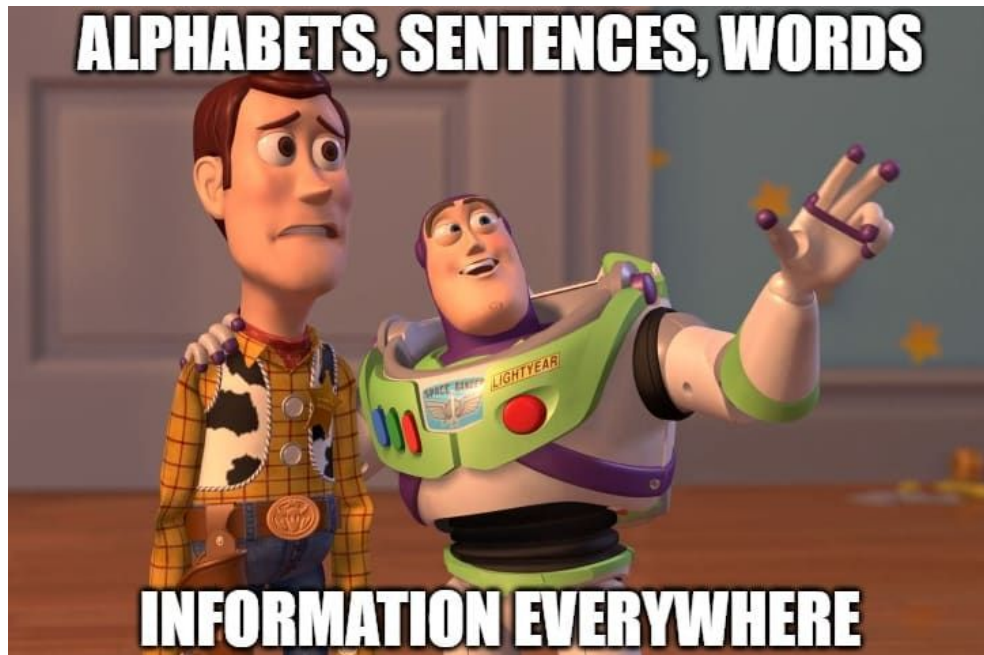
Nearest points in the original space:

economy	0.697
demographics	0.705
political	0.708
geography	0.714
affairs	0.731
nationalism	0.737
philosophy	0.738
wing	0.745
democracy	0.747
imperialism	0.748
socialist	0.752
religion	0.753
ideology	0.757
socialism	0.765
diplomacy	0.772
democratic	0.775
culture	0.779
politically	0.783
presidency	0.783
liberal	0.795

Let's take a few deep breaths now!



# Stanford gloVe



# gloVe: Idea



Fundamental idea behind gloVe?

Word2vec: *relies only on **local information** of language. The semantics learnt for a given word, is only affected by the surrounding words*

Eg: “The cat sat on the mat”

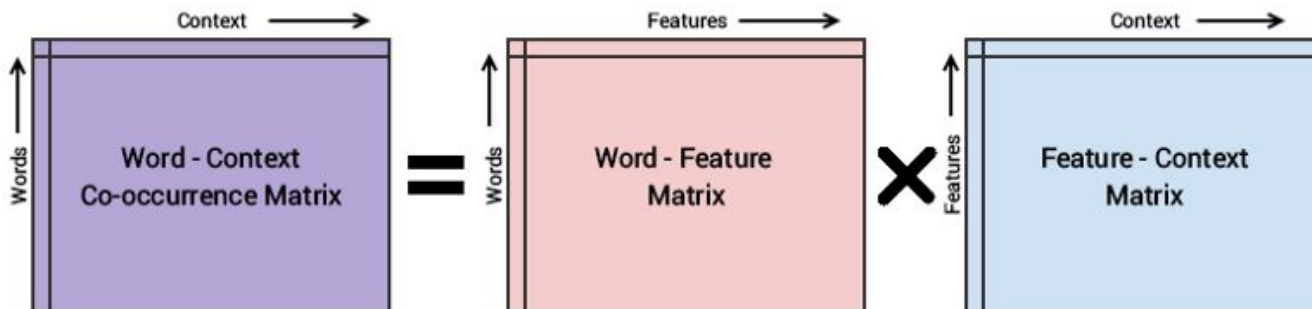
Con: is “the” a special context of the words “cat” and “mat”?

**gloVe**: *You can derive semantic relationships between words from the co-occurrence matrix.*

It's all about **Global Context**. GloVe stands for “**Global Vectors**”

# gloVe: Definition

- Unsupervised learning algorithm for obtaining vector representations for words.
- Training is performed on aggregated global word-word co-occurrence statistics from a corpus
- Resulting representations showcase linear substructures of the word vector space



# gloVe: Process



The relationship of these words can be examined by studying the ratio of their co-occurrence probabilities with various probe words,  $k$ .

**Eg:**

$P(k|w)$ : probability that the word  $k$  appears in the context of *word  $w$*

Consider a word strongly related to *ice*, but not to *steam*, such as **solid**

$P(\text{solid} | \text{ice})$  will be relatively high

$P(\text{solid} | \text{steam})$  will be relatively low

Ratio of  $P(\text{solid} | \text{ice}) / P(\text{solid} | \text{steam})$  will be large

# gloVe: Process

Consider the word **gas** that is related to **steam** but not to **ice**,

Ratio of  $P(\text{gas} \mid \text{ice}) / P(\text{gas} \mid \text{steam})$  will be small

For a word **related** to both **ice** and **steam**, such as **water**, the ratio would be close to one

For words related to **neither** ice nor steam, such as **fashion**, the ratio would be close to one

Probability and Ratio	$k = \text{solid}$	$k = \text{gas}$	$k = \text{water}$	$k = \text{fashion}$
$P(k \text{ice})$	$1.9 \times 10^{-4}$	$6.6 \times 10^{-5}$	$3.0 \times 10^{-3}$	$1.7 \times 10^{-5}$
$P(k \text{steam})$	$2.2 \times 10^{-5}$	$7.8 \times 10^{-4}$	$2.2 \times 10^{-3}$	$1.8 \times 10^{-5}$
$P(k \text{ice})/P(k \text{steam})$	8.9	$8.5 \times 10^{-2}$	1.36	0.96

gloVe falls under count based embeddings capturing global co-occurrences and needs an upfront pass of full data during training

# Word2vec vs gloVe

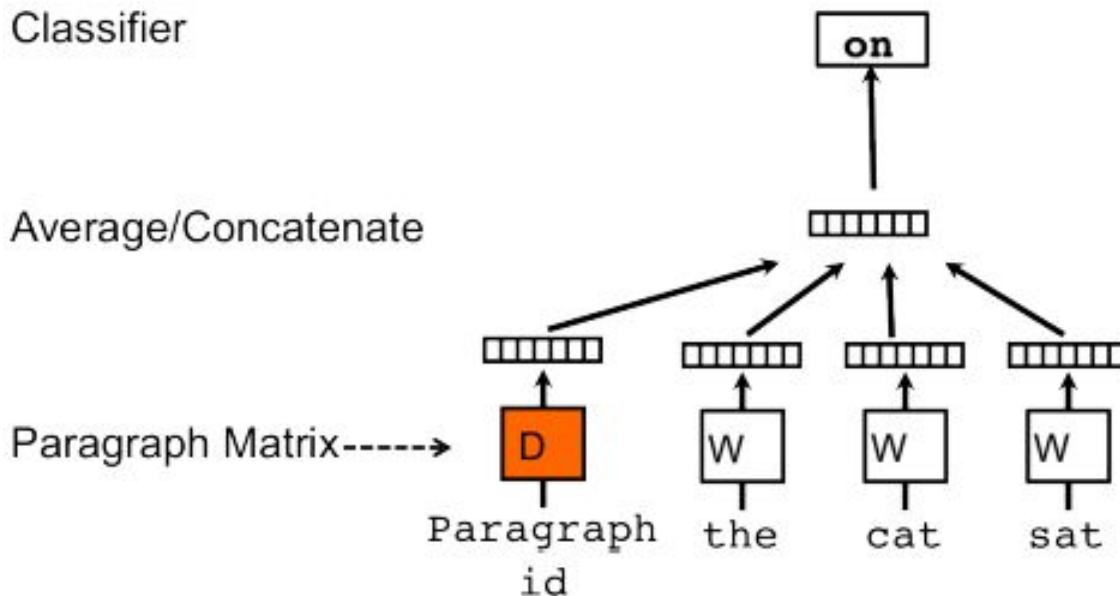


PLEASE  
NOTE...

