



# AI × ECON

## Summer School @ CUHK

AUGUST 7-8 2024

# Transformer

---

*Almost as if it's intelligent.*

# Conversational AI



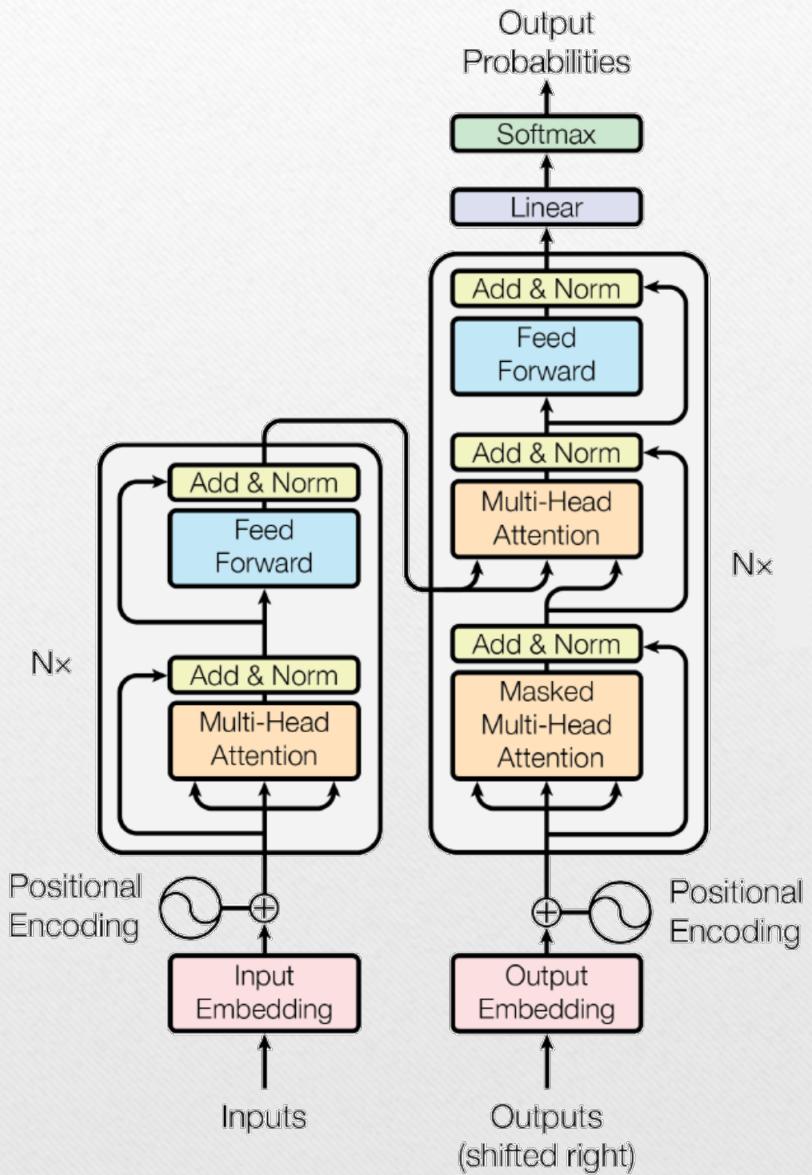
# Image Generation



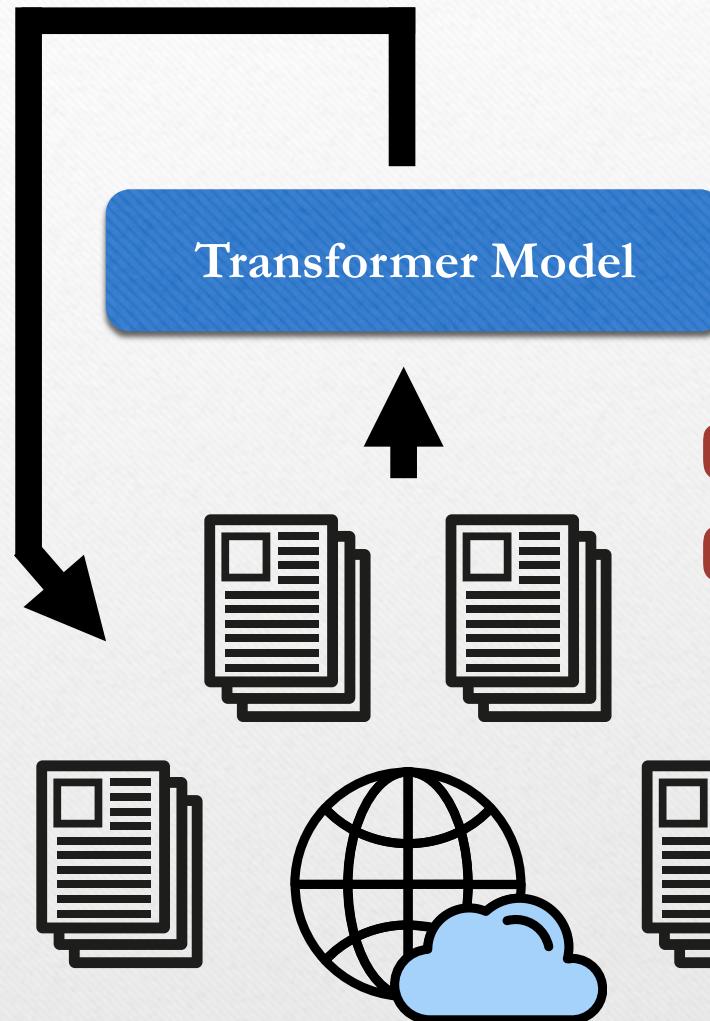
A composite of three DALL-E 3 AI art generations: an oil painting of Hercules fighting a shark, a photo of the queen of the universe, and a marketing photo of "Marshmallow Menace" cereal. *Benj Edwards@Ars Technica*

# Transformer

- Transformer-based models such as GPT are the current state of the art in field of natural language processing.
- Original paper:  
[All You Need is Attention](#)  
(NeurIPS 2017)
- Transformer is the basis of some of the most advanced AI currently in existence, including the Starcraft-playing [AlphaStar](#), text-generating [GPT-4](#) and image-generating [Stable Diffusion](#).



**Unsupervised:**  
No manual labelling



**computationally  
expensive!**

Stage 1: Pre-training

約翰 吃 了 蘋果。

↑ ↑ ↑ ↑  
762 45 3 442

tokens



Transformer Model



135 27 14 267  
↑ ↑ ↑ ↑  
John ate an apple

tokens

# Input and Output

Transformer is a **sequence-to-sequence** model,  
meaning that the input and output are both vectors.

約翰 吃 了 蘋果。

↑ ↑ ↑ ↑  
762 45 3 442



Transformer Model



135 27 14 267  
↑ ↑ ↑ ↑  
John ate an apple

# Output

**Translation:** the model predicts words in another language. This was what the original paper did.

ate an apple and  
[ 27 , 14 , 267, 8 ]



Transformer Model

[ 135, 27, 14 , 267 ]  
John ate an apple

# Output

**The next word:** the model predicts the next word after the input. This is how GPT models are pre-trained.

**Unsupervised:**  
No manual labelling

ate apple  
[ 27 , 267 ]

### Transformer Model

[ 135, 1 , 14 , 1 ]  
John [X] an [X]  
John ate an apple

**Unsupervised:**  
No manual labelling

# Output

**Fill-in-the-blanks/denoising:** the model predicts missing words in a passage from which the input was taken. This is how T5 and BERT models are pre-trained.

# Stage 2: Fine-tuning

The model is then **fine-tuned** to a specific task. This stage can be done with relatively little data and computational power.

With fine-tuning, the model can adapt to tasks that it was not pre-trained for:

- Instruction following
- Sentiment analysis
- Text classification

Fine-tuning typically only takes a few epochs to achieve good results, because the model already understand language from pre-training.

Positive      Negative

[ 0.7 , 0.3 ]



Transformer Model



[ 135, 1 , 14 , 1 ]

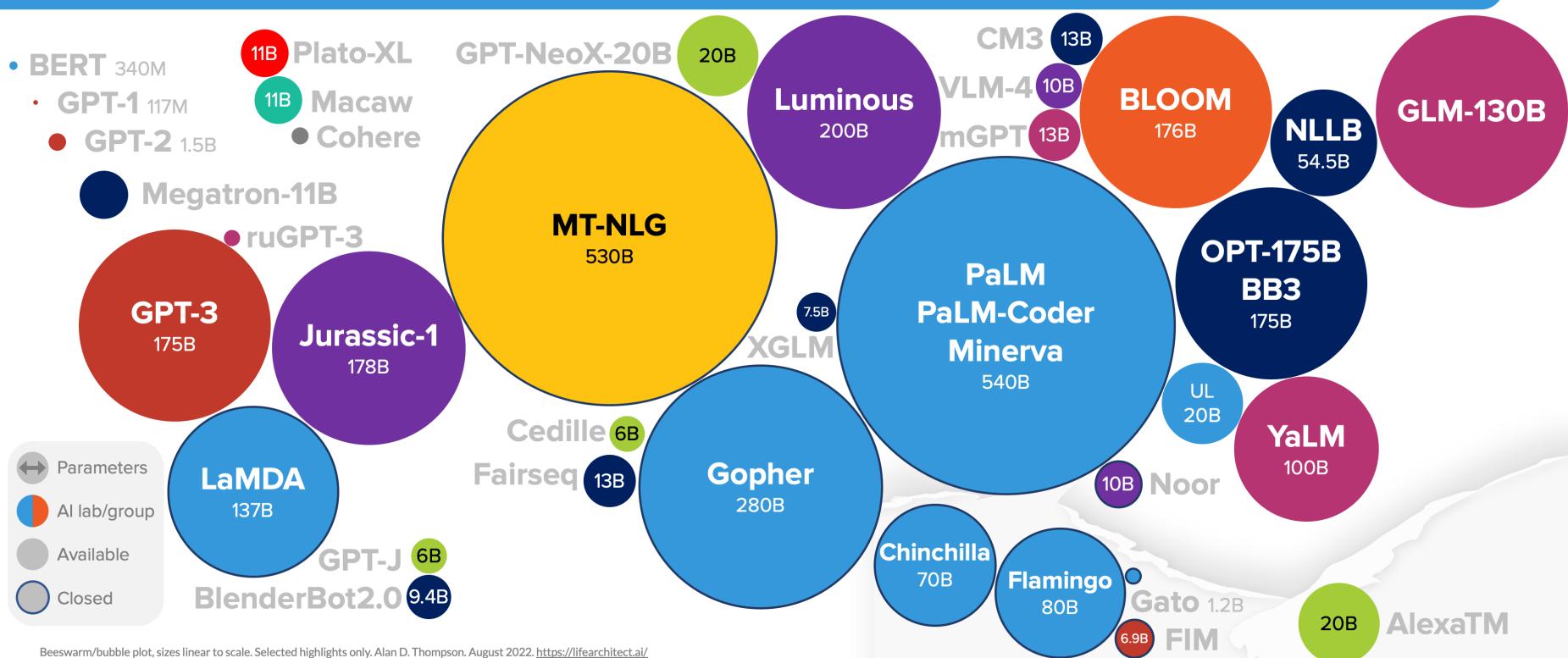
↑      ↑      ↑      ↑  
John [X] an [X]



