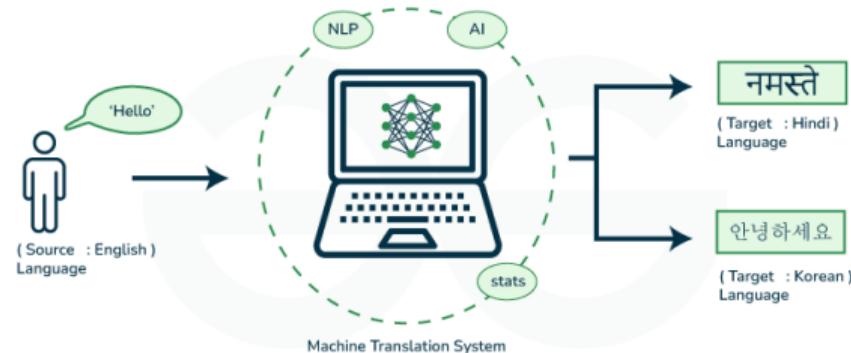


Natural Language Processing

- Language is central to human interaction; many of our daily activities revolve around text and language.
 - Natural Language Processing (NLP) enables computers to understand and generate human language.

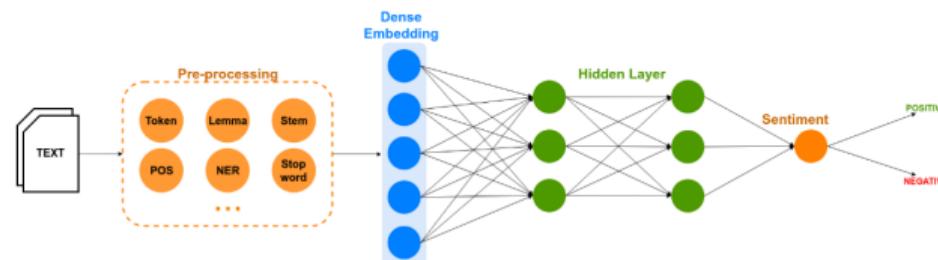
Translation



- NLP helps translate text from one language to another.

Figure adapted from [geeksforgeeks.org/machine-translation-of-languages-in-artificial-intelligence/](https://www.geeksforgeeks.org/machine-translation-of-languages-in-artificial-intelligence/)

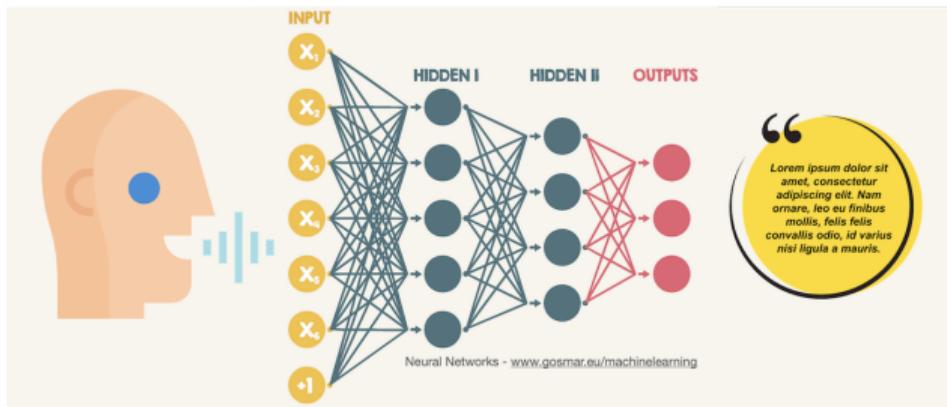
Sentiment Analysis



- Determines the sentiment (e.g., positive or negative) expressed in a text.

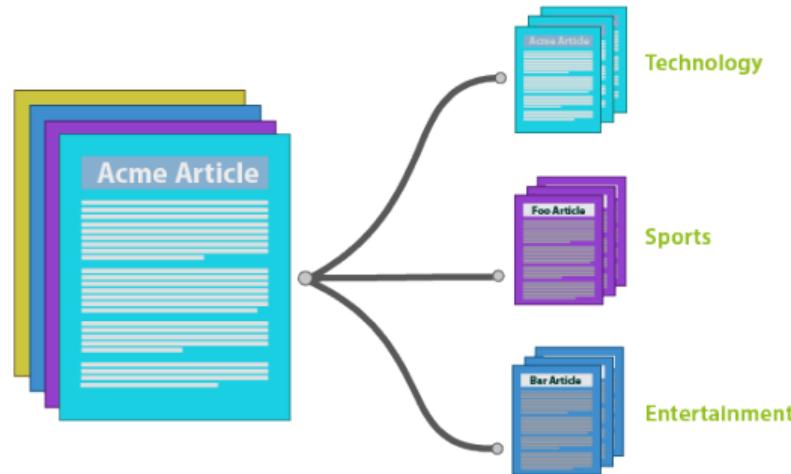
Figure adapted from www.mdpi.com/2079-9292/9/3/483

Speech Recognition



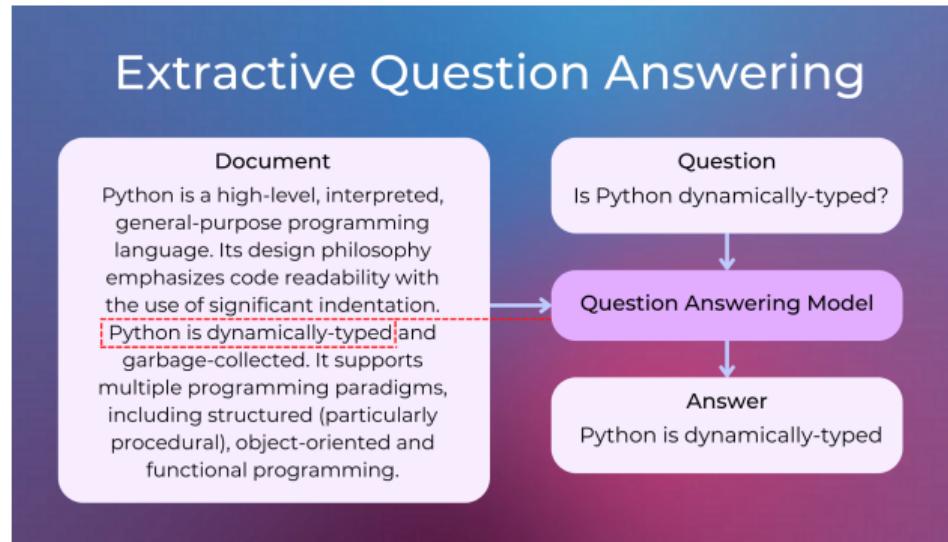
- Converts spoken language into text.