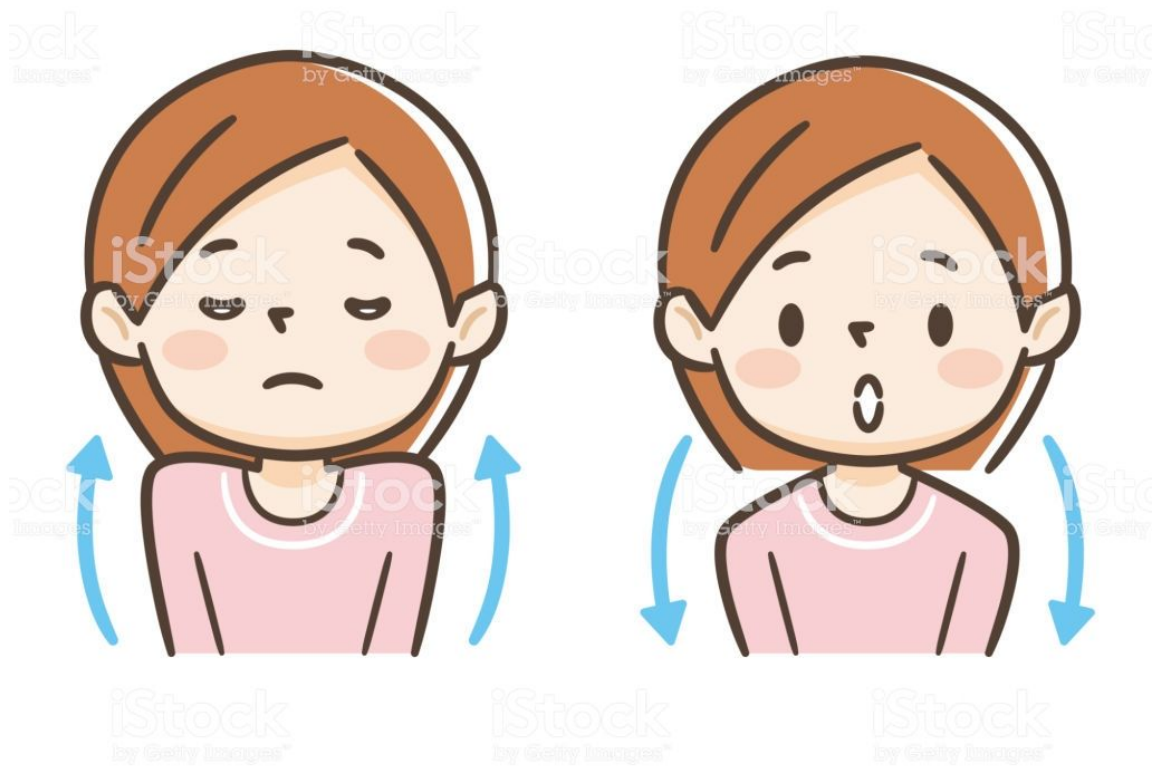
- Word2Vec and GloVe models are very similar in how they work
- Both aim to build a vector space where the position of each word is influenced by its neighboring words based on their context and semantics.
- Word2Vec starts with local individual examples of word co-occurrence pairs
- GloVe starts with global aggregated co-occurrence statistics across all words in the corpus.



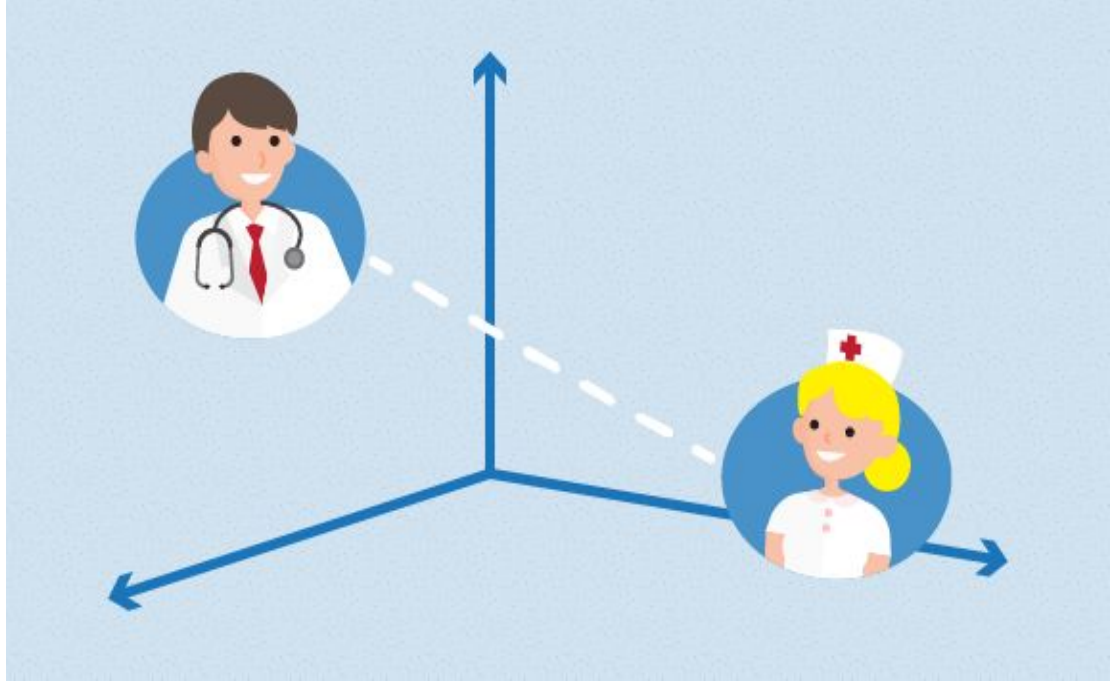
# One step up - Doc2vec



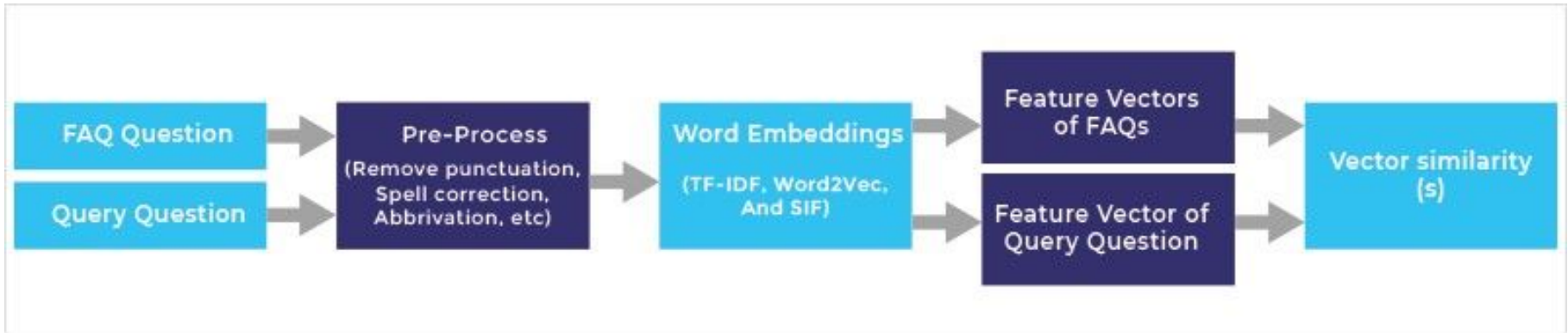
Let's take a few deep breaths now!