John ate an apple

# 2022 Big

## LANGUAGE MODEL SIZES TO AUG/2022



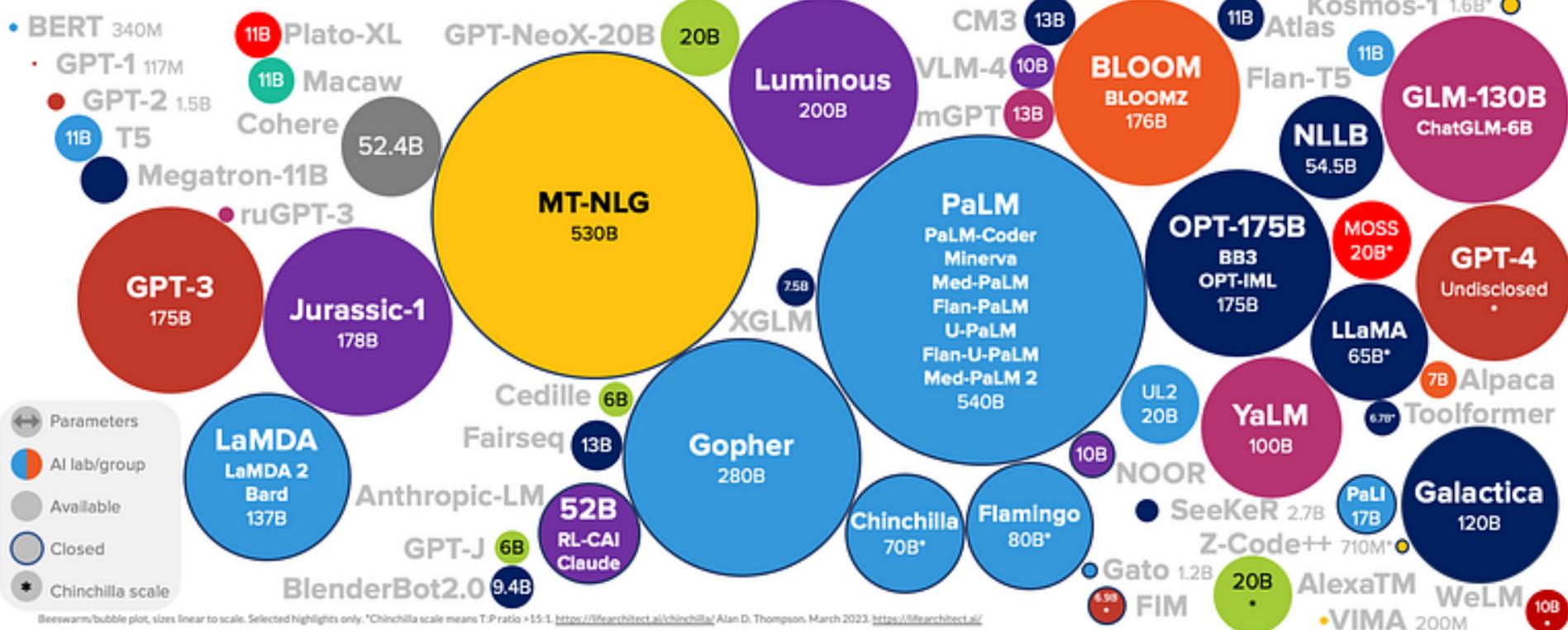
Beeswarm/bubble plot, sizes linear to scale. Selected highlights only. Alan D. Thompson. August 2022. <https://lifearchitect.ai/>



LifeArchitect.ai/models

# 2023 Big

## LANGUAGE MODEL SIZES TO MAR/2023



# Model Categorization

- Transformer-based models, particularly LLMs, are usually categorized by their architecture and the number of parameters they have.
  - E.g. LLaMA 7B means the model uses the LLaMA architecture and has 7 billion parameters.
- Models with more parameters are more powerful, but also more expensive to run.
  - Difference is small for tasks such as text classification, but huge for math complex reasoning.

## Meta Llama 3 Instruct model performance

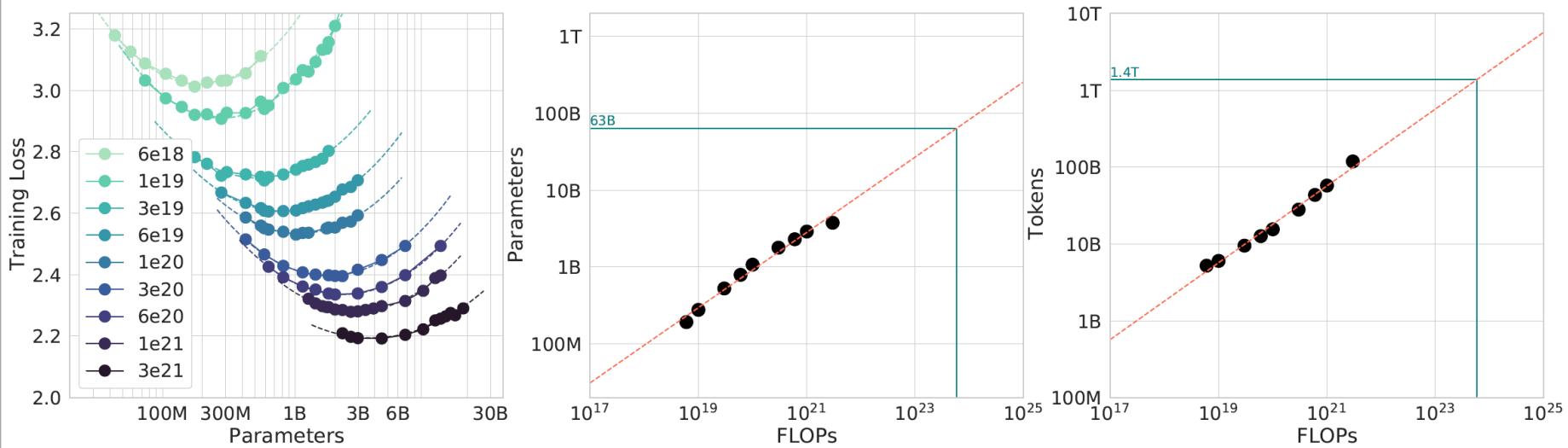
	Meta Llama 3 8B	Gemma 7B - It Measured	Mistral 7B Instruct Measured	Meta Llama 3 70B	Gemini Pro 1.5 Published	Claude 3 Sonnet Published
MMLU 5-shot	<b>68.4</b>	53.3	58.4	<b>82.0</b>	81.9	79.0
GPQA 0-shot	<b>34.2</b>	21.4	26.3	39.5	<b>41.5</b> CoT	<b>38.5</b> CoT
HumanEval 0-shot	<b>62.2</b>	30.5	36.6	<b>81.7</b>	71.9	73.0
GSM-8K 8-shot, CoT	<b>79.6</b>	30.6	39.9	<b>93.0</b>	91.7 11-shot	<b>92.3</b> 0-shot
MATH 4-shot, CoT	<b>30.0</b>	12.2	11.0	50.4	<b>58.5</b> Minerva prompt	40.5

Source: [Meta](#)

# Why Bigger Model?

Large language models are likely undertrained. In other words, the amount of data and computational power spent on training is insufficient relative to model size.

...but we have basically run out of data, so training larger models for longer it is.



# How to Use Transformer-based models?

- Training large Transformer-based models from scratch is very expensive, so in practice you will need to use pretrained models.
- **Closed-source models:** OpenAI and Anthropic provide an API to their models, with different payment tiers. These models are most suitable if your task can be framed as questions. Inaccessible in Hong Kong without VPN.
- **Open-weight models:** Pretrained models available for download. The easiest service to use are [Hugging Face](#) and [GPT4All](#), with the most popular models being Meta's [LLaMa](#), MosaicML's [MPT](#), TII's [Falcon](#) and [BERT](#). BERT is particularly suitable for fine-tuning for non-conversational task such as text classification, due to its relatively small size and encoder-only architecture.

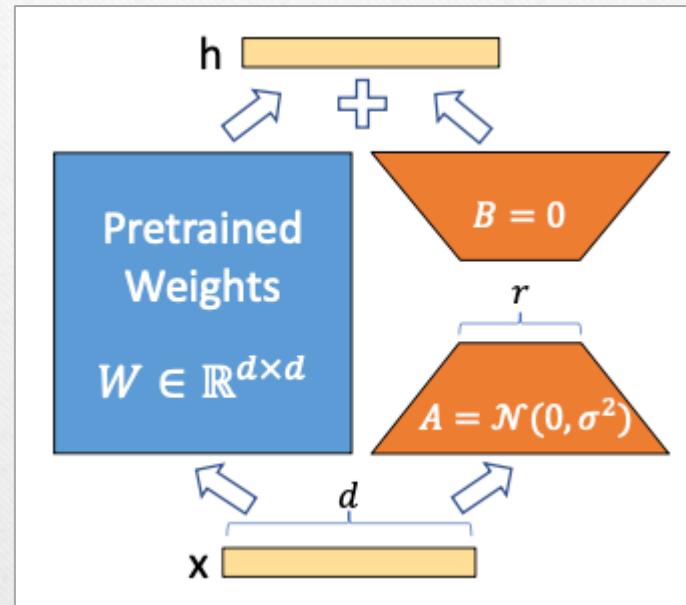
# Memory Requirement for Open-Weight Models

- Models are generally saved in 16-bit numeric format, which means each parameter takes 2 bytes of memory.
  - LLaMA 7B would take approximately 14GB of memory.
- If the model requires more memory than your system—particularly your GPU—can provide, you can lower the numeric precision. This is called **quantization** and will result in lower model performance.
  - LLaMA 7B takes 7GB of memory under 8-bit and 3.5GB under 4-bit.
- Note that the above is just the memory requirement of loading the model. Training the model could require multiples of that, due to the need to save the gradient for every weight and the activation for every neuron for back propagation.

# LoRA: Low-Rank Adapter

LoRA deals with the memory problem by only training a low-rank matrix that is then added to the original LLM. The original LLM can be loaded in lower precision and without tracking gradients.

Pre-trained LoRA addons to popular open-weight models are commonly available.



$$W = S \times A$$

# How Does Transformer Work?

---

*Not necessary if you just  
want to use pre-trained models!*

# Why Transformer?

John ate an apple. It was delicious.

John ate an apple. He likes it.

- Consider the following above passages.

# Self-Attention

John ate an apple. **It** was delicious.

John ate an apple. **He** likes it.

- The most important idea of a Transformer is **self-attention**.
- When the model reads a particular word, the output depends on *all* words, weighed by how important they are given the word in concern.
- We will now go through what a **self-attention layer** does.

John ate an apple. It was delicious.

This is one input sample.

[start] John ate an apple. It was delicious. [end]

During preprocessing, special characters are added to the input and target text phrases. These characters denote the start, the end and other contextual information and are treated just like another any other word.