Text Classification



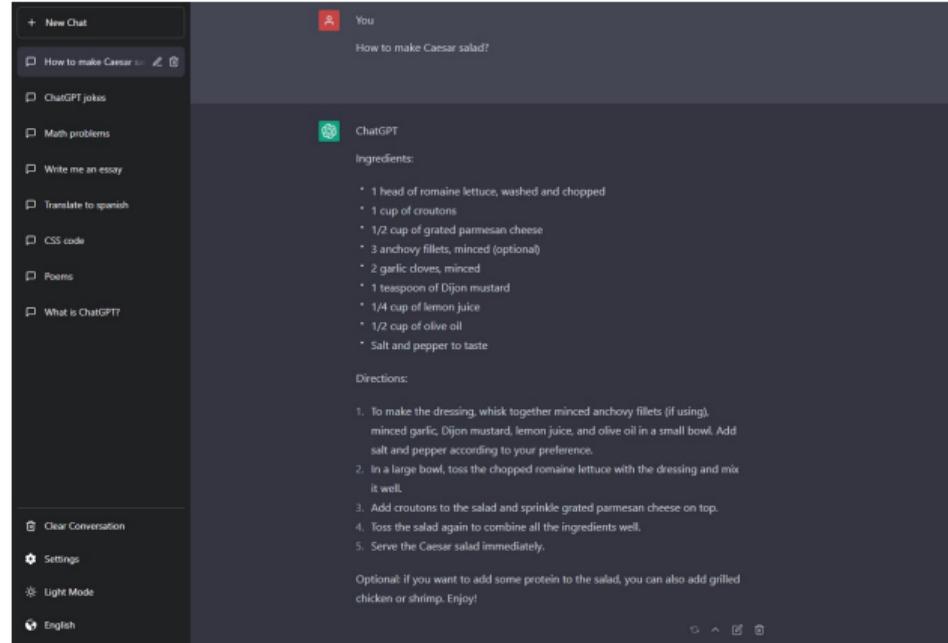
- Categorizes text into predefined groups or topics.

Question Answering



- Answers questions based on a given text or dataset.

Chatbots and Dialogue Systems



- NLP powers chatbots that can interact with users through text or speech.

Motivation and One-Hot Encoding

- Traditional models like one-hot encoding lack semantic understanding and fail to capture word relationships.
- One-hot encoding represents words as binary vectors where only one position is 1, making the representation sparse and non-semantic.
- We need dense, semantic word representations to improve natural language processing tasks.

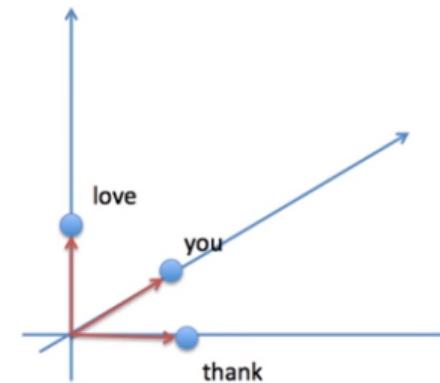
The diagram illustrates the process of generating one-hot encoding from a sentence. On the left, a text box contains the sentence: "Hello, I'm Ironman. I have Friday AI". An arrow labeled "One Hot Encoding" points from this text box to a table on the right. The table has two parts: a header row with words and a body of 8 rows, each corresponding to a word in the sentence. The header row contains: AI, Ironman, Friday, have, Hello, I, and I'm. The body of the table shows binary values (0 or 1) indicating the presence of each word in the sentence. The values are: AI (1, 0, 0, 0, 0, 0, 0), Ironman (0, 1, 0, 0, 0, 0, 0), Friday (0, 0, 1, 0, 0, 0, 0), have (0, 0, 0, 1, 0, 0, 0), Hello (0, 0, 0, 0, 1, 0, 0), I (0, 0, 0, 0, 0, 1, 0), and I'm (0, 0, 0, 0, 0, 0, 1).

AI	1	0	0	0	0	0	0
Ironman	0	1	0	0	0	0	0
Friday	0	0	1	0	0	0	0
have	0	0	0	1	0	0	0
Hello	0	0	0	0	1	0	0
I	0	0	0	0	0	1	0
I'm	0	0	0	0	0	0	1

Example: One-Hot Encoding for Different Words

	thank	you	love
thank	1	0	0
you	0	1	0
love	0	0	1

unique word	encoding
thank	[1, 0, 0]
you	[0, 1, 0]
love	[0, 0, 1]



- Different words can be perpendicular to each other in space regardless of their meanings.

