

# CSCI-567: Discussion Session 09/22

Joshua Robinson, Oliver Liu



The background is a solid orange color. In the top-left corner, there are three vertical bars of varying heights, each composed of several overlapping semi-transparent orange circles. In the bottom-right corner, there are four vertical bars of increasing height from left to right, each also composed of several overlapping semi-transparent orange circles.

# The Transformer



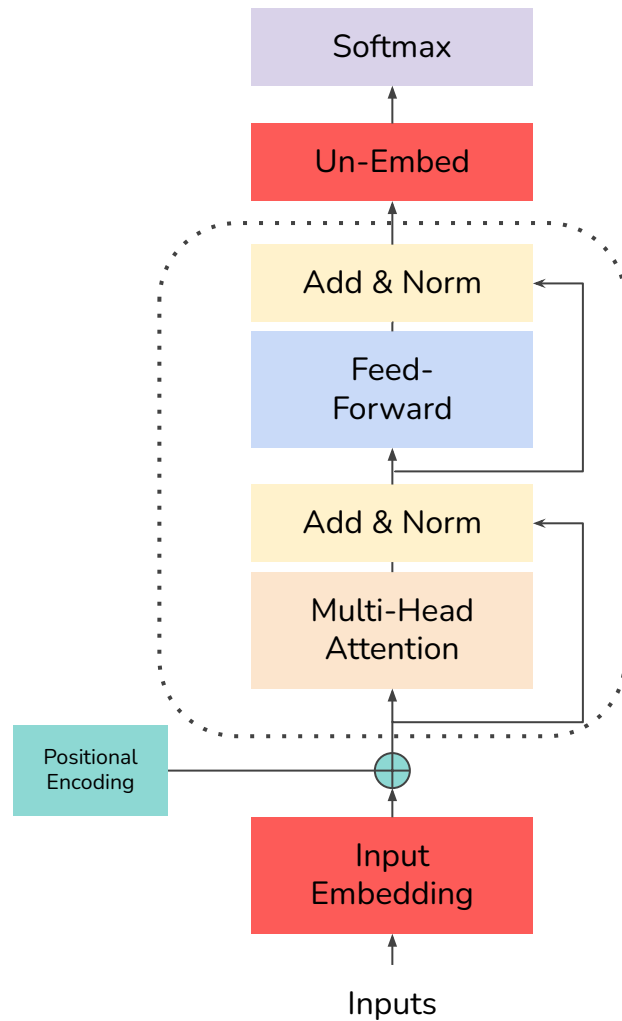
# The Transformer Architecture

- Designed for **Sequence Learning**;
- Very **scalable** (up to trillions of parameters and training tokens);
- Built on stack of **self-attention layers**.
- Useful resources:
  - [The Illustrated GPT-2](#),
  - [The Annotated Transformer](#).



# Model Overview

- Inputs
- Embedding Matrix
- Positional Encoding
- Self-Attention Layer
- Un-embed Matrix
- Softmax