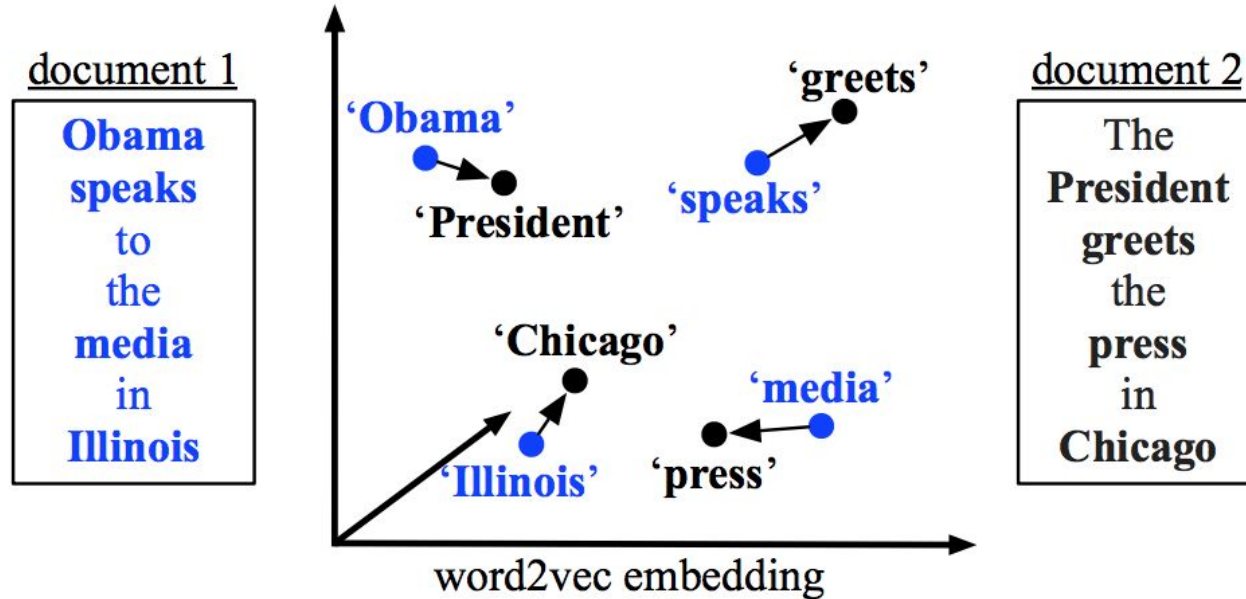
# Use-case: Similarity Scoring



# Building Phrase/ Document Similarity Model



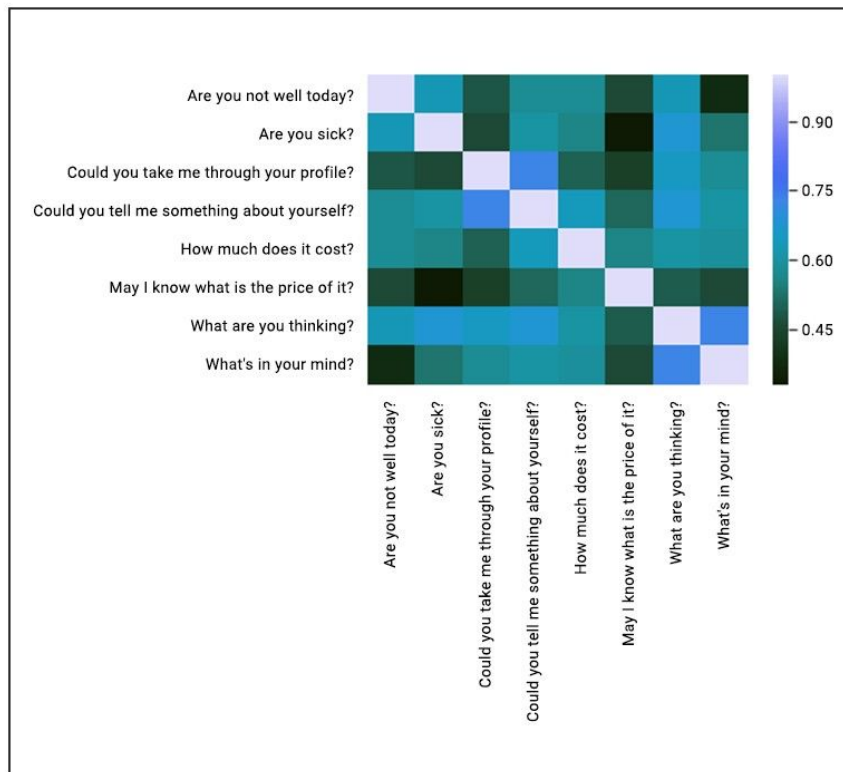
# Similarity Metrics - Word Mover Distance



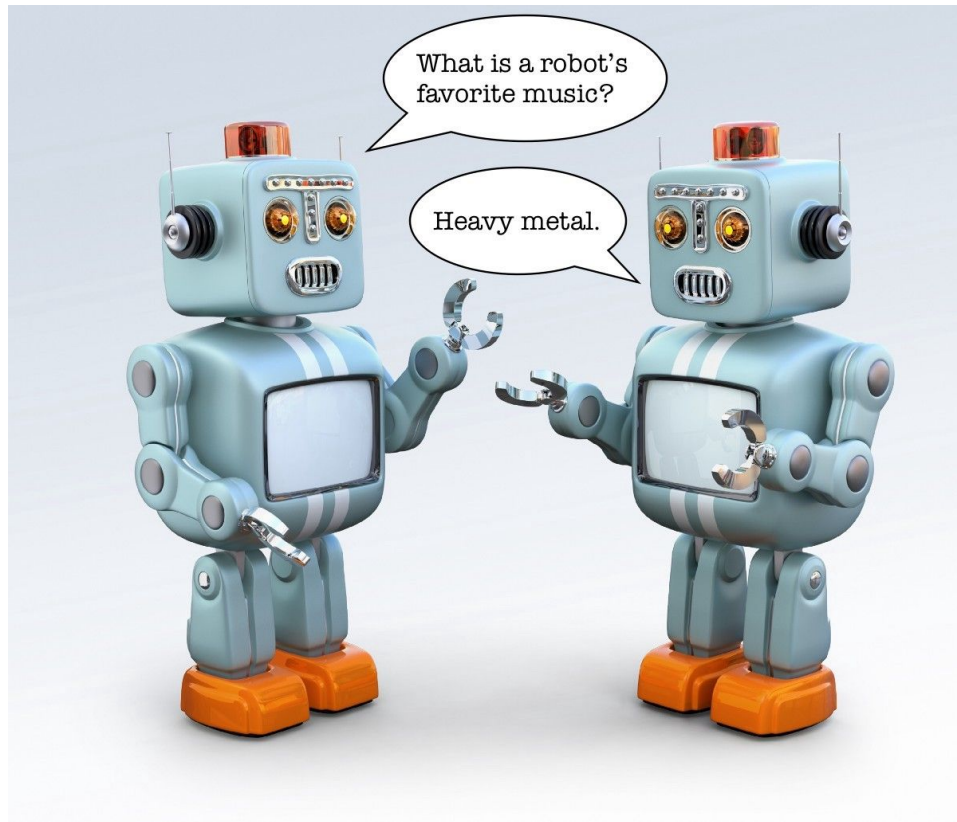
# Combination of various pre-processing & word-embedding techniques