Example: Word Embedding

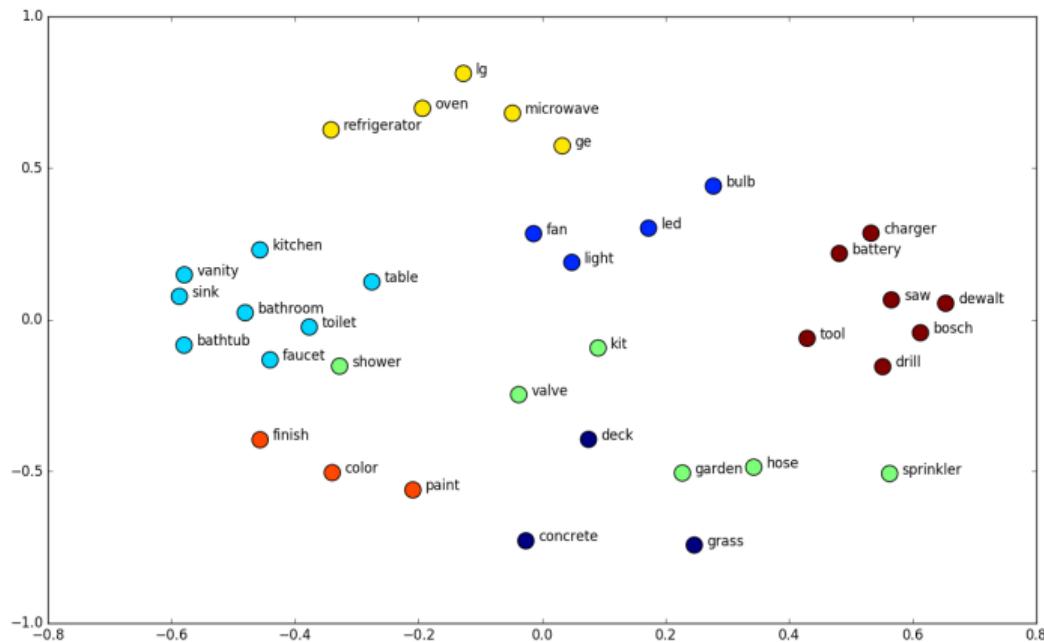


Figure adapted from classes.engr.oregonstate.edu/eecs/fall2023/ai534-400/unit4/wordembeddings.html

Skip-gram Example

Window Size	Text	Skip-grams
2	[The wide road shimmered] in the hot sun.	wide, the wide, road wide, shimmered
	The [wide road shimmered in the] hot sun.	shimmered, wide shimmered, road shimmered, in shimmered, the
	The wide road shimmered in [the hot sun].	sun, the sun, hot
3	[The wide road shimmered in] the hot sun.	wide, the wide, road wide, shimmered wide, in
	[The wide road shimmered in the hot] sun.	shimmered, the shimmered, wide shimmered, road shimmered, in shimmered, the shimmered, hot
	The wide road shimmered [in the hot sun].	sun, in sun, the sun, hot

Different window sizes and samples drawn from context words and their target

Word2Vec as a Neural Network

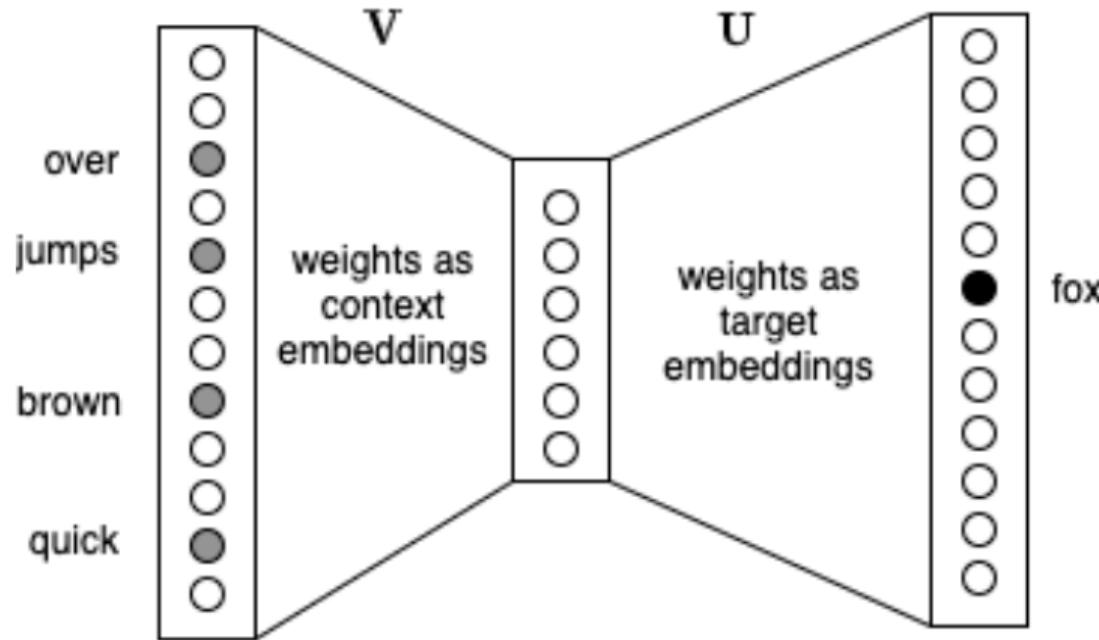
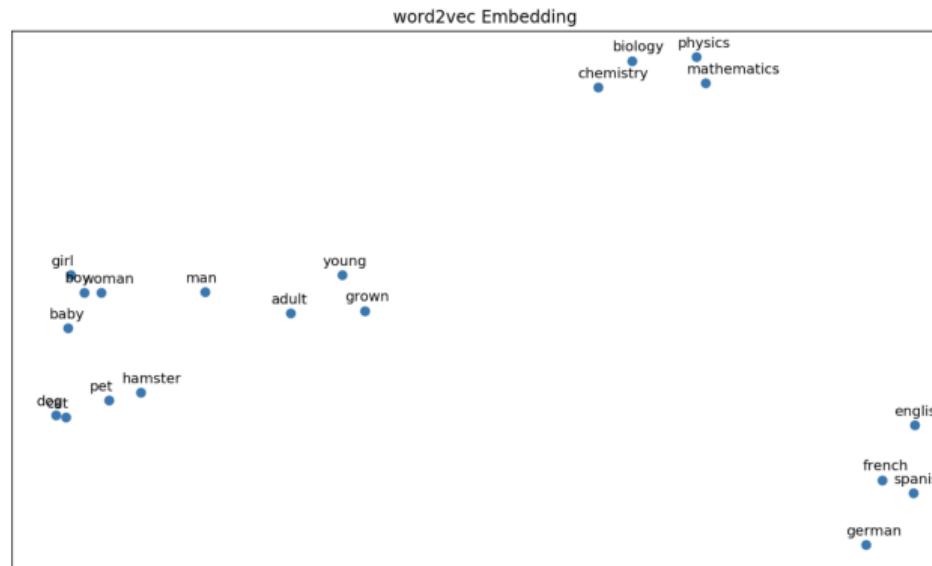