[start] John ate an apple. It was delicious. [end]

For space reasons I will omit them in latter slides.

$$\begin{matrix} \vec{e}_{john} & \vec{e}_{ate} & \vec{e}_{an} & \vec{e}_{apple} & \vec{e}_{it} & \vec{e}_{was} & \vec{e}_{delicious} \\ \uparrow & \uparrow & \uparrow & \uparrow & \uparrow & \uparrow & \uparrow \\ \text{John} & \text{ate} & \text{an} & \text{apple.} & \text{It} & \text{was} & \text{delicious.} \end{matrix}$$

As in all NLP neural network models, we start with assigning each word an embedding—a vector learnt through training to represent the nature of each word.

e.g. Noun Verb Male

$$\vec{e}_{john} = [1 \quad 0 \quad 1 \quad \dots]$$

$$\vec{e}_{ate} = [0 \quad 1 \quad 0 \quad \dots]$$

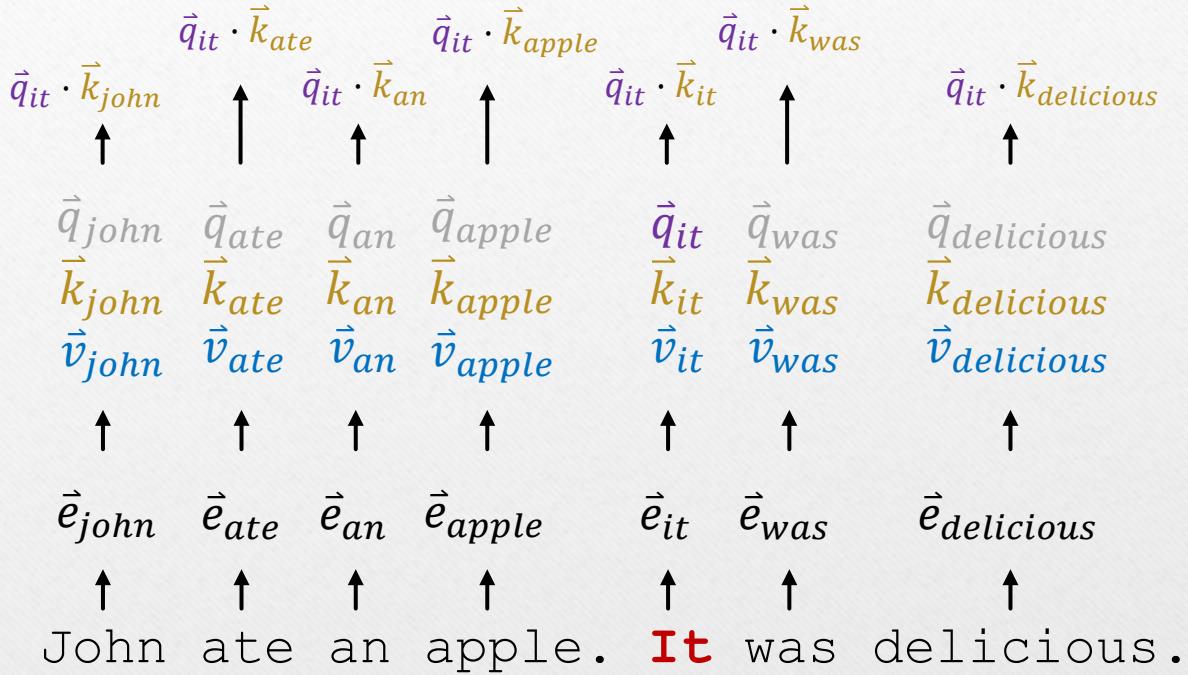
$$\vec{e}_{he} = [0 \quad 0 \quad 1 \quad \dots]$$

$$\begin{array}{ccccccc} \vec{e}_{john} & \vec{e}_{ate} & \vec{e}_{an} & \vec{e}_{apple} & \vec{e}_{it} & \vec{e}_{was} & \vec{e}_{delicious} \\ \uparrow & \uparrow & \uparrow & \uparrow & \uparrow & \uparrow & \uparrow \\ \text{John} & \text{ate} & \text{an} & \text{apple.} & \text{It} & \text{was} & \text{delicious.} \end{array}$$

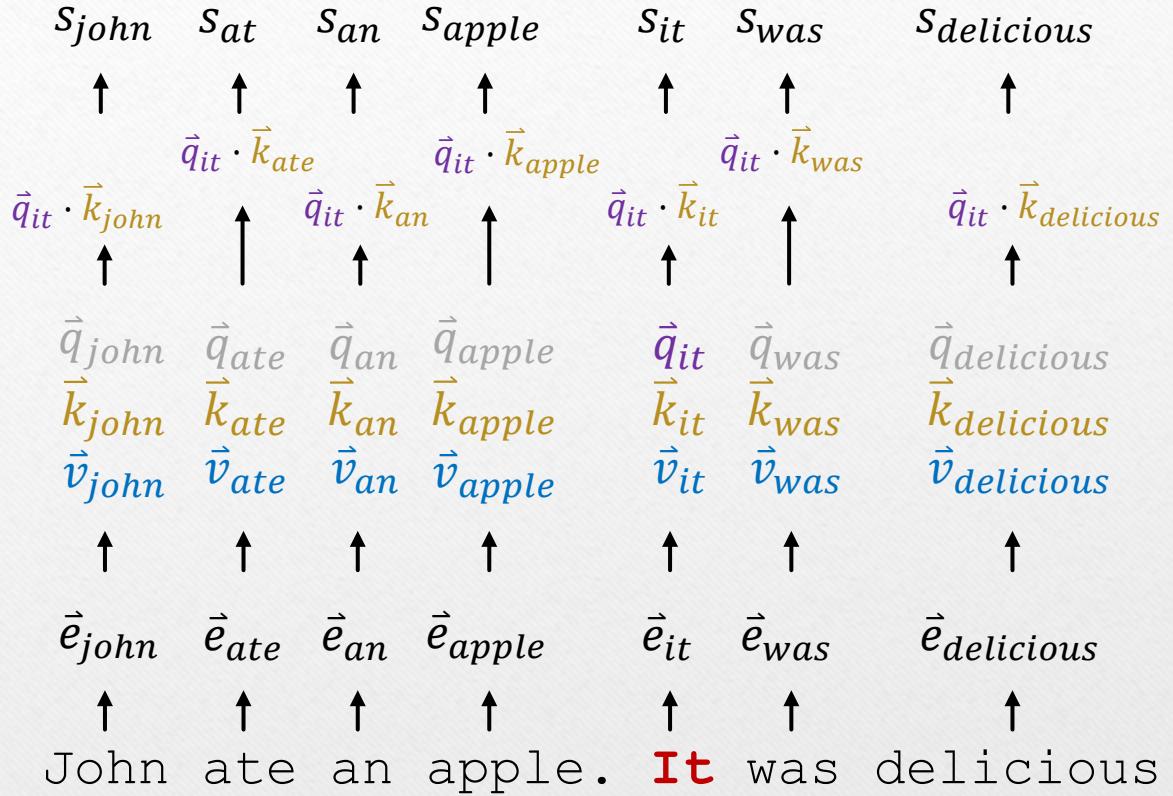
As in all NLP neural network models, we start with assigning each word an embedding—a vector learnt through training to represent the nature of each word.

$\vec{q}_{john}$	$\vec{q}_{ate}$	$\vec{q}_{an}$	$\vec{q}_{apple}$	$\vec{q}_{it}$	$\vec{q}_{was}$	$\vec{q}_{delicious}$
$\vec{k}_{john}$	$\vec{k}_{ate}$	$\vec{k}_{an}$	$\vec{k}_{apple}$	$\vec{k}_{it}$	$\vec{k}_{was}$	$\vec{k}_{delicious}$
$\vec{v}_{john}$	$\vec{v}_{ate}$	$\vec{v}_{an}$	$\vec{v}_{apple}$	$\vec{v}_{it}$	$\vec{v}_{was}$	$\vec{v}_{delicious}$
↑	↑	↑	↑	↑	↑	↑
$\vec{e}_{john}$	$\vec{e}_{ate}$	$\vec{e}_{an}$	$\vec{e}_{apple}$	$\vec{e}_{it}$	$\vec{e}_{was}$	$\vec{e}_{delicious}$
↑	↑	↑	↑	↑	↑	↑
John	ate	an	apple.	<b>It</b>	was	delicious.

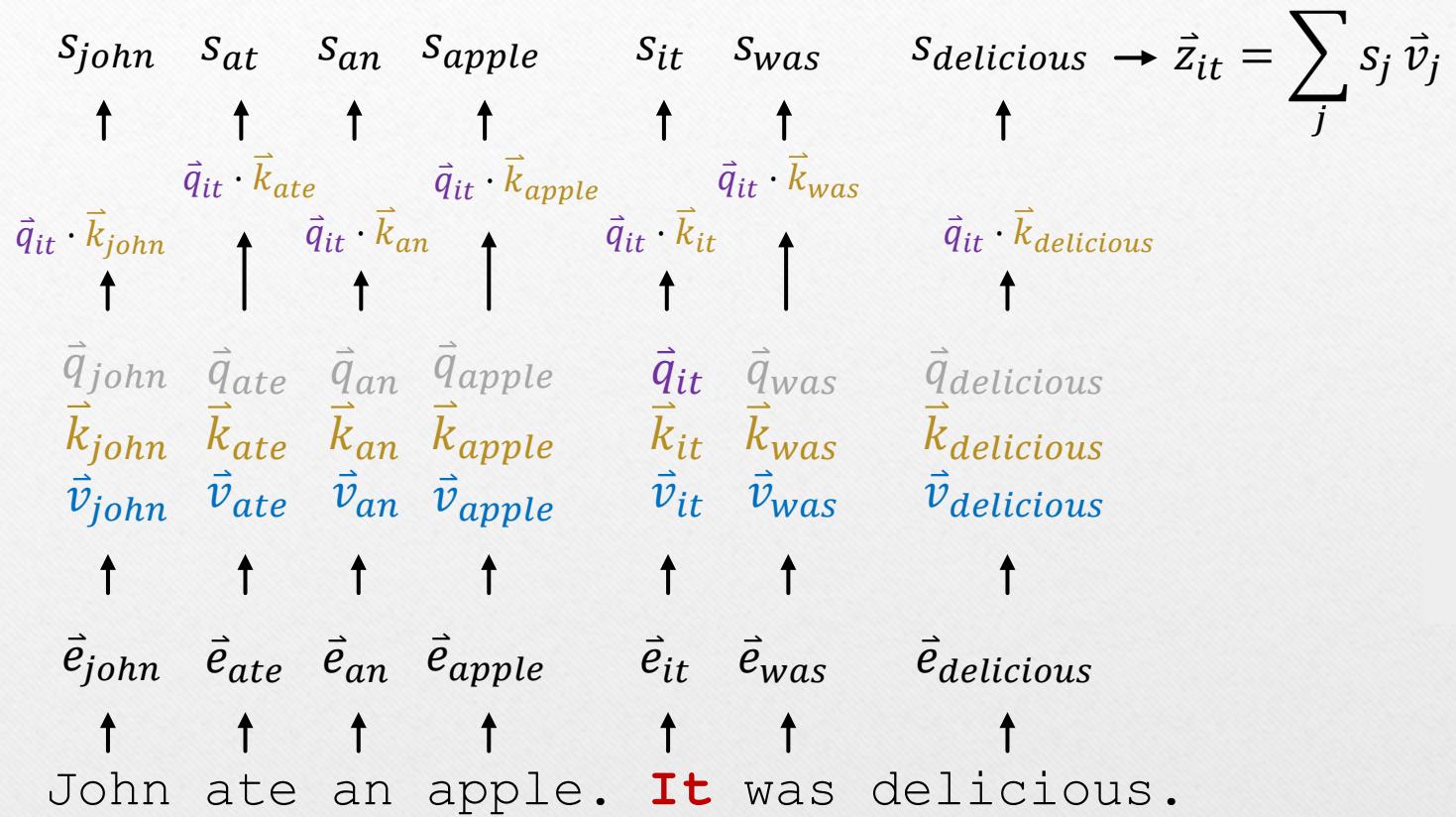
The attention mechanism requires three vectors as inputs: **query**, **key** and **value**. In the simplest case, they are simply the output from the previous layer.