Pre-Process	Feature Extraction	Similarity Threshold	Questions Answered	Precision
Stop words removal	TF-IDF (KB vocab)	0.7	56.85%	90.40%
+ Punctuation removal	word2vec	0.7	27.60%	97.80%
+ Lemmatization	(custom trained)			
+ Spell Correction	TF-IDF (KB vocab)	0.7	40.12%	96.48%
+ Abbreviations	+ word2vec (custom trained)			
	SIF	0.7	87.50%	73.72%
	SIF + Clustering	0.8	68.75%	90.62%
	SIF with word2vec or custom trained phrases	0.8	64.11%	93.08%
+ Custom stop words.	SIF	0.8	74.60%	92.70%
	+ TF-IDF (KB vocab)			
	+ word2vec(custom trained)			

# Evaluation - Confusion Matrix

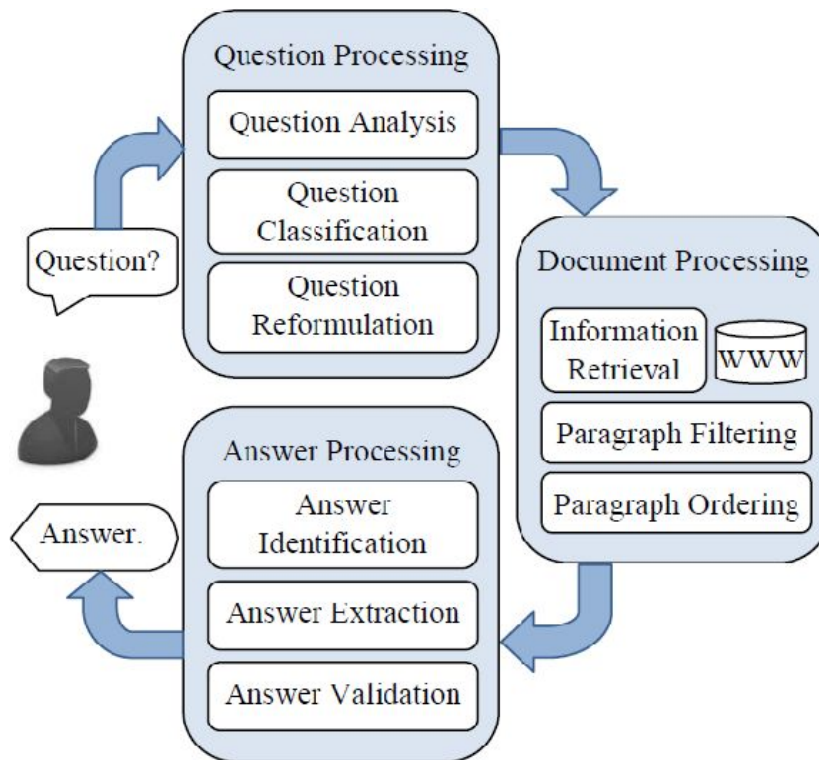


# Use-case: Question Answering

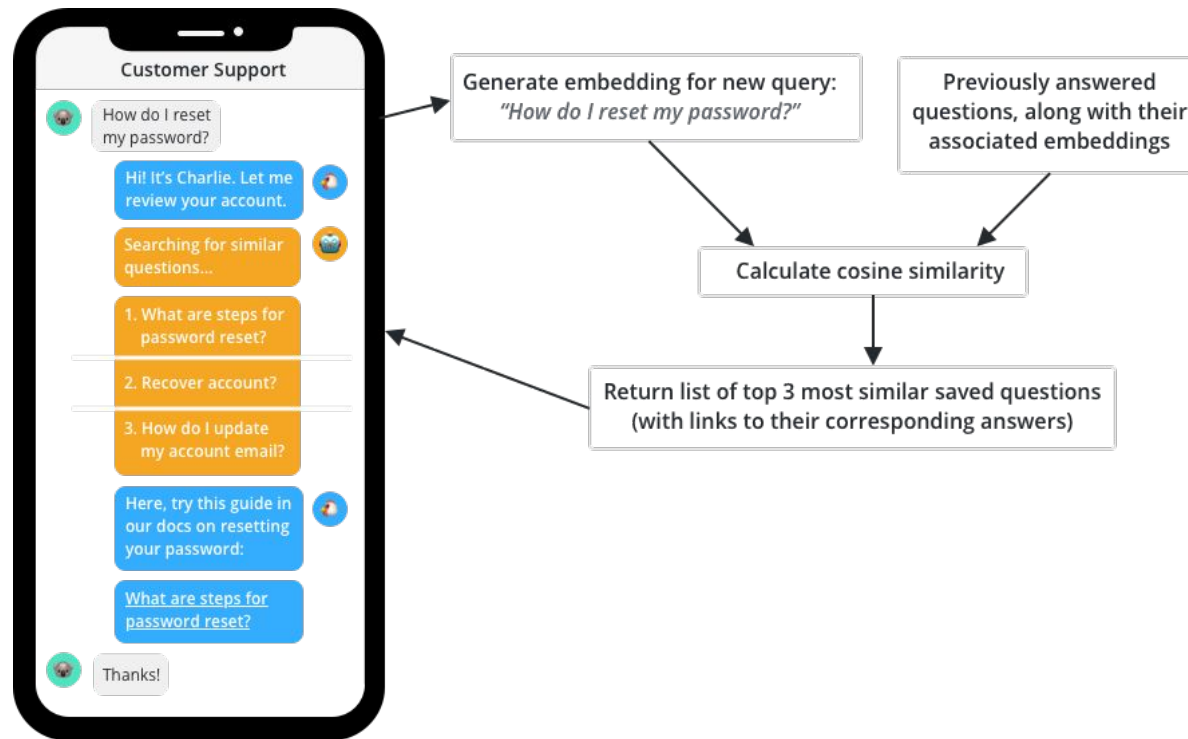




# Building QA Model



# Applications - Chat Assistant



# Applications - Information Retrieval

Sentence having the right answer

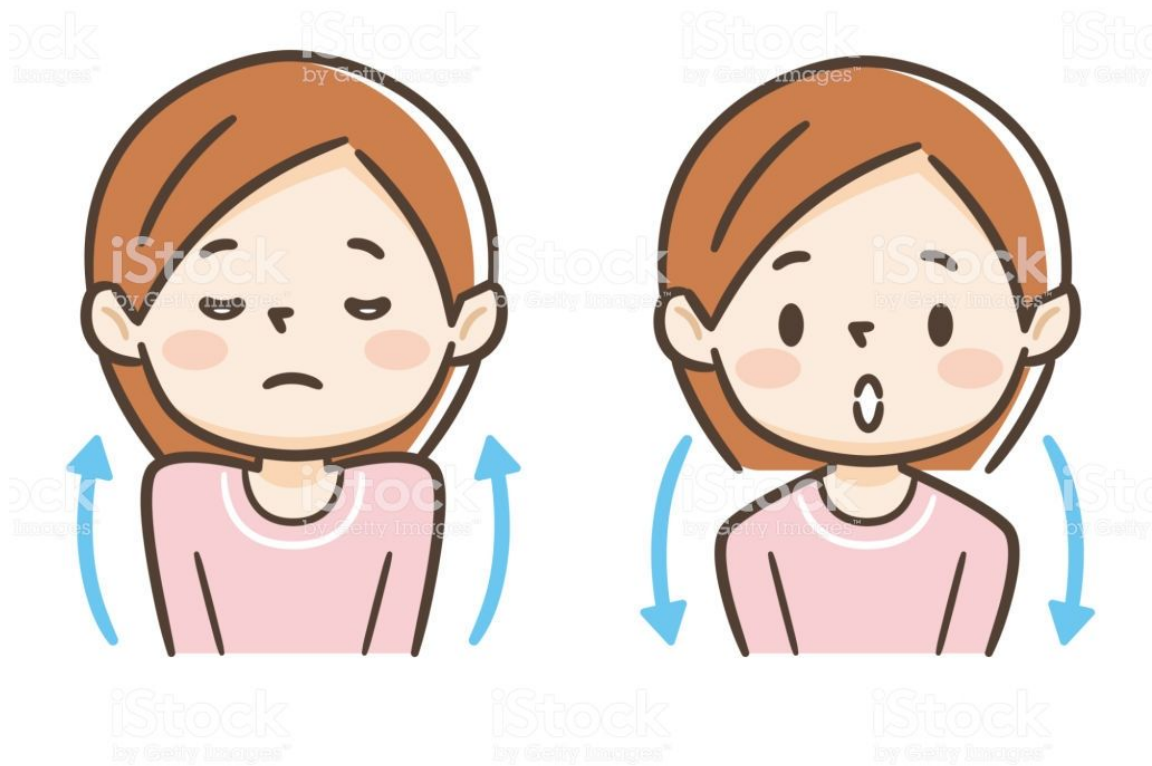
**'context':** 'Beyoncé Giselle Knowles-Carter (/bi:ˈjɒnsər/ bee-YON-say) (born September 4, 1981) is an American singer, songwriter, record producer and actress. Born