Figure adapted from machinelearningcaban.com/tablbook/chembedding/word2vec.html

Visualizing Words in 2D



Words represented in a 2D space after dimensionality reduction

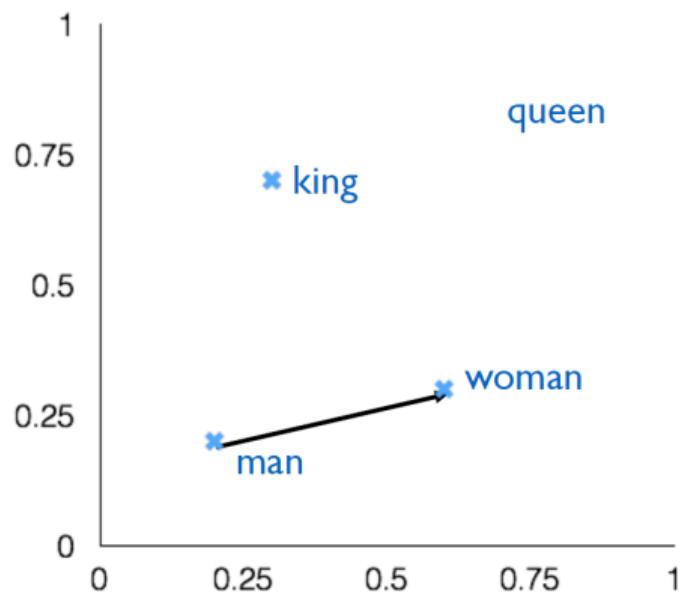
Word Analogy: Vector Arithmetic in Word2Vec

- Word2Vec embeddings can solve analogy tasks by performing vector arithmetic.
- The analogy task takes the form:

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

- The analogy is solved by finding the word vector closest to $\mathbf{v}_{\text{king}} - \mathbf{v}_{\text{man}} + \mathbf{v}_{\text{woman}}$.

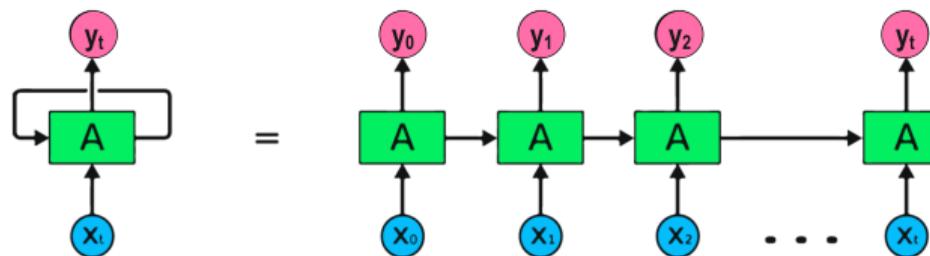
Word Analogy Example



Word analogy example

Recurrent Neural Networks

- To process each word, we need to remember the previous words.
- How do we create this memory? By maintaining a **hidden state**.
- Each hidden state captures information about:
 - Current word
 - All previous words in the sequence



where h_t serves as contextualized representation, containing memory of previous words.

RNN's Hidden State Update

- Process sequences by maintaining a hidden state.
 - Update state sequentially: $h_t = f_W(h_{t-1}, x_t)$
 - This recurrence formula is applied at each time step to process a sequence of vectors x .
 - The same function and the same set of parameters are used for each word.

$$h_t = f_W(h_{t-1}, x_t)$$

new state / old state input vector at
 some function some time step
 with parameters W

Contextualized Word Embeddings

Compute contextual vector:

$$\mathbf{c}_k = f(w_k \mid w_1, w_2, \dots, w_n) \in \mathbb{R}^d$$

Examples:



$f(\text{light} | \text{Please turn off the } \overbrace{\text{light}}^{\text{light}})$

#

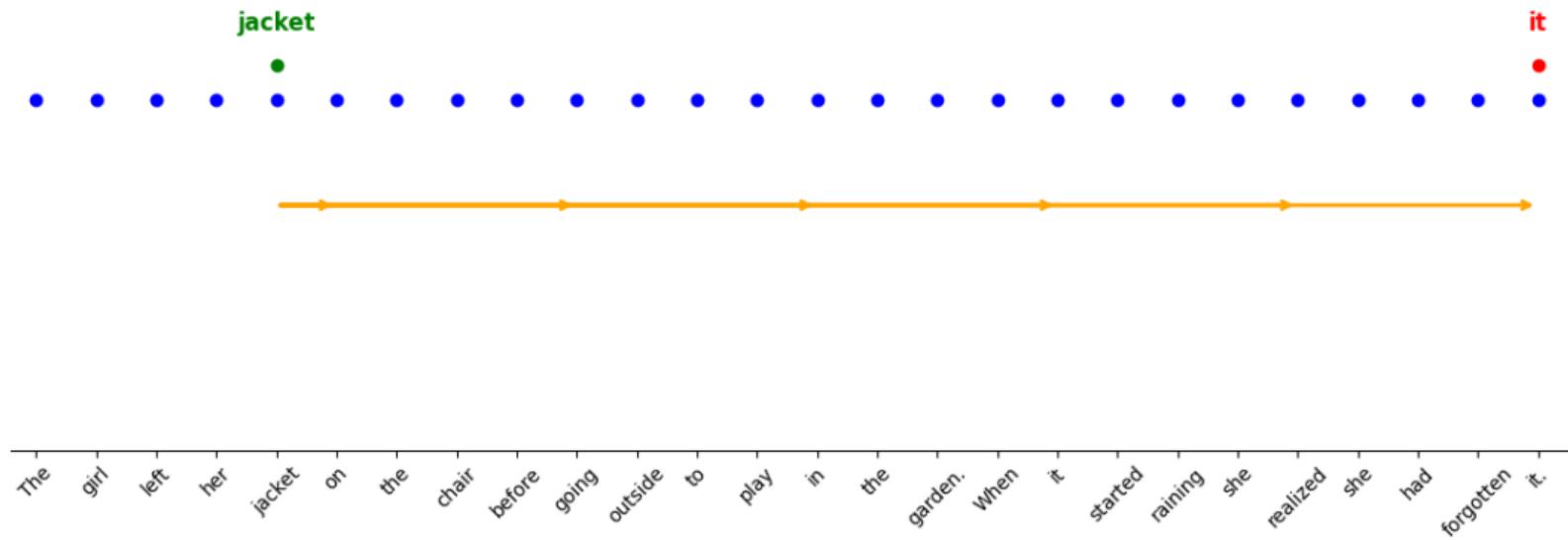
$f(\text{light} | \text{This box is very light to carry.})$



How do we implement the context function f ?