When we process a word  $i$ , we compute a **score** for *all* words by taking the dot product of  $i$ 's query and each word's key.



A weight is assigned to each word by feeding the score to a softmax function, i.e. multinomial logit.



The output  $z_i$  of the attention mechanism for word  $i$  is the sum of the multiple of each word's weight with its value.

# Why Does Self-Attention Make Sense?

- Reconsider our example:

$$\vec{h}_{john} = [1 \quad 0 \quad 1]$$

$$\vec{h}_{ate} = [0 \quad 1 \quad 0]$$

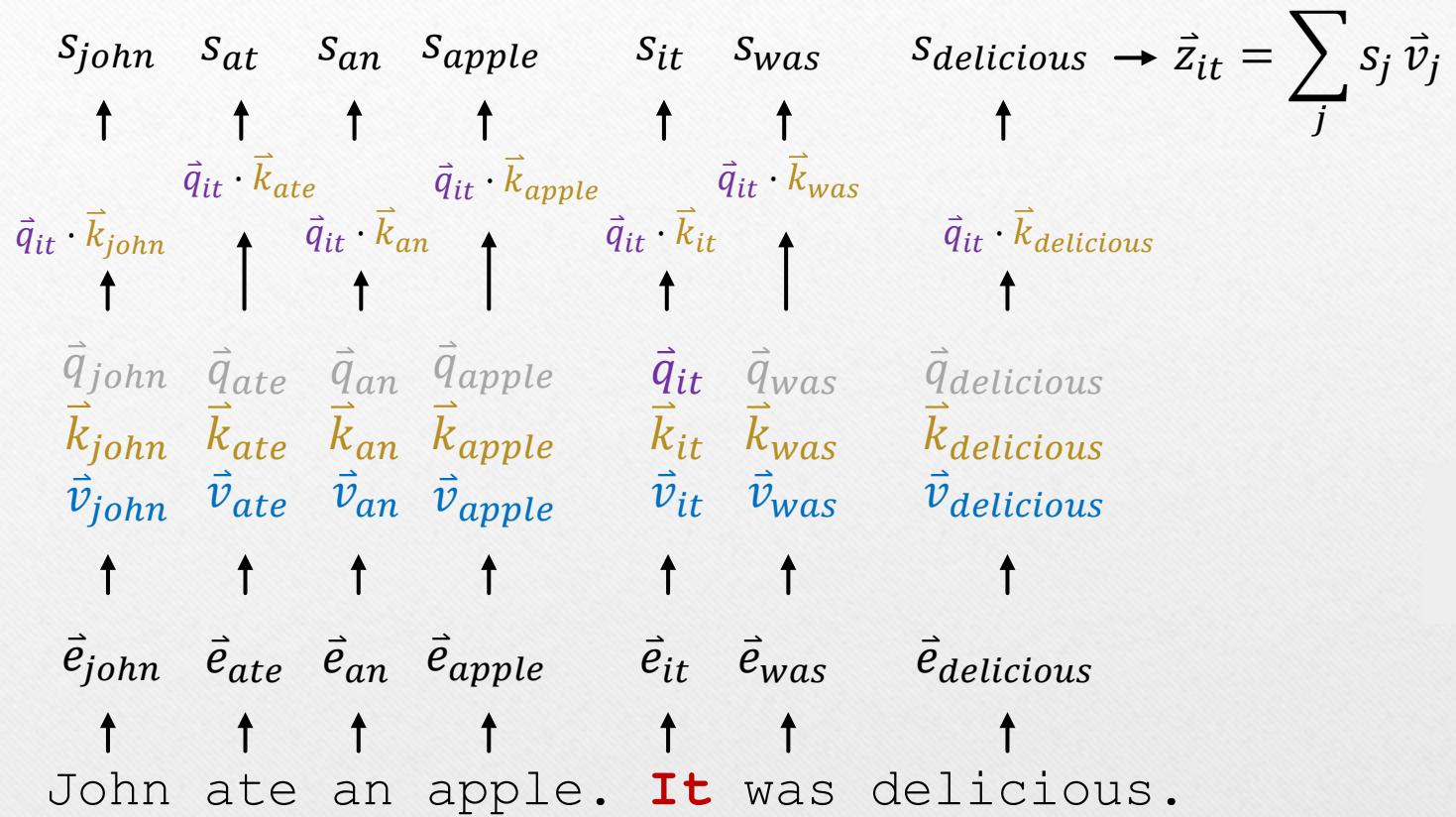
$$\vec{h}_{he} = [0 \quad 0 \quad 1]$$

- When processing the word “he”, the attention mechanism takes the dot product of the vector representations of “he” and each word:

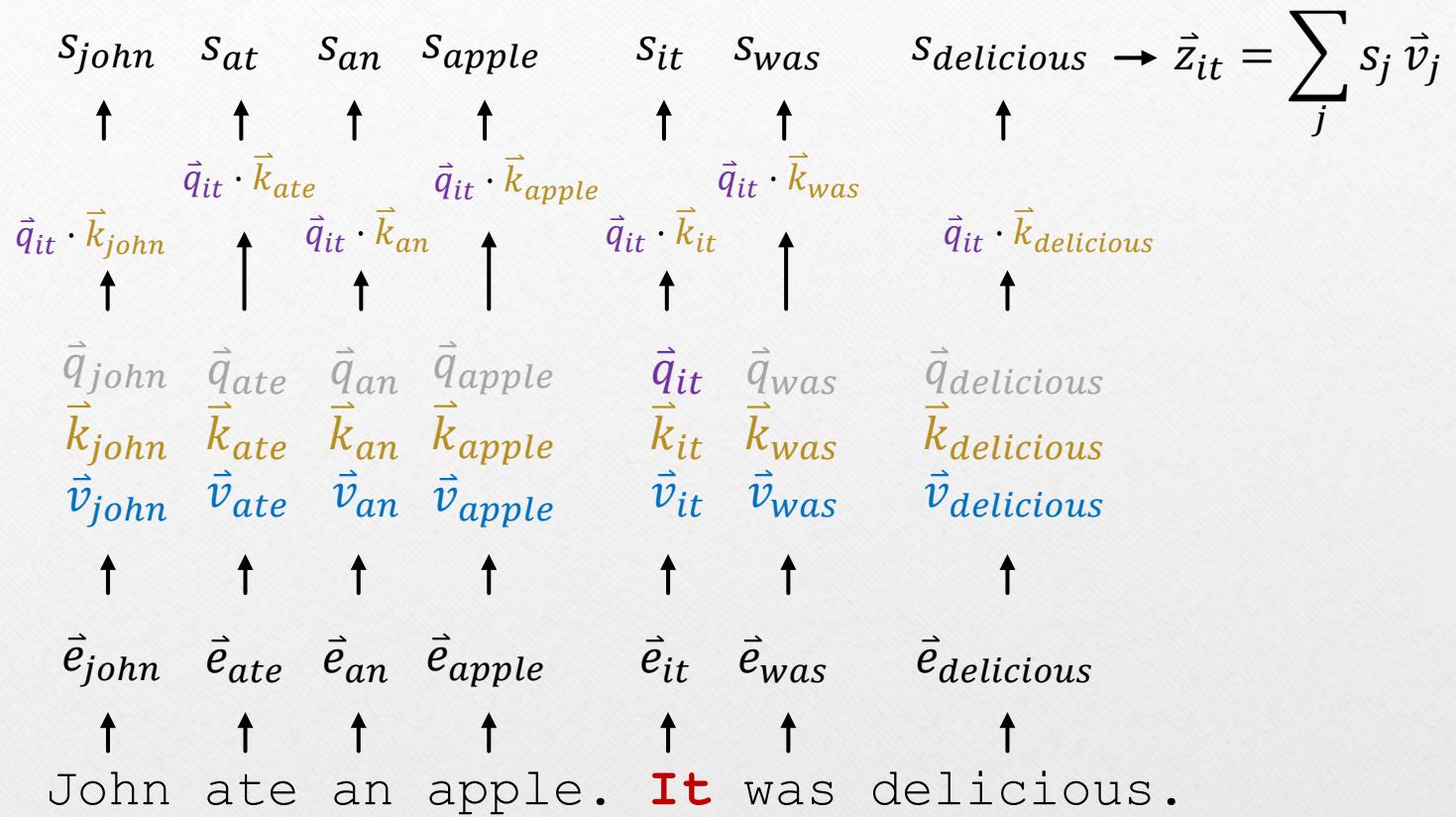
$$s_{ate} = \vec{h}_{he} \cdot \vec{h}_{ate} = 0$$

$$s_{john} = \vec{h}_{he} \cdot \vec{h}_{john} = 1$$

- Words with similar vector representations get a higher weight. This allows the model to “pay attention” to relevant words.