and raised in Houston, Texas, she performed in various singing and dancing competitions as a child, and rose to fame in the late 1990s as lead singer of R&B girl-group Destiny's Child. Managed by her father, Mathew Knowles, the group became one of the world's best-selling girl groups of all time. Their hiatus saw the release of Beyoncé's debut album, *Dangerously in Love* (2003), which established her as a solo artist worldwide, earned five Grammy Awards and featured the Billboard Hot 100 number-one singles "Crazy in Love" and "Baby Boy".',

**'text':** 'in the late 1990s'

**'question':** 'When did Beyonce start becoming popular?'

Exact Answer

Let's take a few deep breaths now!



# State of NLP today



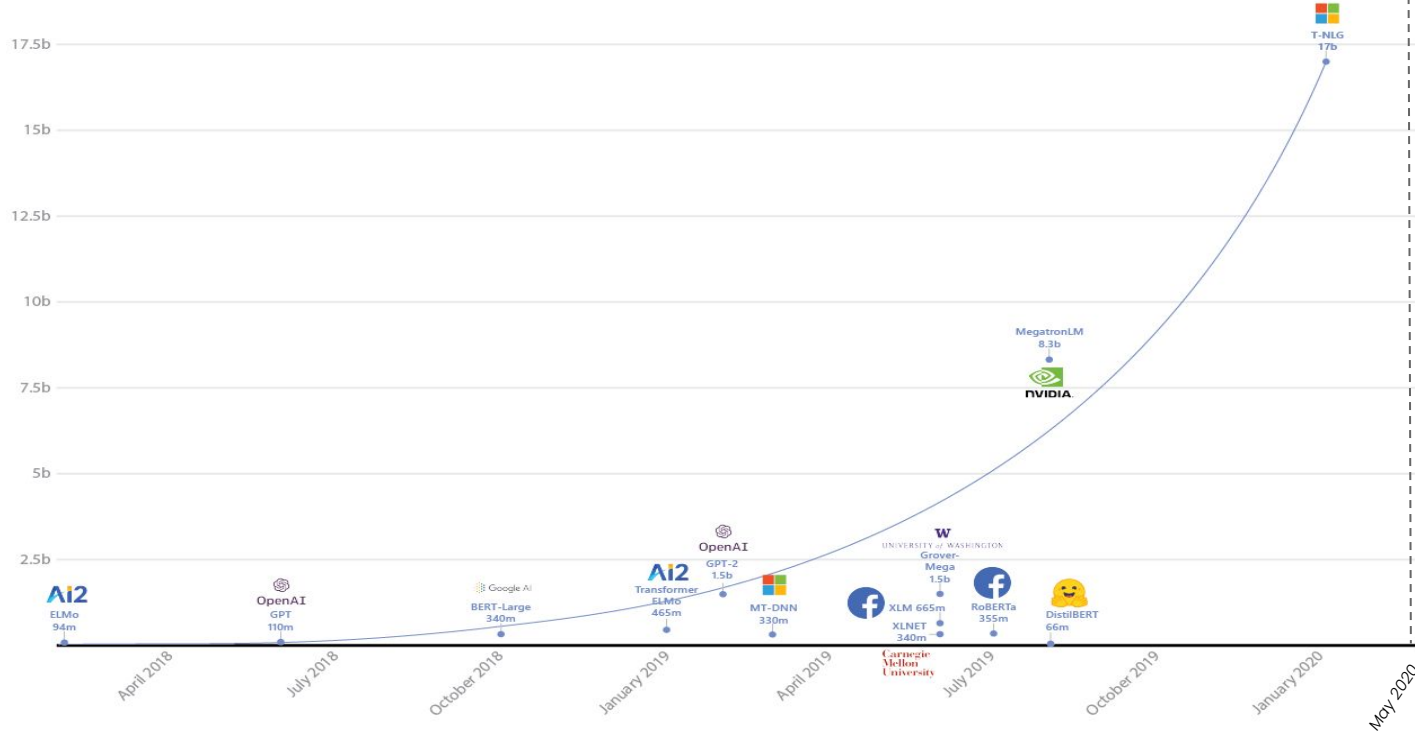


Open AI  
175B

WOMEN WHO  
**CODE**  
DATA SCIENCE

# NLP Growth Curve

Parameters  
(in B)



Timeline

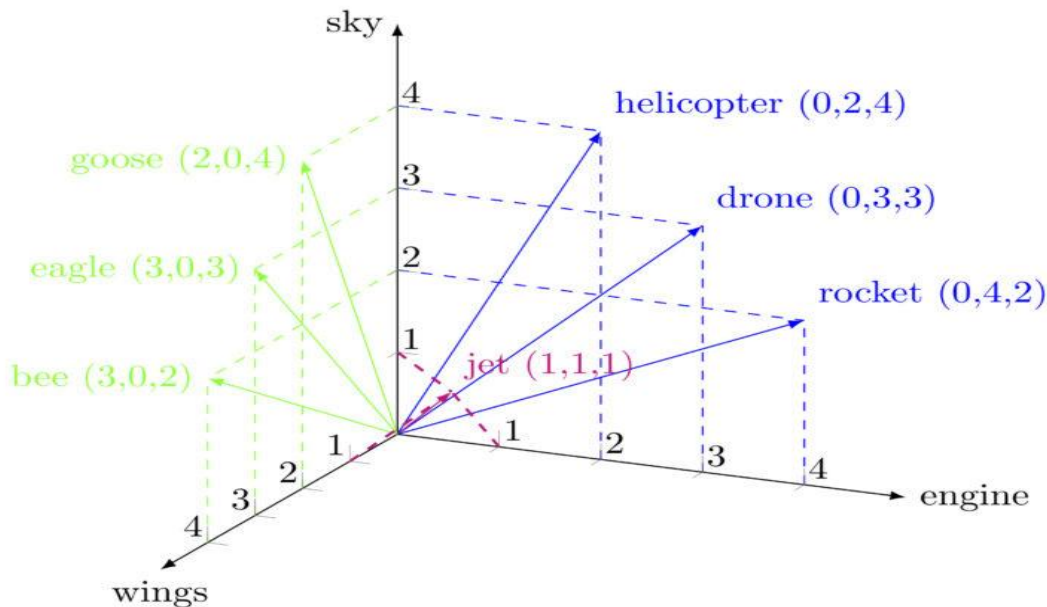
Source: <https://www.microsoft.com/en-us/research/blog/turing-nlg-a-17-billion-parameter-language-model-by-microsoft/>