Challenges of RNNs

RNNs face challenges in retaining information over long sequences, such as when pronouns refer to distant words in the sequence.



So, what solution do we have?

Why Attention?

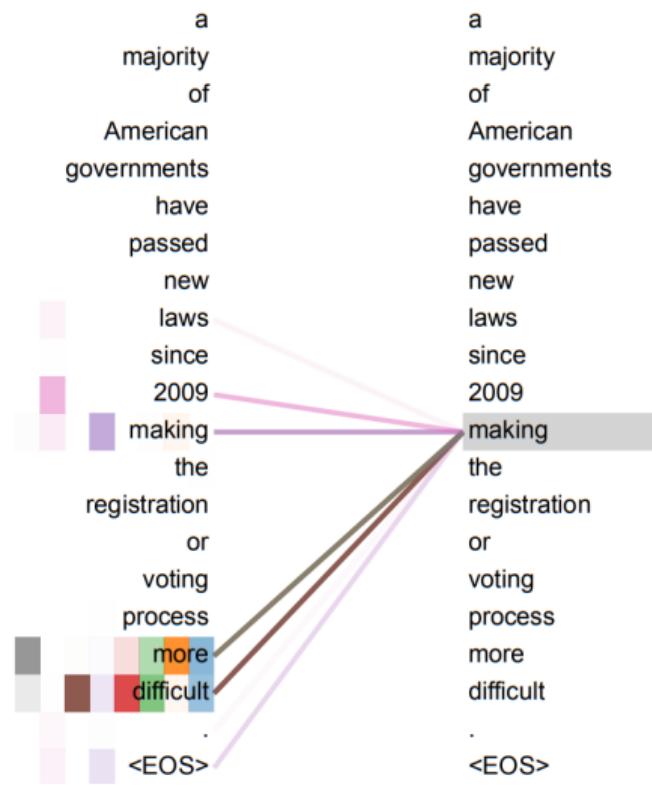
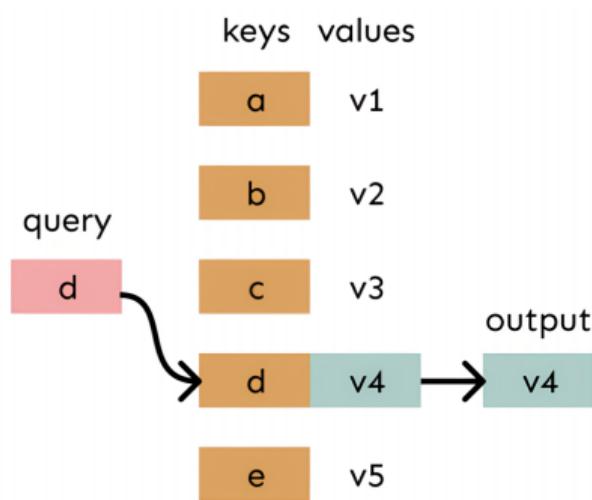
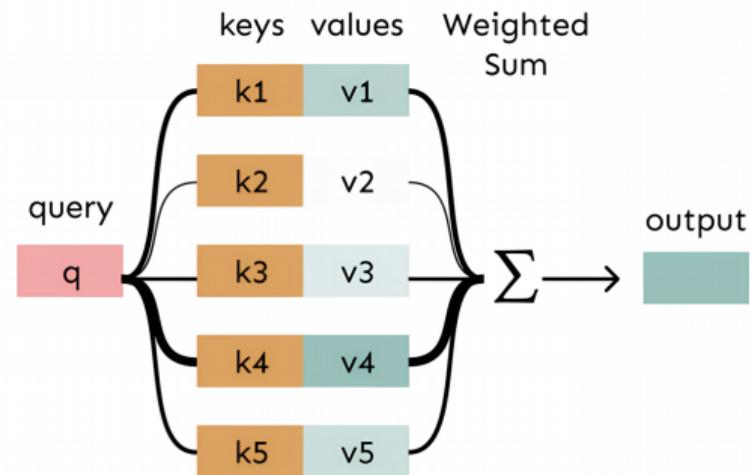


Figure adapted from Ashish Vaswani et al. Attention is All You Need paper

Attention as a Soft, Averaging Lookup Table



In a **lookup table**, we have a table of keys that map to values. The query matches one of the keys, returning its value.

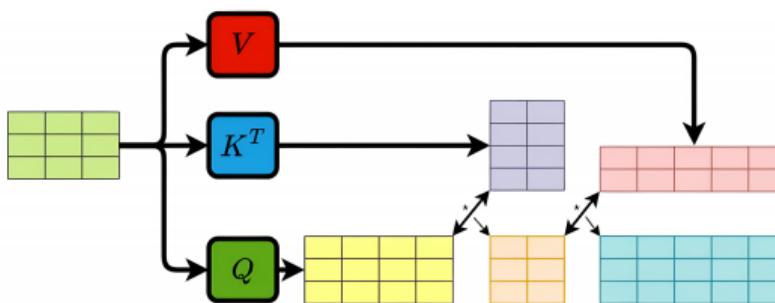


In the **attention**, the query matches all keys softly, to a weight between 0 and 1. The keys' values are multiplied and summed by the weights.

Self-Attention vs Cross-Attention

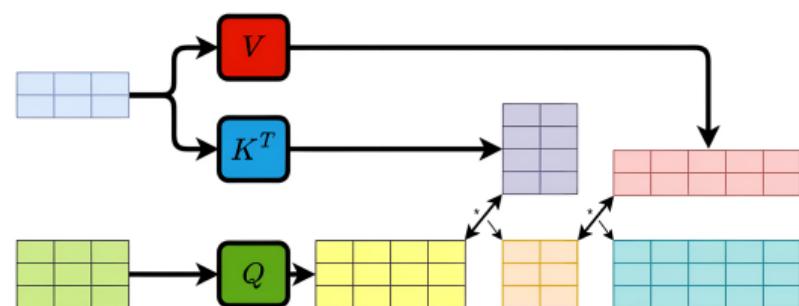
Self-Attention

- Queries, Keys, and Values are derived from the same sequence.
- Each position attends to all other positions in the sequence.