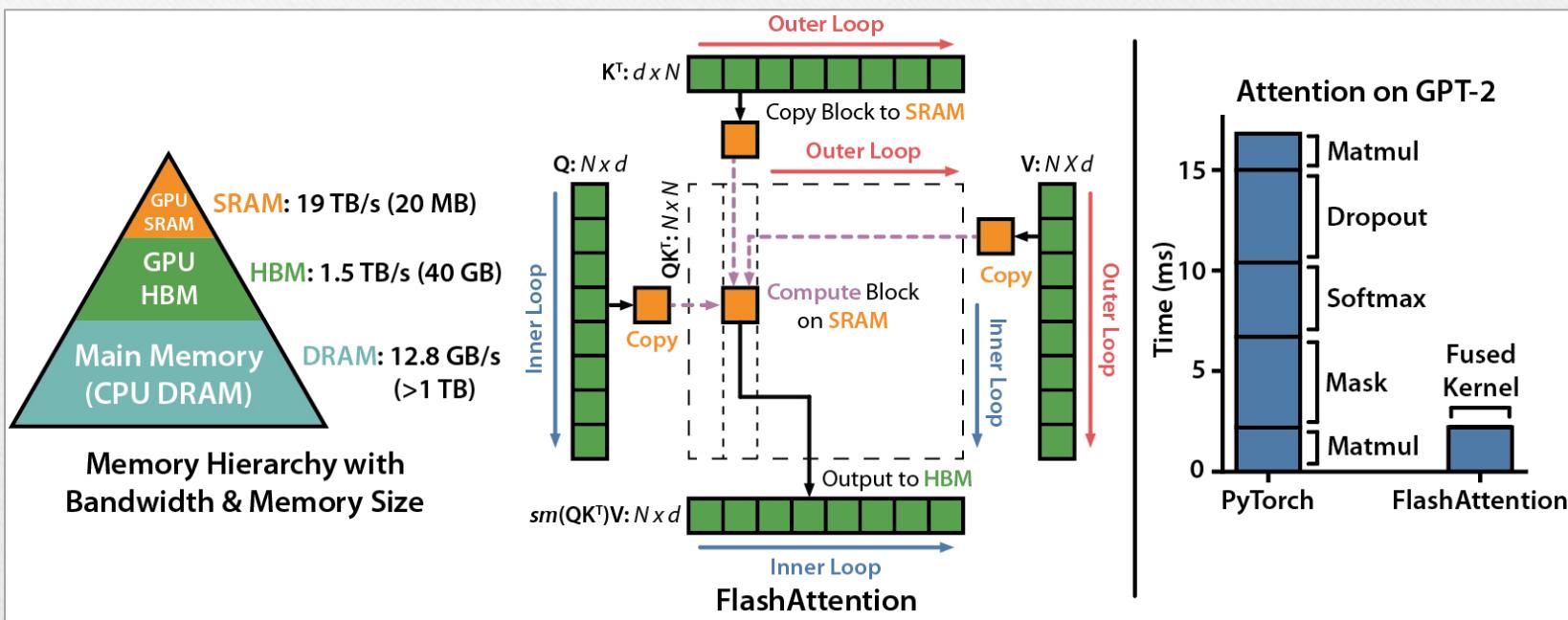
As self-attention is applied to all words, the mechanism is bidirectional—or more accurately, non-directional. It makes no distinction between words in different position.

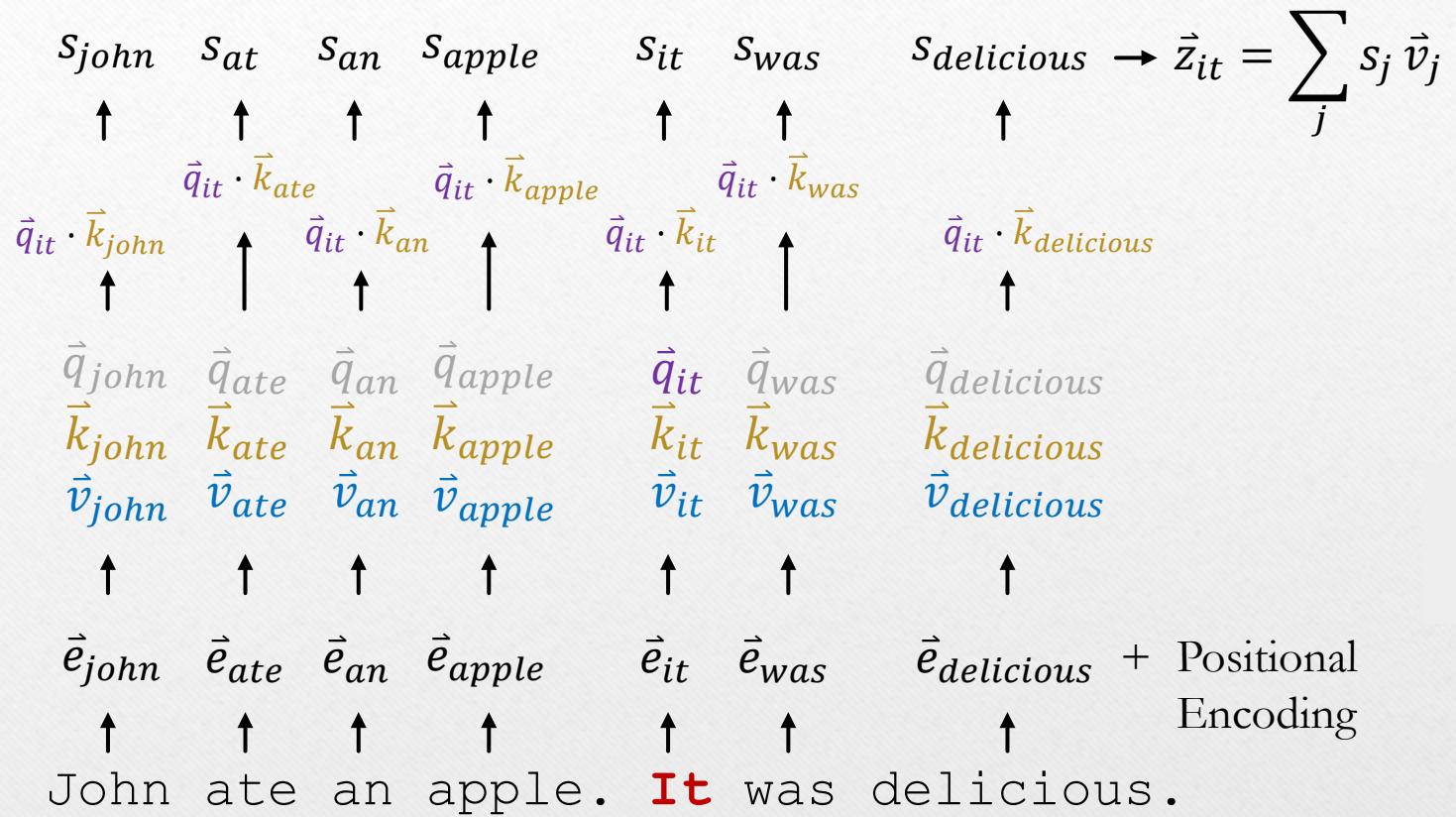


The fact that self-attention is applied to all words means that it takes  $n^2$  operations and memory to process  $n$  words. This is costly and is the main reason why older transformer models cannot handle long sequence of words.

There are different ways to make self-attention less costly, including:

- Computing an approximate value instead of an exact value
- Limiting self-attention to nearby words
- Better algorithm to take better advantage of GPU architecture:  
**Flash Attention**