# Overview of BERT



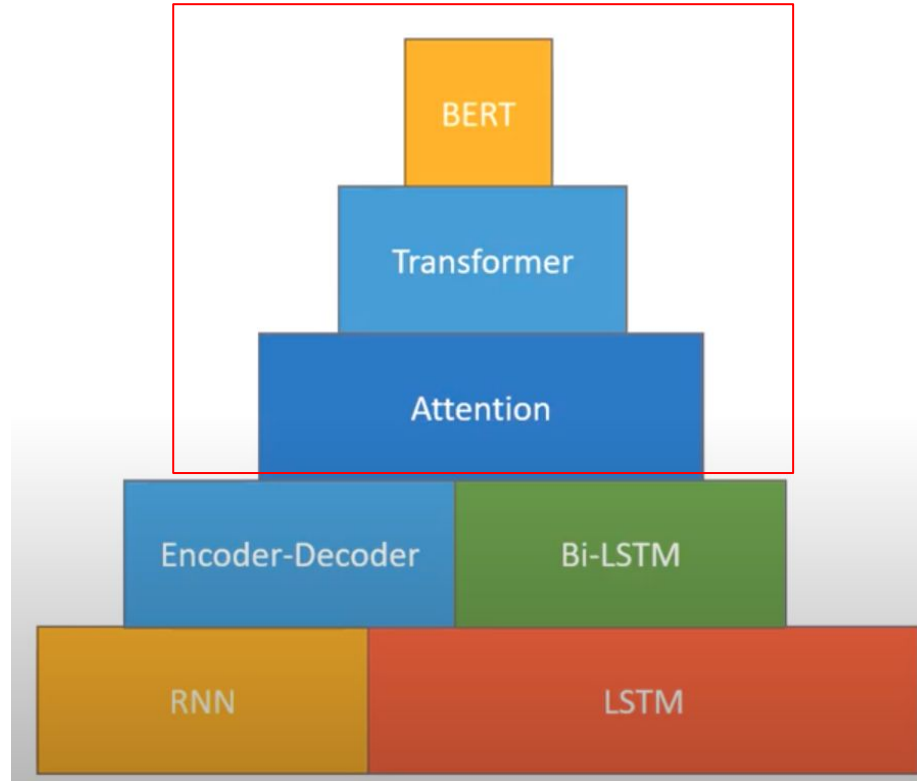
# Base Concept



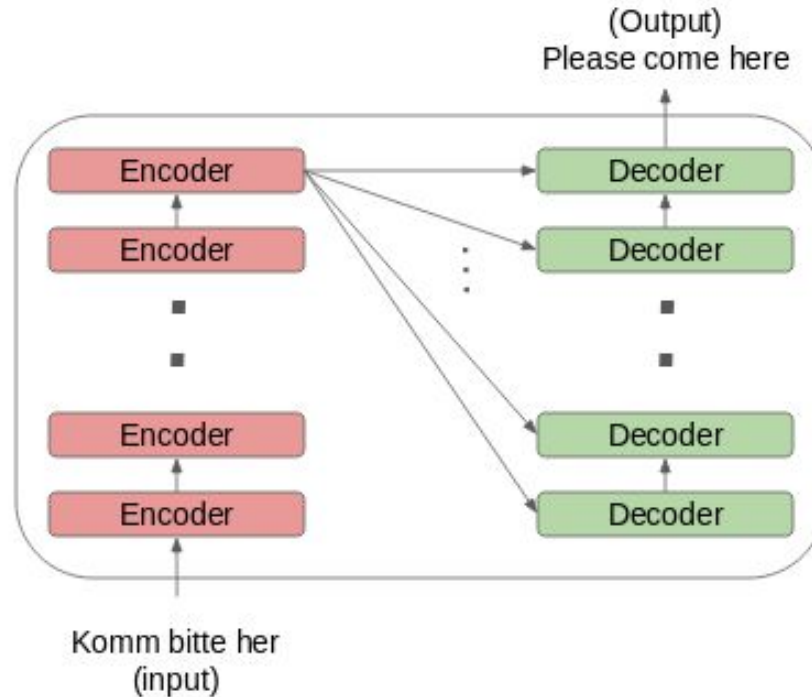
An important concept we learnt - **Word Embeddings** is the base  
Word Embedding = Feature Vector representation of a word



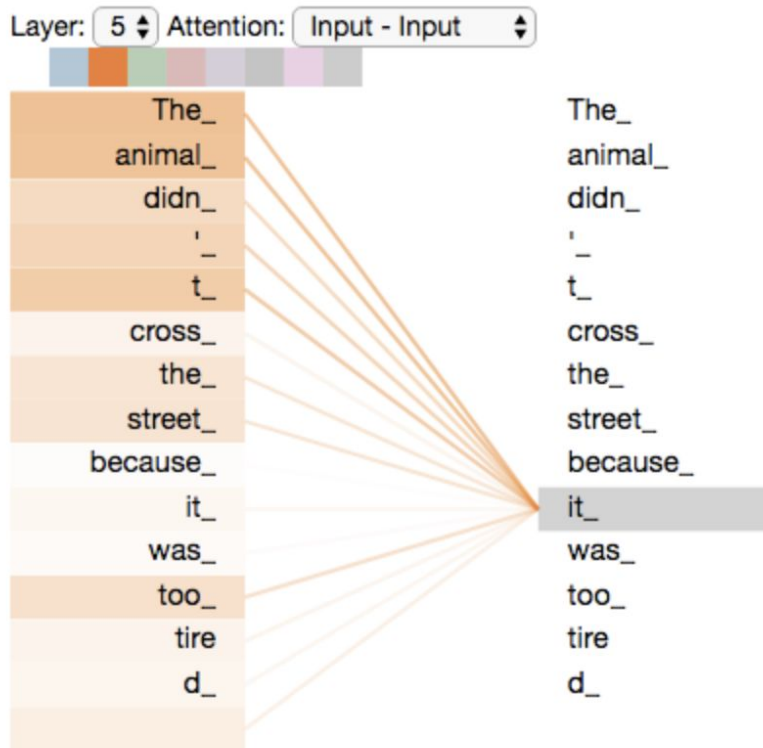
# BERT Mountain - by Chris McCormick




# General Idea of Transformers



# Attention is all you need!



# BERT: Bidirectional Encoder Representations from Transformer



## Trained Model:

Wikipedia  
(2,500M words),  
Book Corpus  
(800M words)

## Concept:

Bidirectional means  
that BERT learns  
information from both  
the left and the right  
side of a token's  
context during training