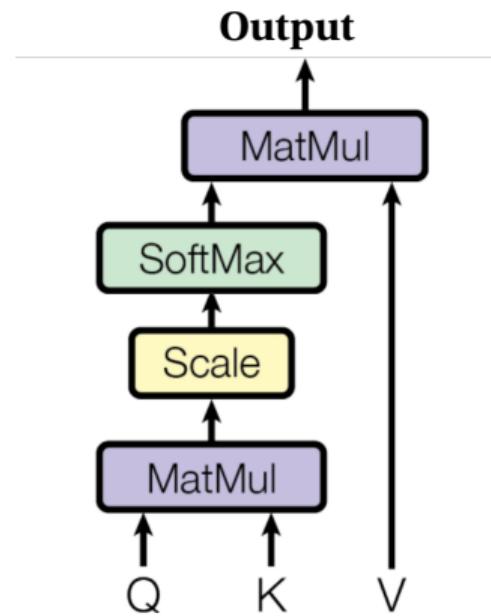


Cross-Attention

- Queries come from one sequence.
- Keys and Values come from another sequence.



Scaled Dot-Product Attention

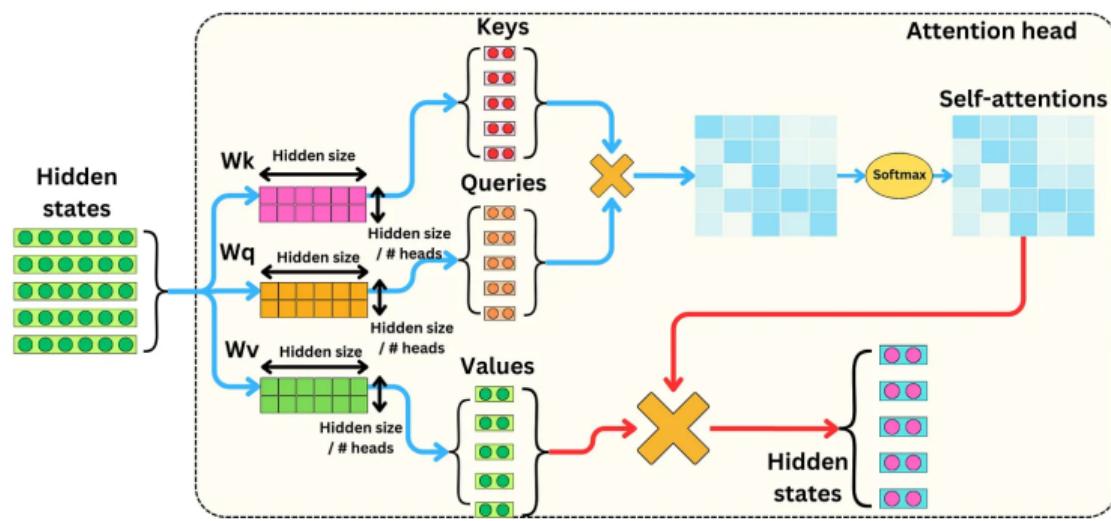


$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Figure adapted from Ashish Vaswani et al. Attention is All You Need paper

Multi-Head Attention

$$\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$



Multi-Head Attention

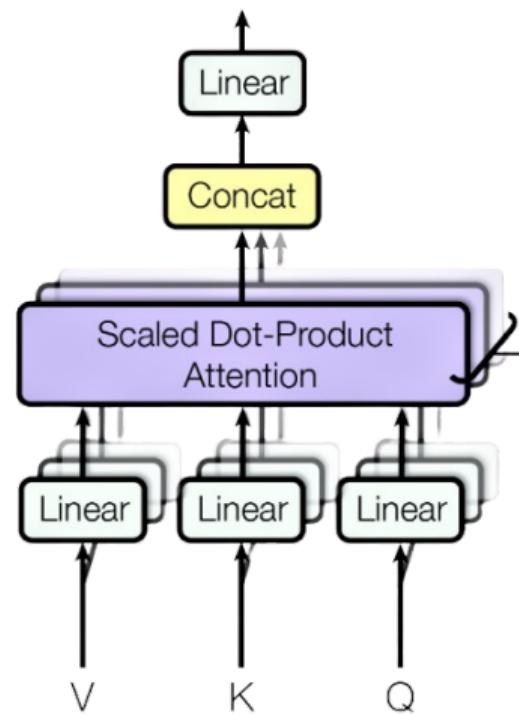
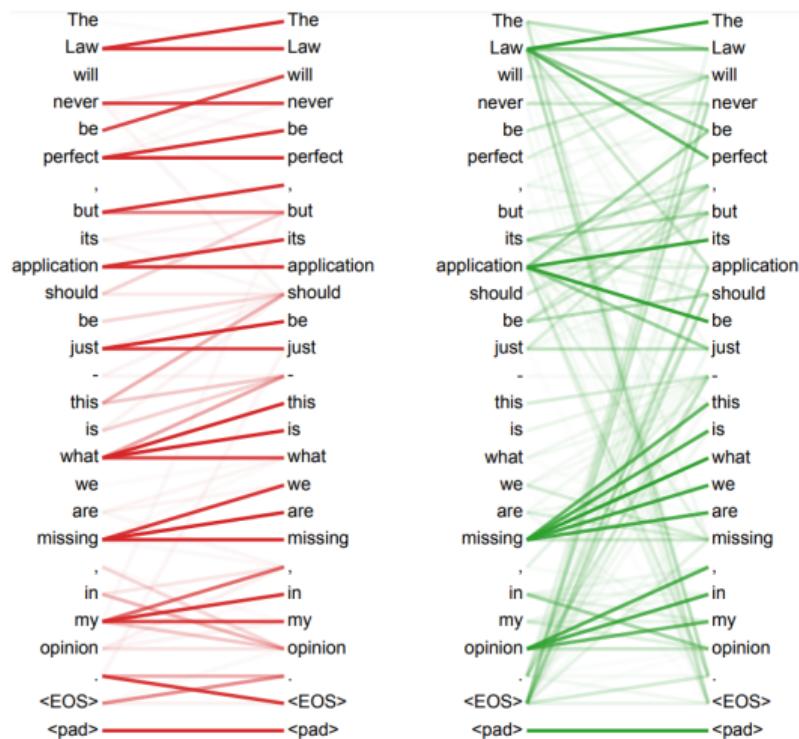