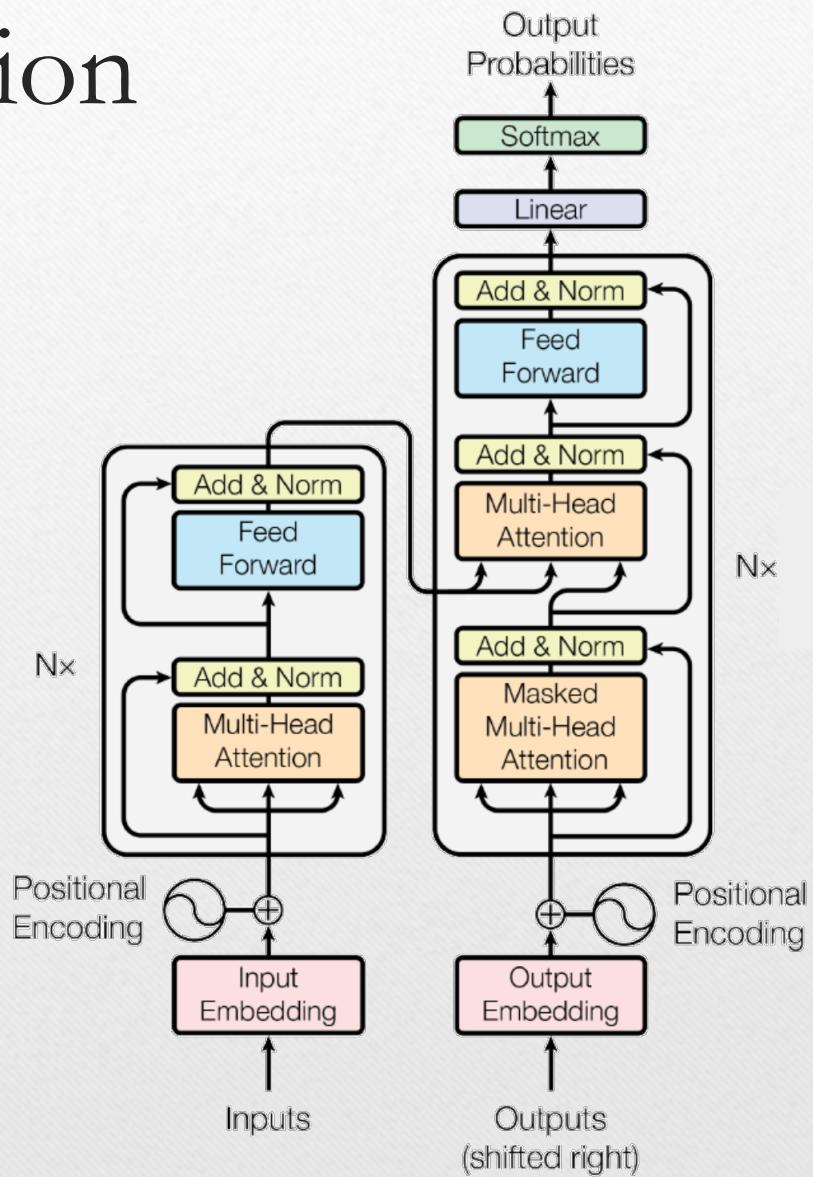


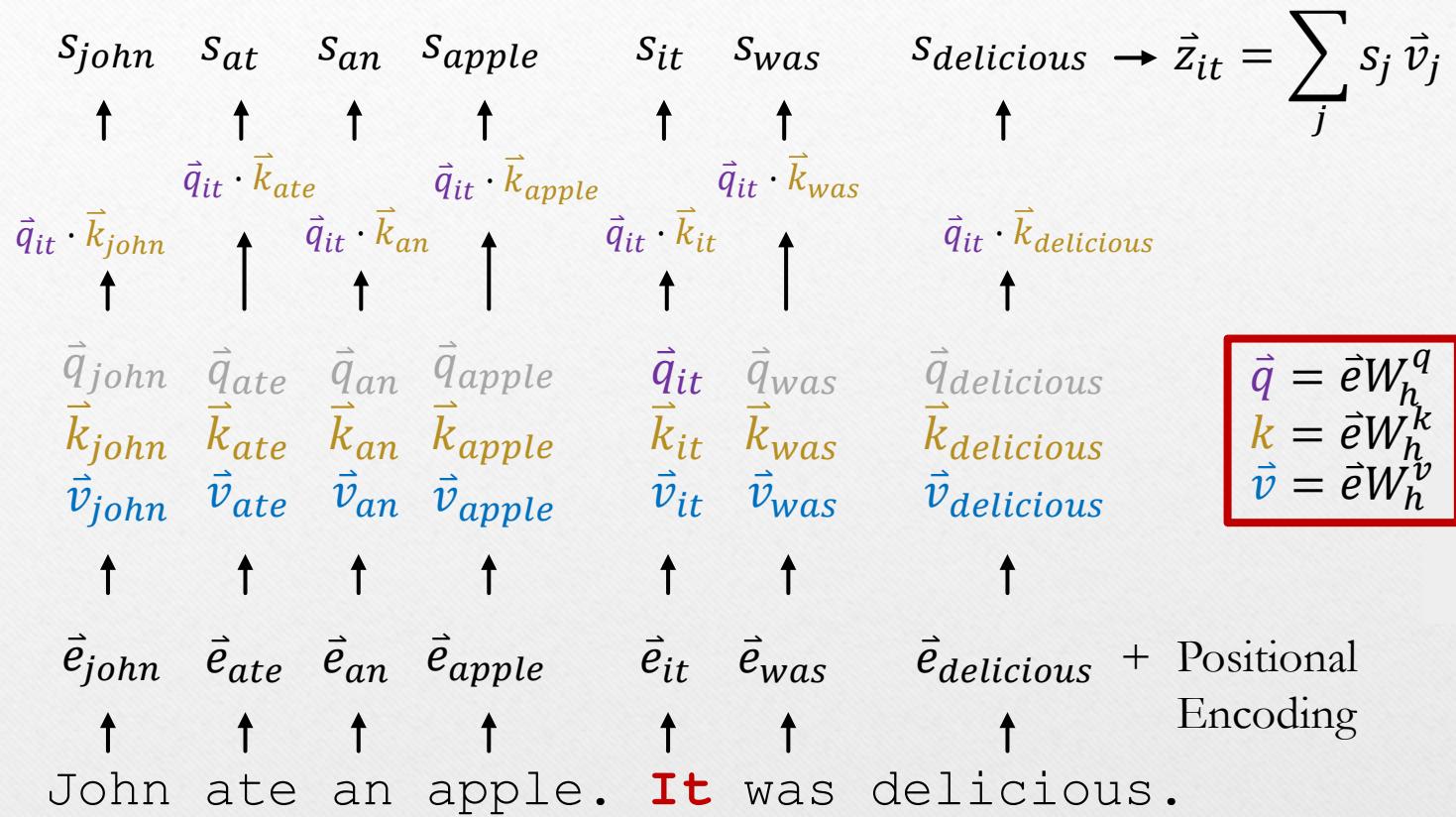


To provide the model with word positions, a **positional encoding** is added to the usual word embedding. The positional encoding can be fixed or learnt.

# Multi-Head Attention

- What we just saw is **single-head (self-)attention**.
- **Multi-head attention** simply means we have multiple copies of the process, each with different weights.





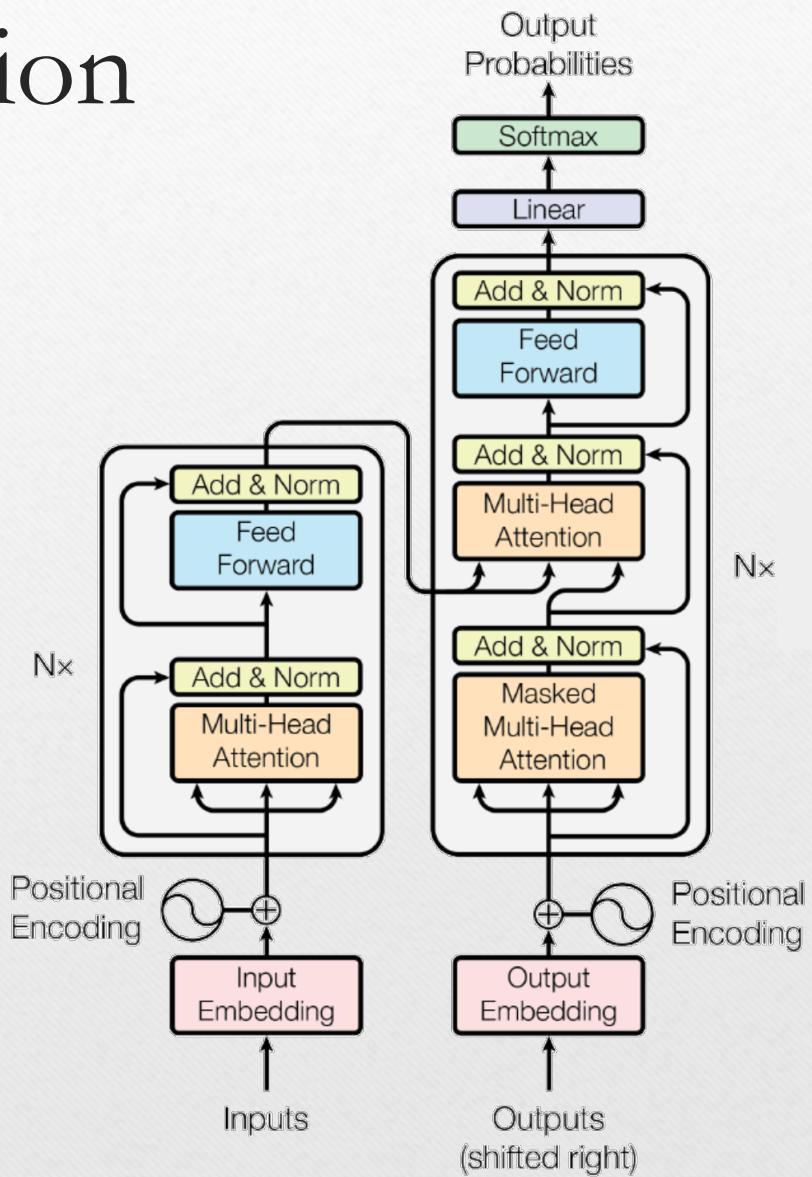
In multi-head attention, query, key and value are projections of the previous layer's output. You can view these as new embeddings capturing specific aspects of a word.

The weights  $W_h^q$ ,  $W_h^k$  and  $W_h^v$  are learnt from training.

# Multi-Head Attention

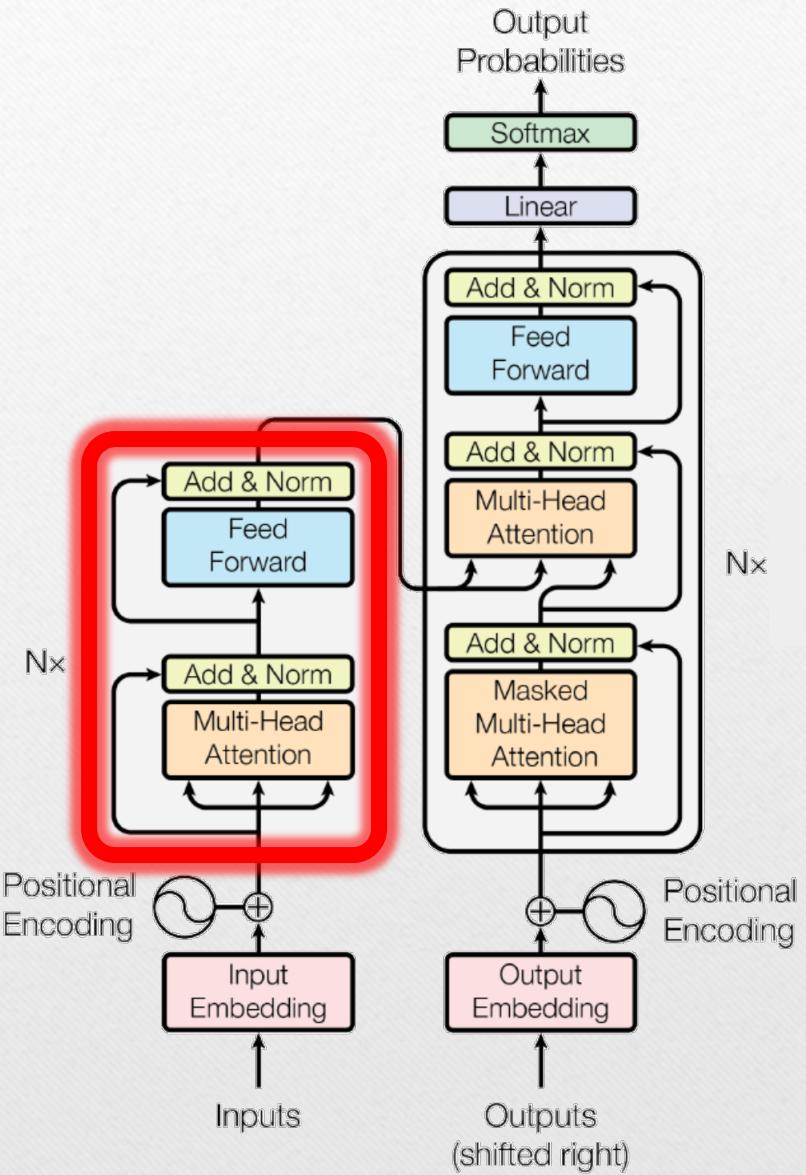
- What we just saw is **single-head (self-)attention**.
- **Multi-head attention** simply means we have multiple copies of the process, each with different weights.
- The weights from the heads are concatenated and multiplied to yet another set of weights, generating the final output of the attention layer:

$$\vec{z}_i = [\vec{z}_{1,i}, \dots, \vec{z}_{h,i}]W_o$$



# Encoder

- The output from the attention mechanism is then
  - Added to the output of the previous layer (i.e. residual connection)
  - Normalized
  - Fed to a fully-connected network
  - Another residual connection added and normalized
- The whole set of computation is called a **Transformer Encoder Layer**.
- Multiple encoder layers stacked together form the **Transformer Encoder**.
- The final output is a vector for each word, just like RNN.

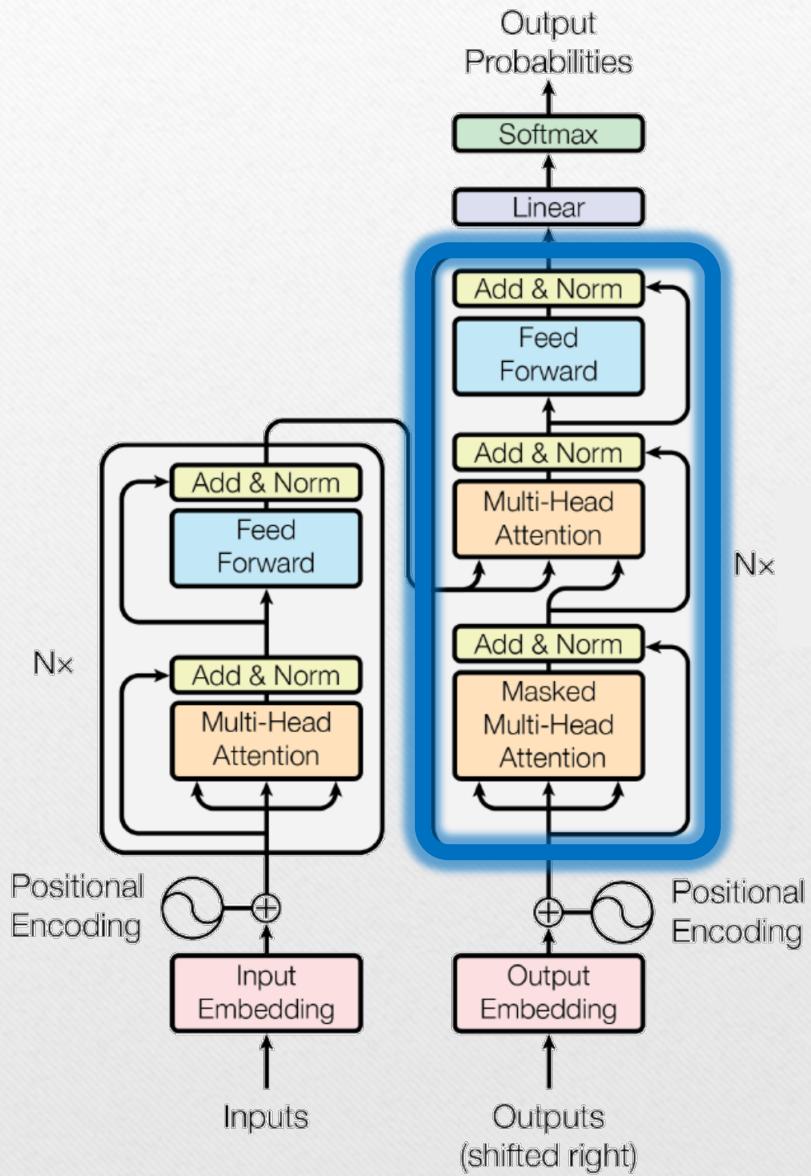


# Decoder

The **Transformer Decoder** is responsible for generating a new string of text based on the input.