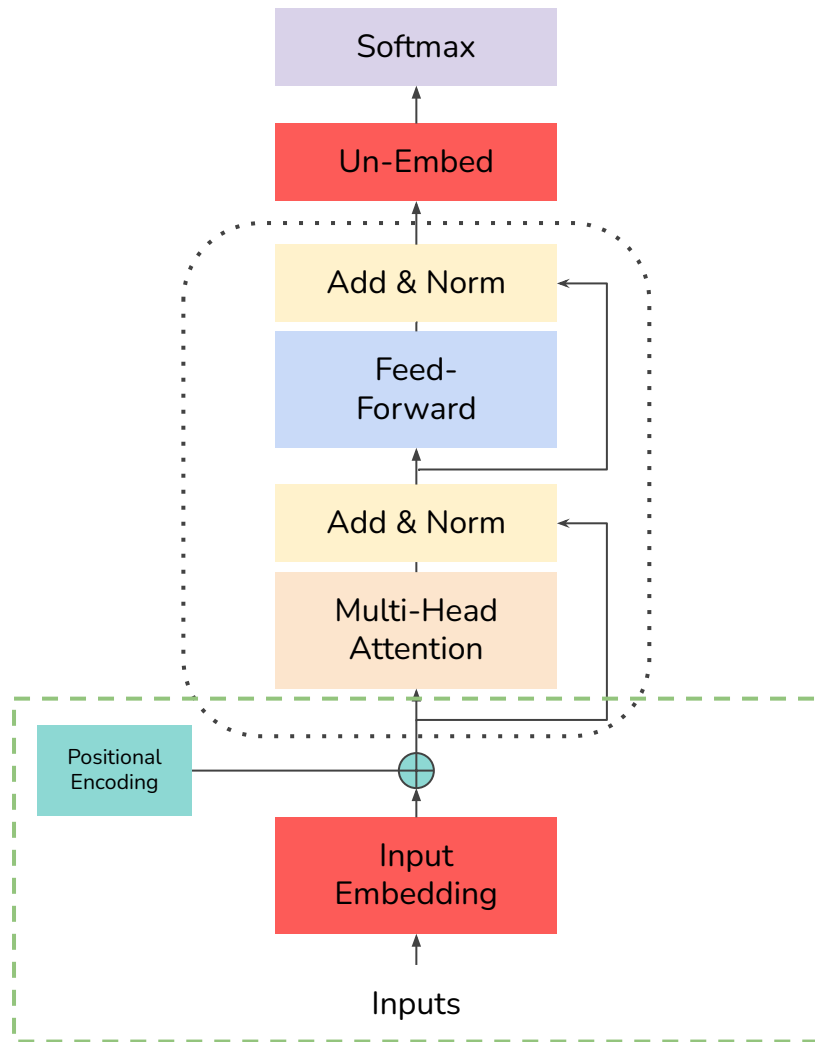
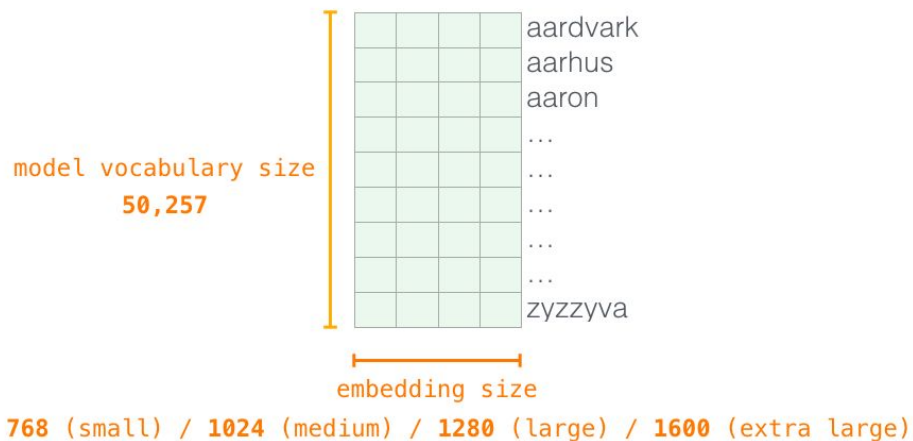
# Inputs and Embedding

- **Input in natural language:** *I have a dream.*
- **Vocabulary:** a mapping between each word to a unique index in  $\{0, 1, \dots, V\}$
- **Tokenized input:**  $[x_1, x_2, x_3, x_4]$  a sequence of length 4
  - $x_i$ : index of the corresponding vocabulary (e.g.  $x_3$  is the index of *dream*)
- We feed tokenized input to the **Input Embedding  $E$** , a matrix of size  $V \times D$ 
  - Input to the first self-attention layer:  $X = [E_{x_1}, E_{x_2}, E_{x_3}, E_{x_4}]$  of shape  $4 \times D$



# Inputs and Embedding

## Token Embeddings (wte)





# Self-Attention

- **Attention:** Given a **Query**, computed a weighted average of **Values** based on the query's similarity with **Keys**.
- **Self-Attention:** **Q**, **K**, **V** are linear transformations of **X**.
  - **Q**, **K**, **V** are of size  $T \times D$
  - What is the size of the attention output?