Figure adapted from Ashish Vaswani et al. Attention is All You Need paper | Kucedra : The 7 headed dragon, MythLok

Ali Sharifi-Zarchi (Sharif University of Technology)

Machine Learning (CE 40477)

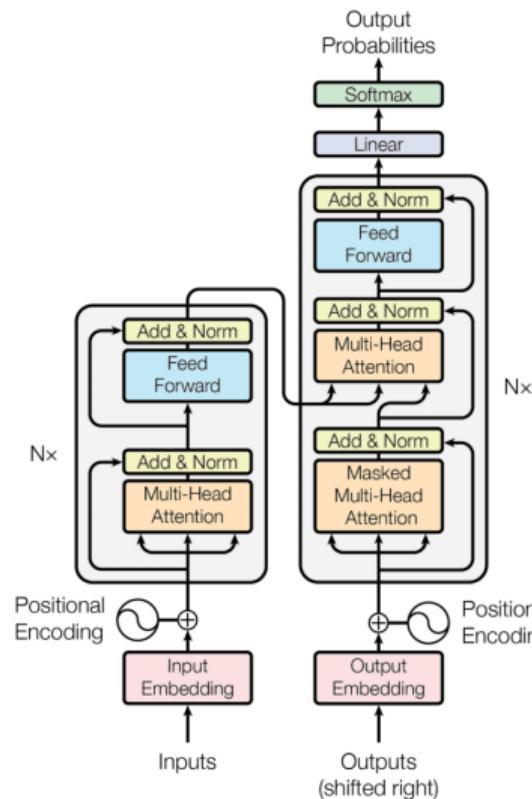
Multi-Head Attention



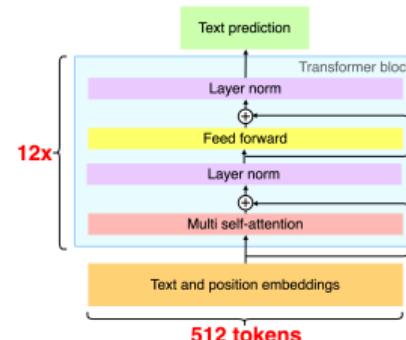
- Different attention heads seem to focus on different parts of the sentence.
- We show two examples from different heads.
- Each head appears to have learned a different task.