A **Transformer Decoder Layer** is very similar to the encoder layer, with two crucial differences:

1. The self-attention mechanism is masked. When we process word  $i$ , the keys and values of words  $> i$  are not used. This is similar to unidirectional RNN.
2. There is a second attention mechanism where the keys and values are the output of the encoder.

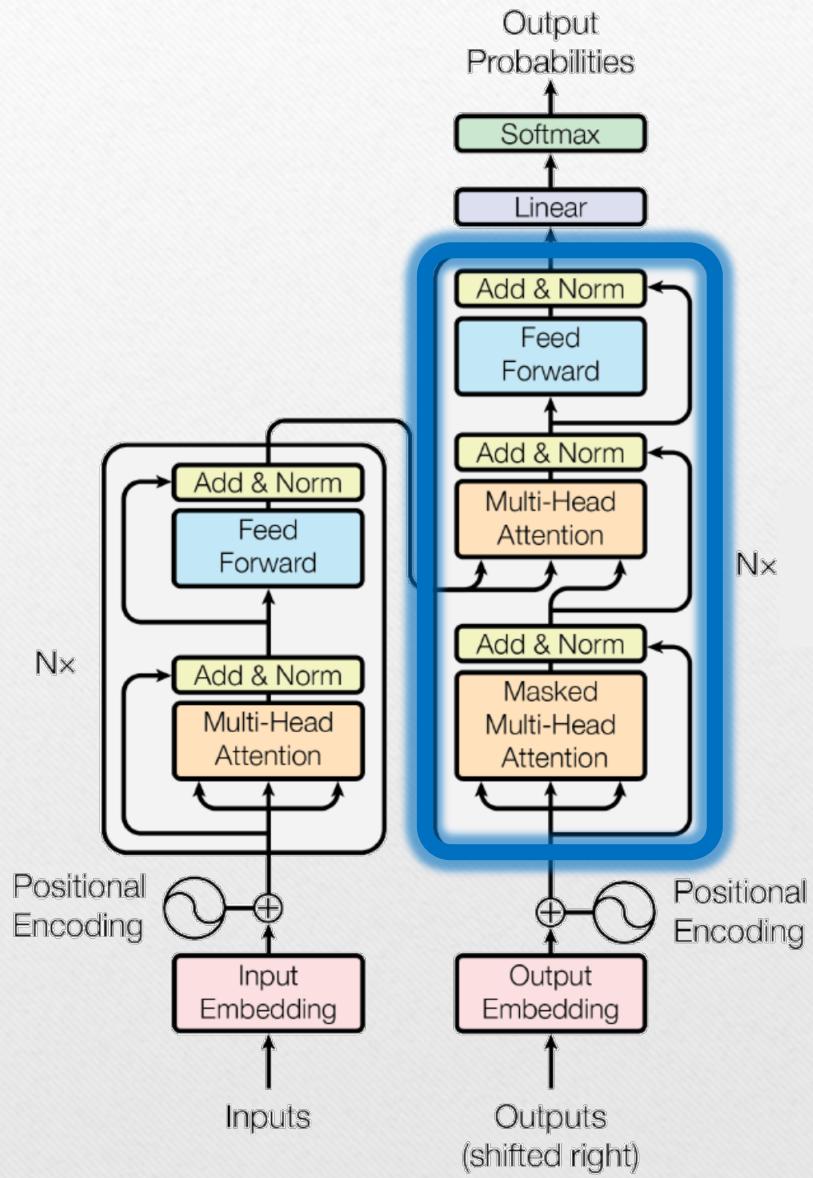


# Decoder

The decoder takes a list of words as input. It predicts what the next word should be.

During training, the input to the decoder is simply the desired output shifted right by one word.

- Desired output:  
John ate an apple [end]
- Input:  
[start] John ate an apple
- When predicting ate, the effective input is [start] John



# Important Families of Transformer Model

There are three major families of Transformer-based models:

- The original Transformer has Transformer Encoder and Decoder. The most important model in this family is Google's [T5](#).
- GPT (**G**enerative **P**re-**T**raining) only uses the Transformer Decoder. Models in this family include OpenAI's GPT [1](#), [2](#), [3](#) and [4](#), Meta's LLaMA [1](#), [2](#) and [3](#), MosaicML's [MPT](#) and TII's [Falcon](#).
- BERT (**B**idirectional **E**ncoder **R**epresentation from **T**ransformers) only uses the Transformer Encoder. Models in this family include [BERT](#), [RoBERTa](#), [DistilBERT](#) and [ALBERT](#).

# Inference of Decoder-Only Model

The decoder takes a list of words as input. It predicts what the next word should be.

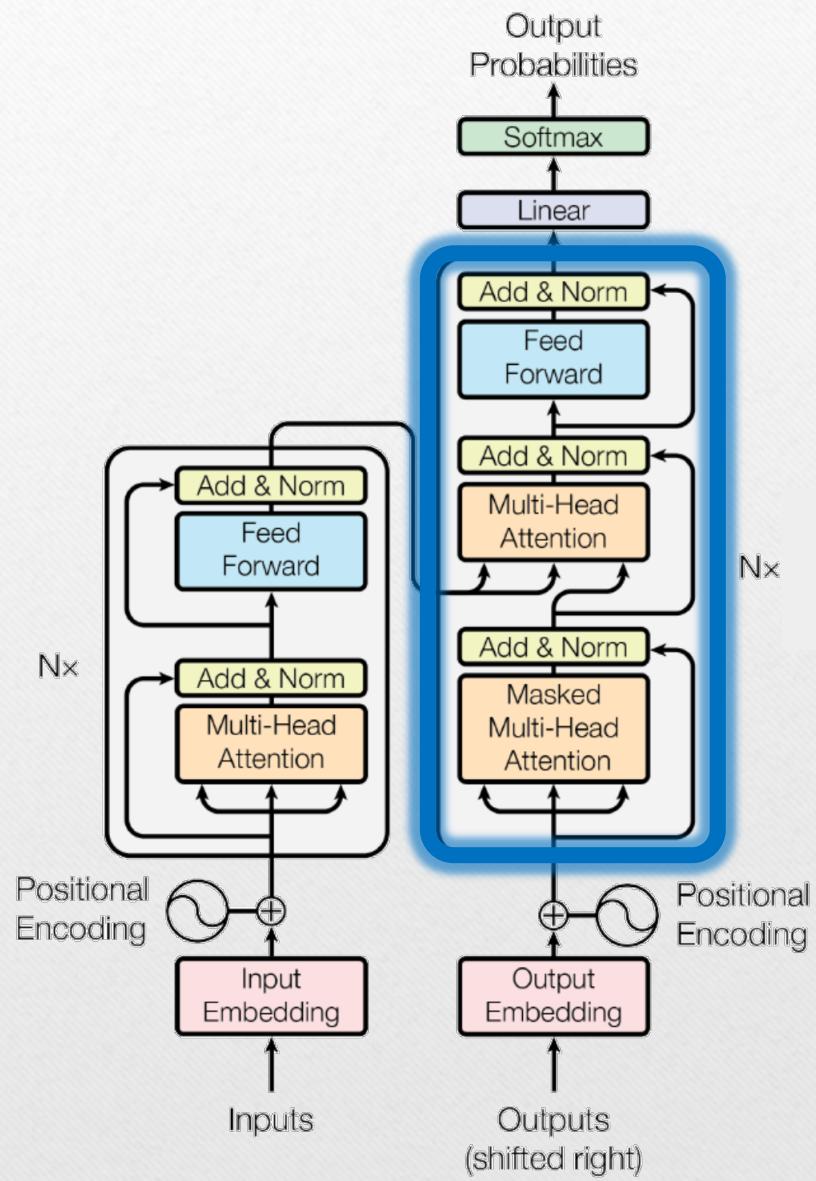
During inference, after the decoder generates the next word, it will be used as part of the input to generate the word after the next.

[start] → ?

Say the model predicts Paul, so next we have:

[start] Paul → ?

The way decoder-only models work means they are mostly used for text generation and not other tasks.

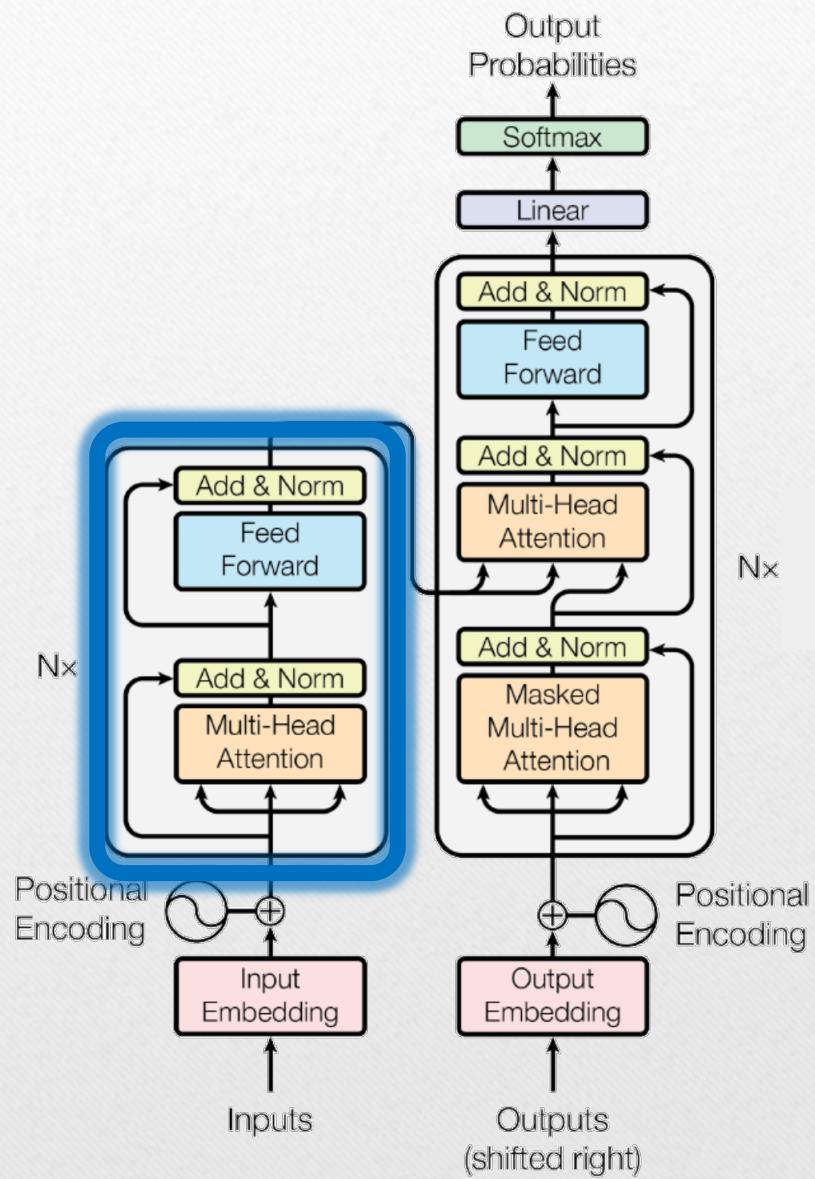


# Inference of Encoder-Only Model

When using a trained model for inference, the following happens:

1. Provide the encoder with the input text, e.g. a sentence in English.
2. Encoder generates a vector for each word.

Because the generated vector for each word utilizes information from all words, it is common to use just one vector—e.g. the first word—for downstream tasks such as text classification.