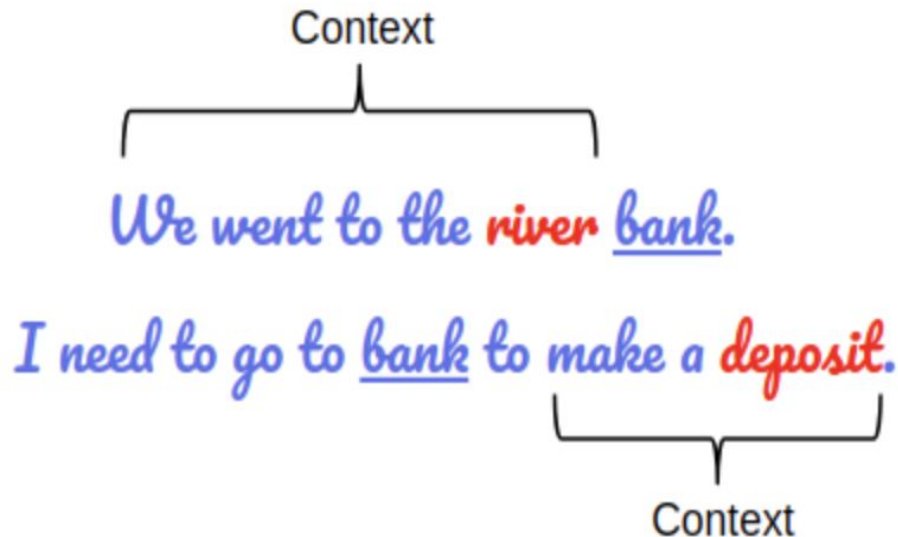
## Architecture:

BERT Base -  
-12 layers (transformer  
blocks),  
-768-hidden  
-12 attention heads,  
-110M parameters

**Trend Setter: First deeply bidirectional, unsupervised, pre-trained on plain text**

# BERT: Goal

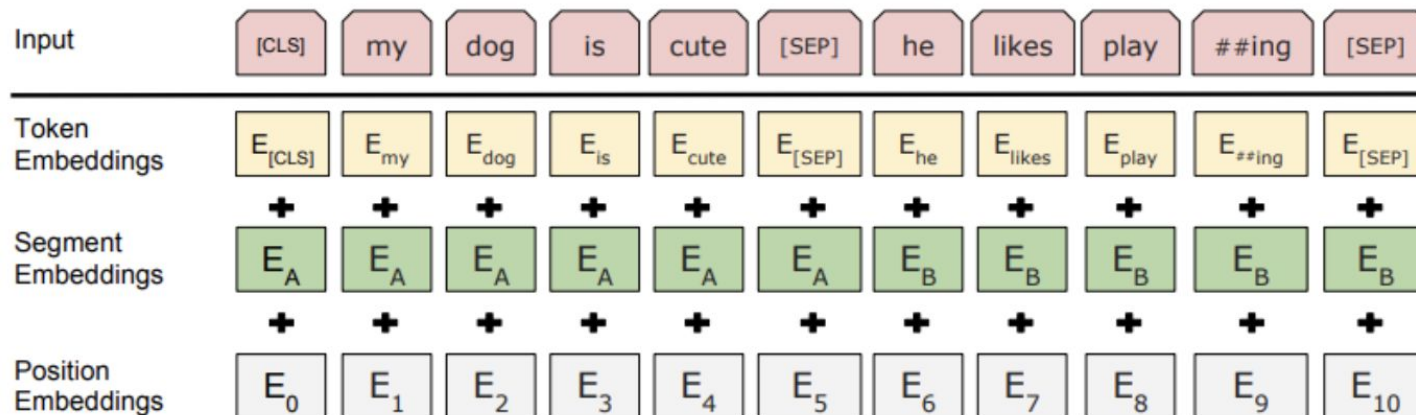
Goal of BERT: Contextual Input Representation and Word Piece Embeddings



# BERT: Layers

Three embedding layers

- Token
- Position - word to word relations
- Segment - Sentence to sentence relations



# BERT: Pre-trained on two NLP tasks

- Masked Language Modeling
- Next Sentence Prediction

## MLM: The Strength of Bidirectionality

**Input:** The man went to the [MASK]<sub>1</sub> . He bought a [MASK]<sub>2</sub> of milk .

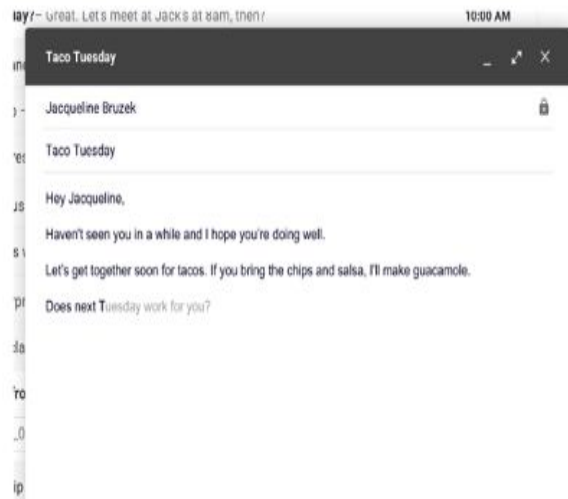
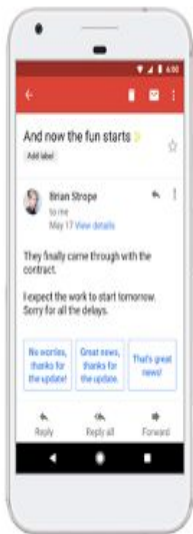
**Labels:** [MASK]<sub>1</sub> = store; [MASK]<sub>2</sub> = gallon

Next Sentence: model relationships between sentences

**Sentence A** = The man went to the store.  
**Sentence B** = He bought a gallon of milk.  
**Label** = IsNextSentence

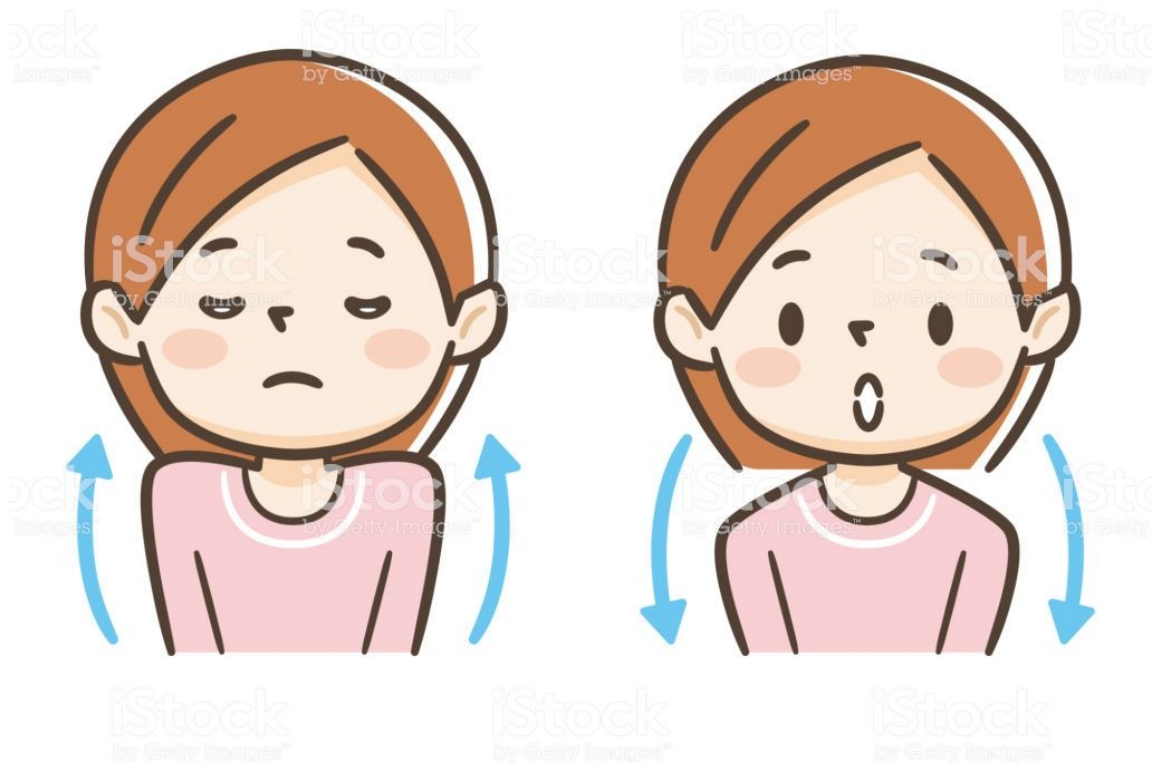
**Sentence A** = The man went to the store.  
**Sentence B** = Penguins are flightless.  
**Label** = NotNextSentence