Figure adapted from Ashish Vaswani et al. Attention is All You Need paper

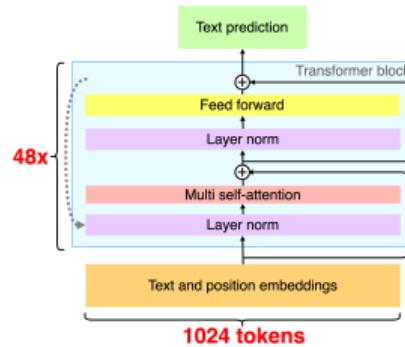
Transformers: A General Architecture



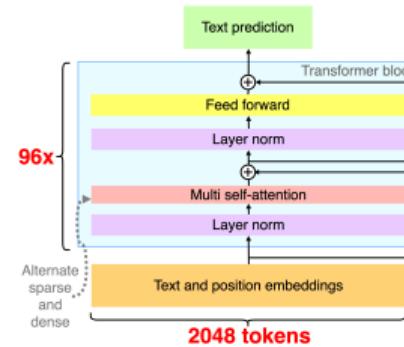
GPTs

GPT-1

512 tokens

GPT-2

1024 tokens

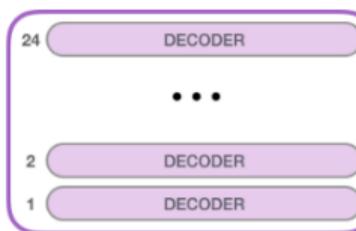
GPT-3

2048 tokens

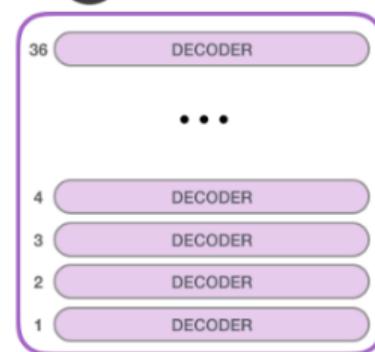
GPT-2

GPT-2
SMALL

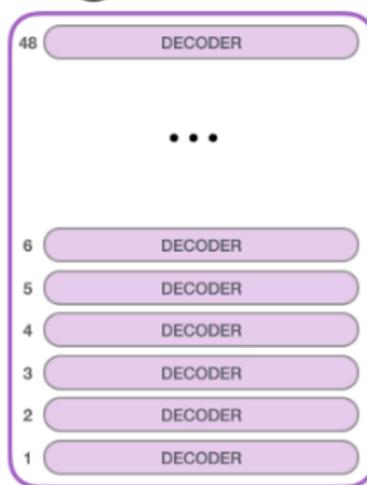
Model Dimensionality: 768

GPT-2
MEDIUM

Model Dimensionality: 1024

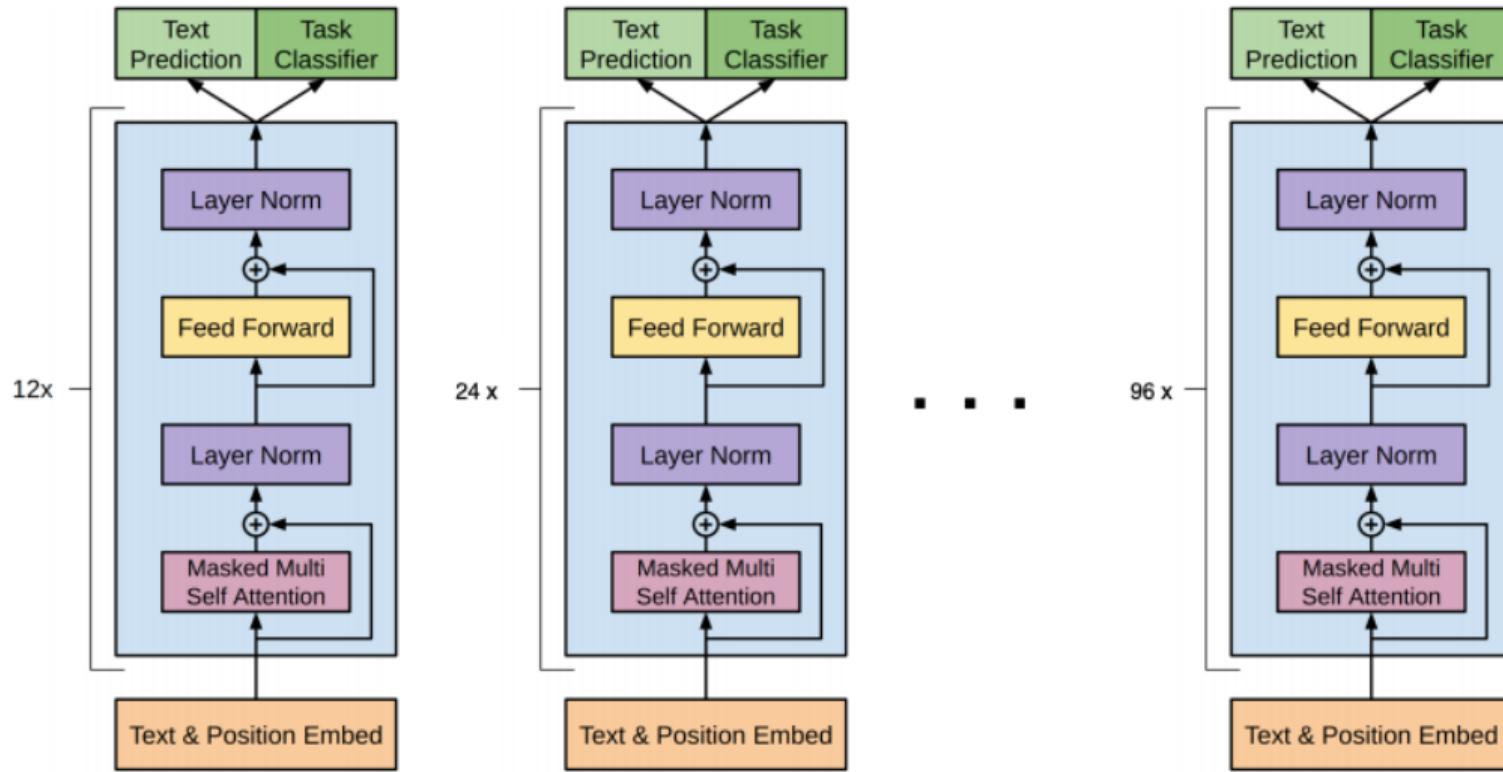
GPT-2
LARGE

Model Dimensionality: 1280

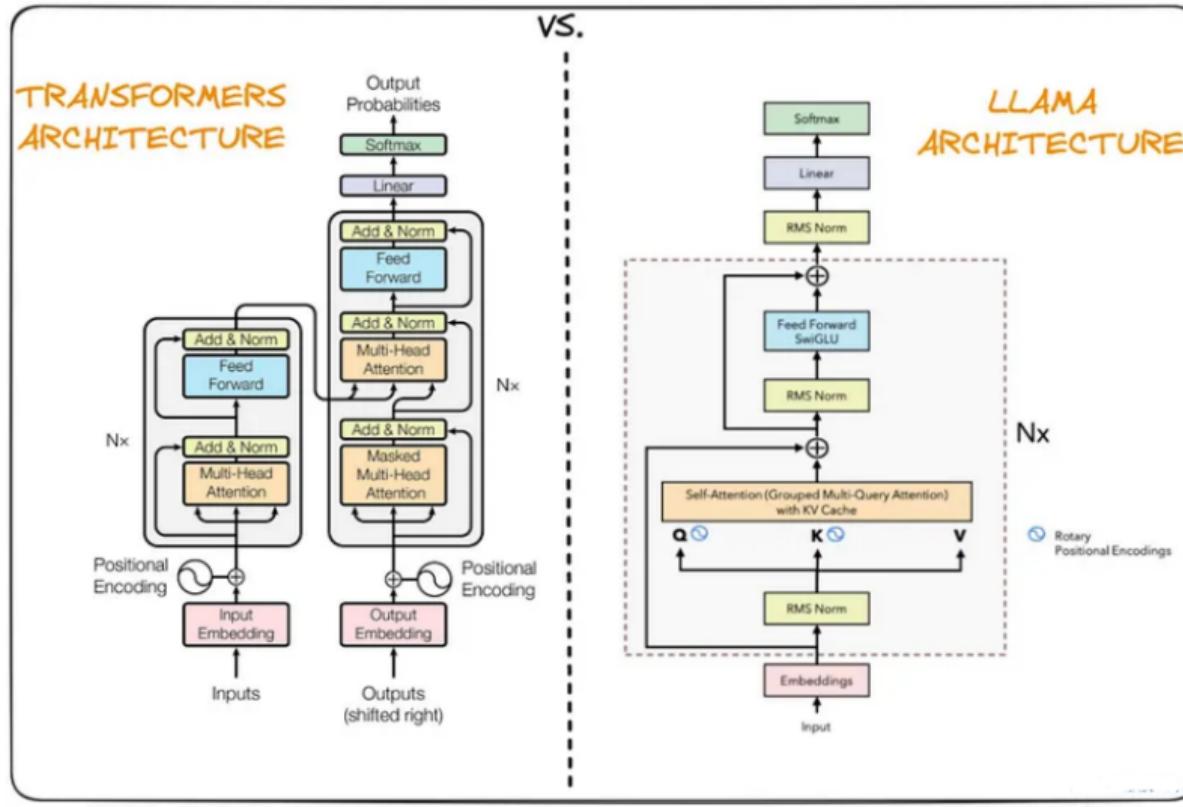
GPT-2
EXTRA
LARGE

Model Dimensionality: 1600

GPT-3



Llama 2 Architecture



Pre-training and adaptation



Figure 1: Overview of Model Training: Pre-training on large, unlabeled data builds foundational knowledge, while fine-tuning on smaller, labeled datasets adapts the model for specific tasks.

Adapters

- **Adapters:** New modules inserted between layers of a pre-trained model, with the original model weights fixed.
- **Training Efficiency:** Only adapter modules are tuned, initialized to ensure their output resembles that of the original model.