# Reinforcement Learning with Human Feedback (RLHF)

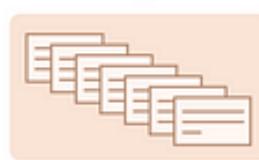
- Instruction-following LLMs—i.e. chatbots—are fine-tuned from instruction-answer datasets.
- Such datasets are relatively small—10,000s samples vs billions of samples in the pre-training stage. How to scale up the fine-tuning process?
- **RLHF** trains a reward model, itself a small LLM, based on a relatively small number of samples of human choices over pairs of generated answers
- This reward model can then be used to fine-tune the chatbot over many *generated* samples

## ① Collect human feedback

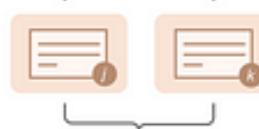
A Reddit post is sampled from the Reddit TL;DR dataset.



Various policies are used to sample a set of summaries.



Two summaries are selected for evaluation.



A human judges which is a better summary of the post.



" $j$  is better than  $k$ "

## ② Train reward model

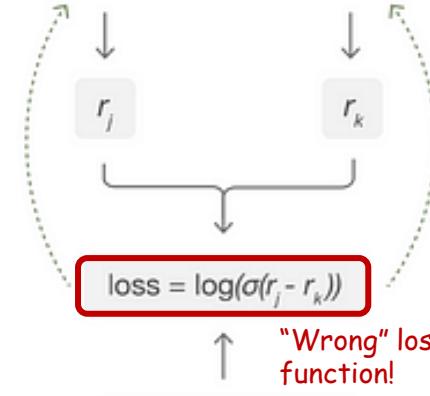
One post with two summaries judged by a human are fed to the reward model.



The reward model calculates a reward  $r$  for each summary.



The loss is calculated based on the rewards and human label, and is used to update the reward model.

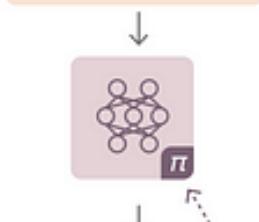


## ③ Train policy with PPO

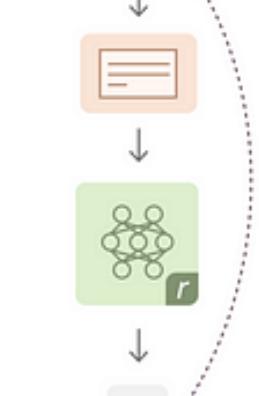
A new post is sampled from the dataset.



The policy  $\pi$  generates a summary for the post.



The reward model calculates a reward for the summary.



The reward is used to update the policy via PPO.

Source: Nisan Stiennon et al (2020). "[Learning to Summarize from Human Feedback](#)". NeurIPS  
Also read: [Proximal Policy Optimization](#)

# Text Generation Methods

- The decoder takes a list of words as input. It predicts what the next word should be.

[start] → ?

- The simplest way to do this is to return the most likely word.

$$x_t = \operatorname{argmax}_{x \in X} P(x|x_{t-1})$$

This is called **greedy search**.

- Greedy search often produce text that is natural, because the sequence of next most likely words is usually not the sequence that has the highest overall probability. i.e.

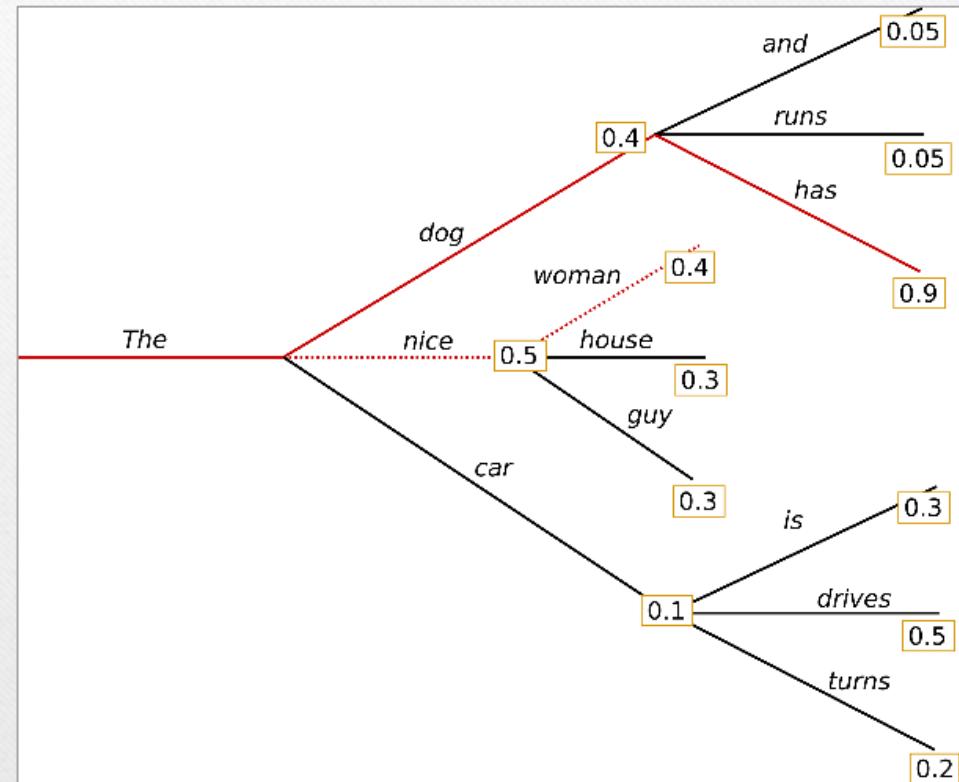
$$x_t = \operatorname{argmax}_{x \in X} P(x|x_{t-1}) \text{ and } x_{t+1} = \operatorname{argmax}_{x \in X} P(x|x_t)$$

does not give you the same answer as

$$x_t, x_{t+1} = \operatorname{argmax}_{x_t, x_{t+1} \in X} P(x, x_{t+1}|x_{t-1})$$

# Text Generation Methods

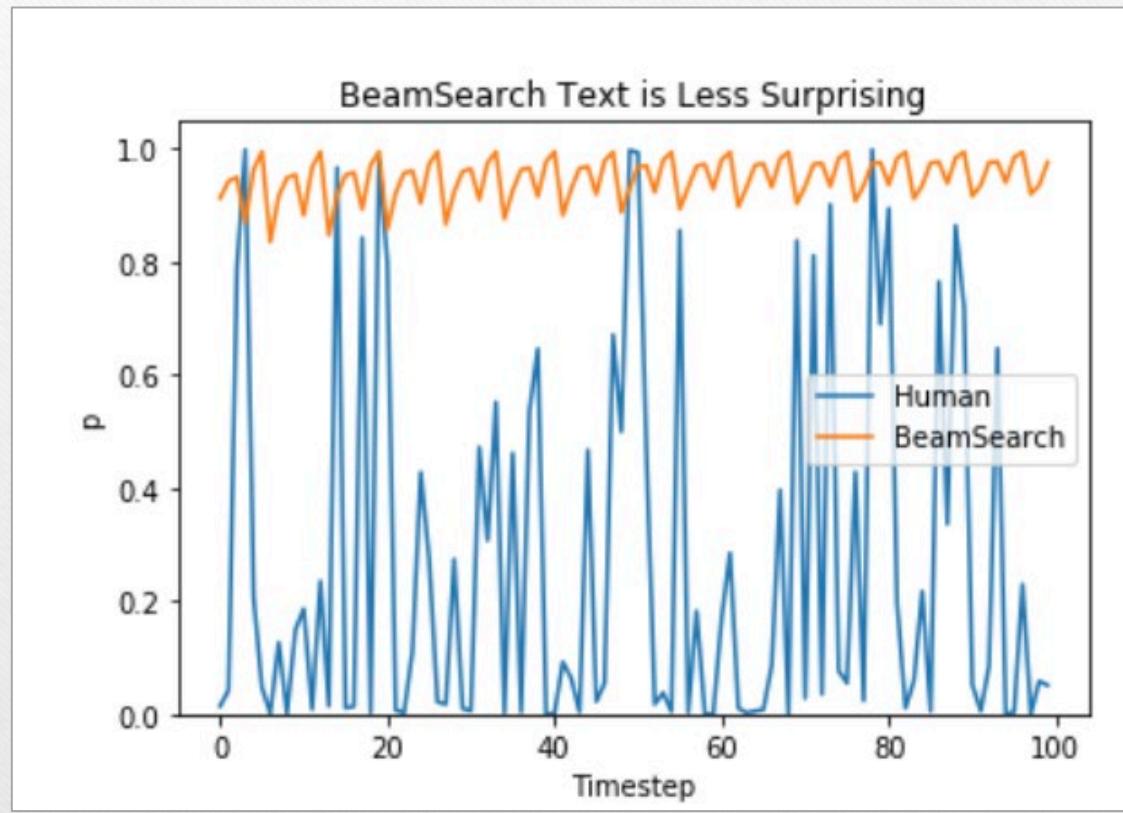
**Beam search** generates multiple sequences of words and pick the sequence that has the highest overall probability. The number of sequences to keep track of is a hyperparameter.



Source: [HuggingFace Blog](#)

# Text Generation Methods

But human speech is not based on uttering the next most likely word:



Source: [HuggingFace Blog](#)

# Text Generation Methods

Regardless of whether we use beam search or not—chatbots most do not due to latency issue—we can also **sample** the next word based on their estimated probability of appearing, instead of always picking the most likely word.

$$P(x_t = x_i) = P(x_i | x_{t-1})$$

Hyperparameters that control sampling:

- **Temperature:** a lower temperature sharpens the distribution. Temperature = 0 is the same as greedy search.
- **Top-K:** Only sample from the top  $k$ -most likely words.
- **Top-P:** Only sample words are within the top cumulative distribution of  $p$
- **Repetition penalty:** lower the probability of sampling words that have already been sampled before.