# Some important applications



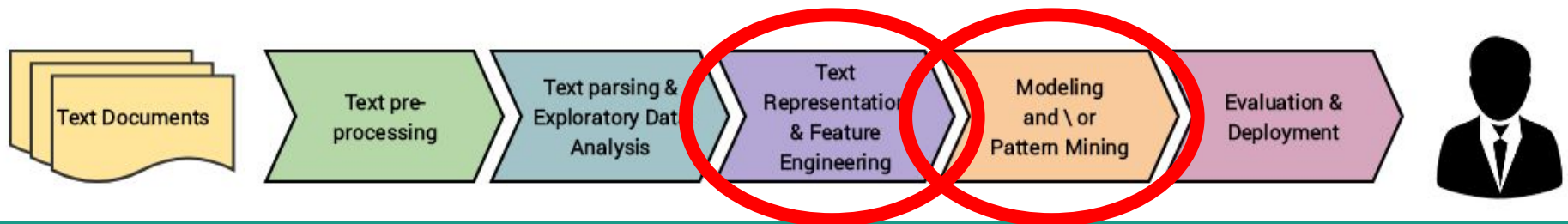


Let's take a few deep breaths now!



# 4

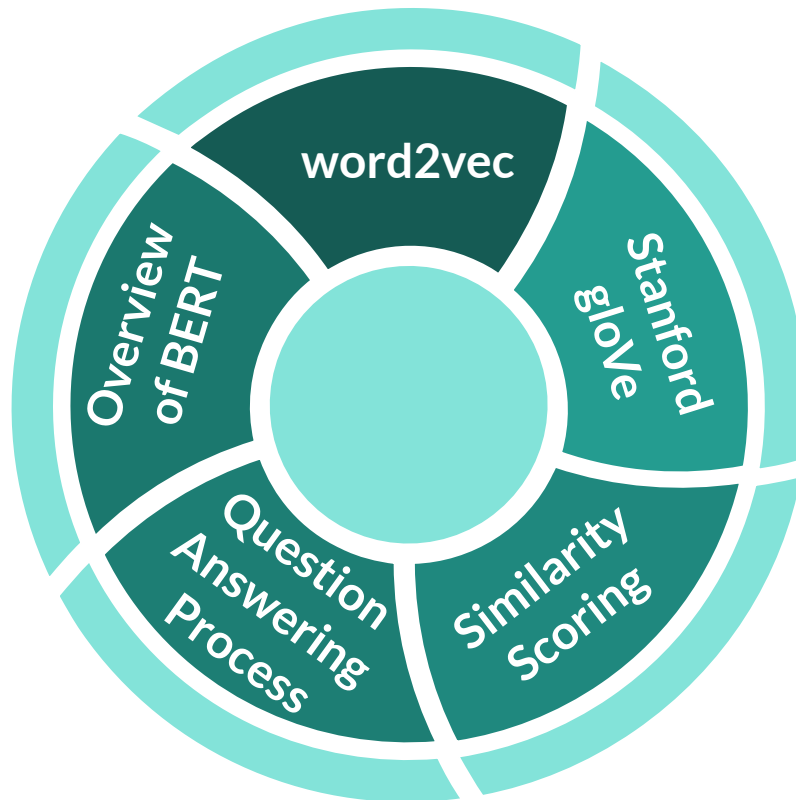
# NLP Workflow



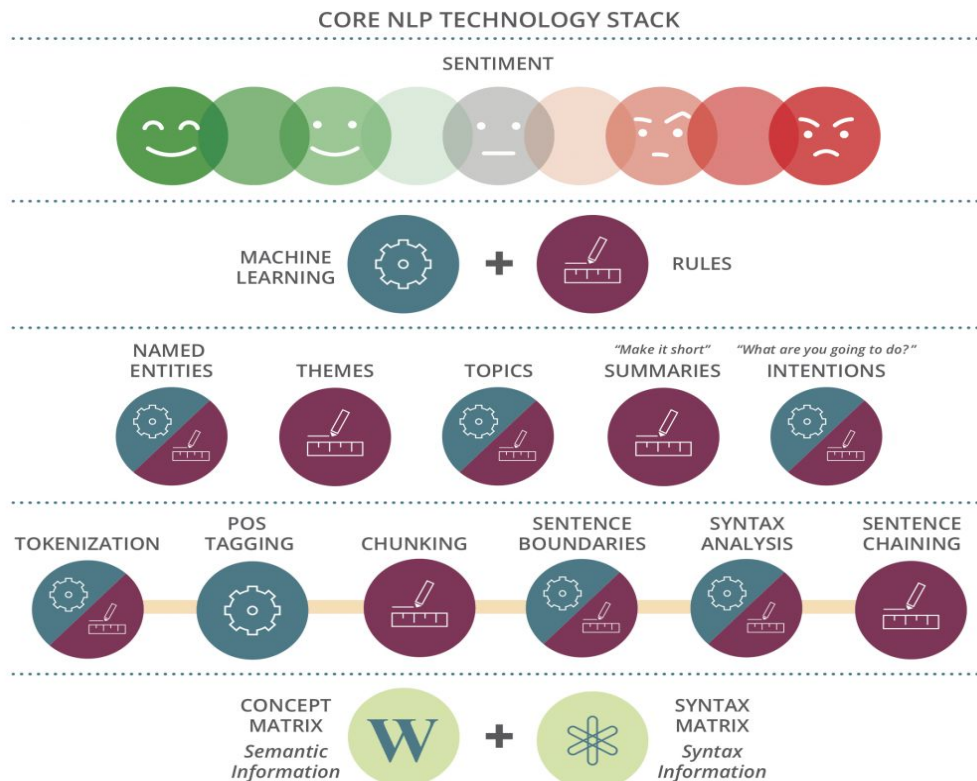
5

# Theory Wrap-up & Next Steps

# Recap



# Recap



# 6

# Google Colab Project

<https://bit.ly/introtonlp-week4-notebook>

# Homework #1

## Additional Resources

- [\(Jalammar\) The Illustrated Word2vec](#)
- [\(TowardsDataScience\) Intuitive Guide to Understanding Word2vec](#)
- [\(Chris McCormick\) Applying word2vec to Recommenders and Advertising](#)
- [\(AnalyticsVidhya\) Word2Vec using Gensim](#)
- [\(TowardsDataScience\) word2vec-Predict Article Success](#)

See you  
next week!

## Questions?

Join us on [Slack](#) and post your questions  
to the #help-me channel