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

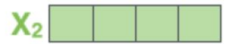
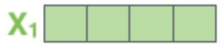


Input

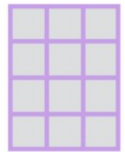
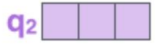
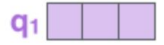
Thinking

Machines

Embedding

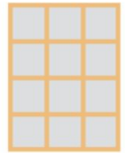
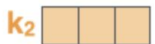
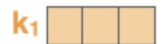


Queries



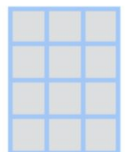
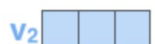
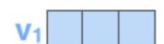
$W^Q$

Keys



$W^K$

Values



$W^V$





Input

Embedding

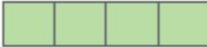
Queries

Keys

Values

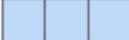
Score

Thinking

$x_1$  

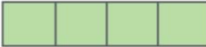
$q_1$  

$k_1$  

$v_1$  

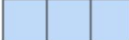
$$q_1 \cdot k_1 = 112$$

Machines

$x_2$  

$q_2$  

$k_2$  

$v_2$  

$$q_1 \cdot k_2 = 96$$



Input

Embedding

Queries

Keys

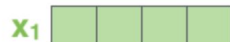
Values

Score

Divide by  $8 (\sqrt{d_k})$

Softmax

Thinking

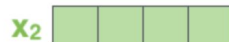


$q_1 \cdot k_1 = 112$

14

0.88

Machines



$q_2 \cdot k_2 = 96$

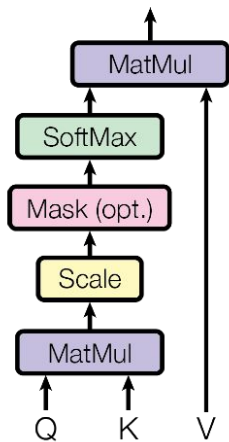
12

0.12

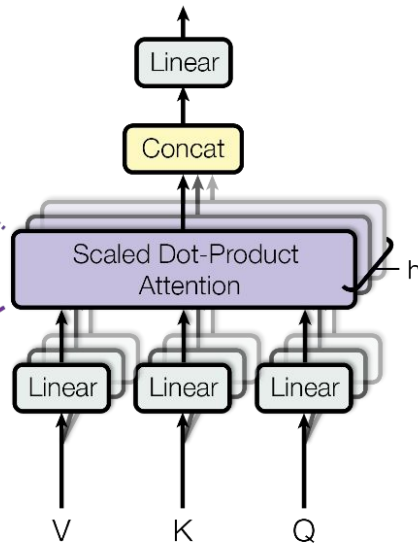


# Multi-Head Attention

Scaled Dot-Product Attention



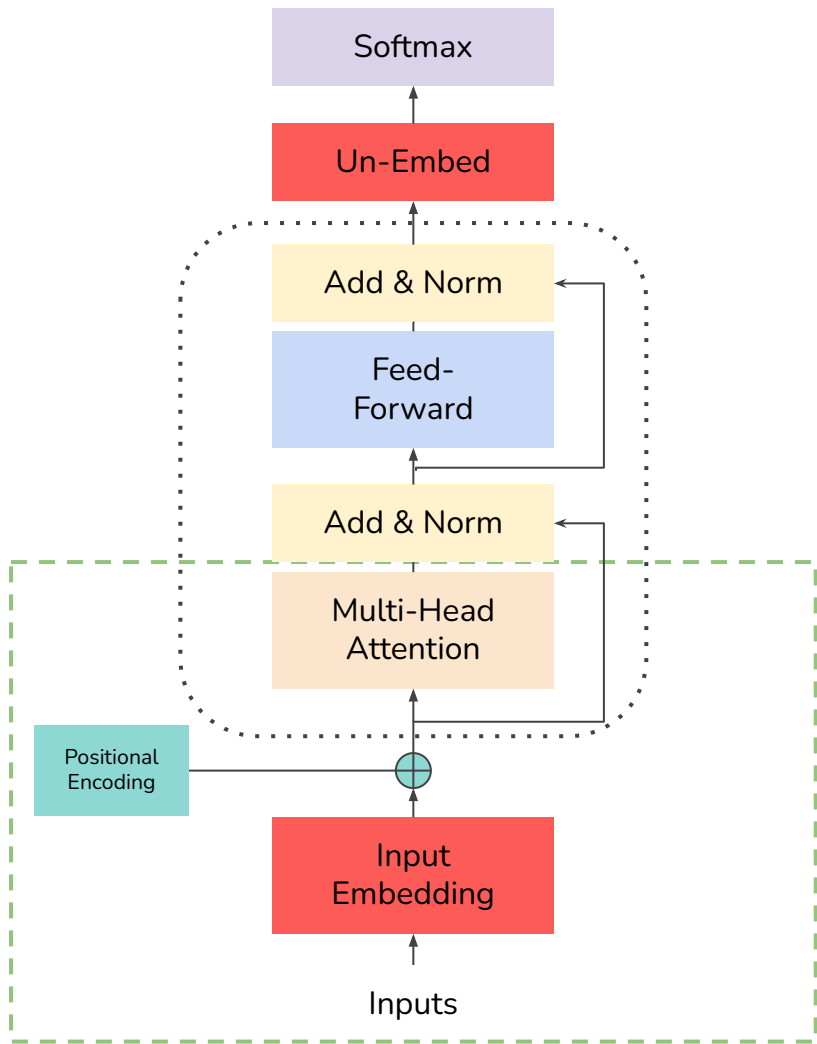
Multi-Head Attention





# Story so far

- **Input embedding** convert a set of discrete tokens to continuous vectors.
- The **attention mechanism** computes a representation of each token position based on their relevance with other tokens.
- The attention output is of size  $T \times D$ .
- We ignore add & norm for now...





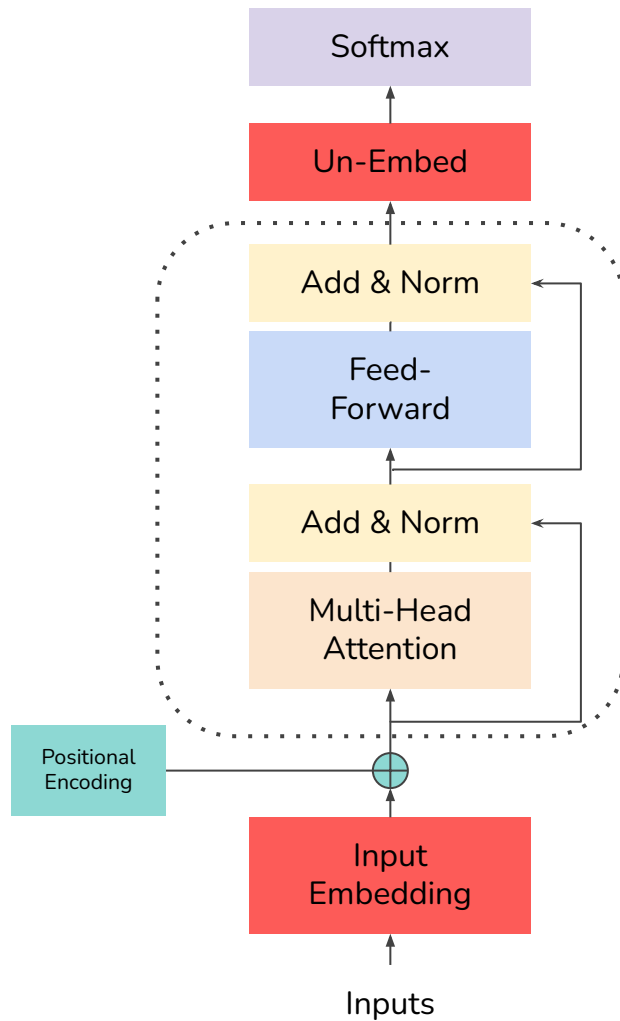
# Positional Feed-forward

- Apply the same MLP at each token position.
- **Input:**  $[E_1, E_2, \dots, E_T]$  of shape in  $T \times D$
- **PFF** computes:  $[MLP(E_1), MLP(E_2), \dots, MLP(E_T)]$  for the same *MLP*!
- **Output** is of the same shape!



# Add & Norm

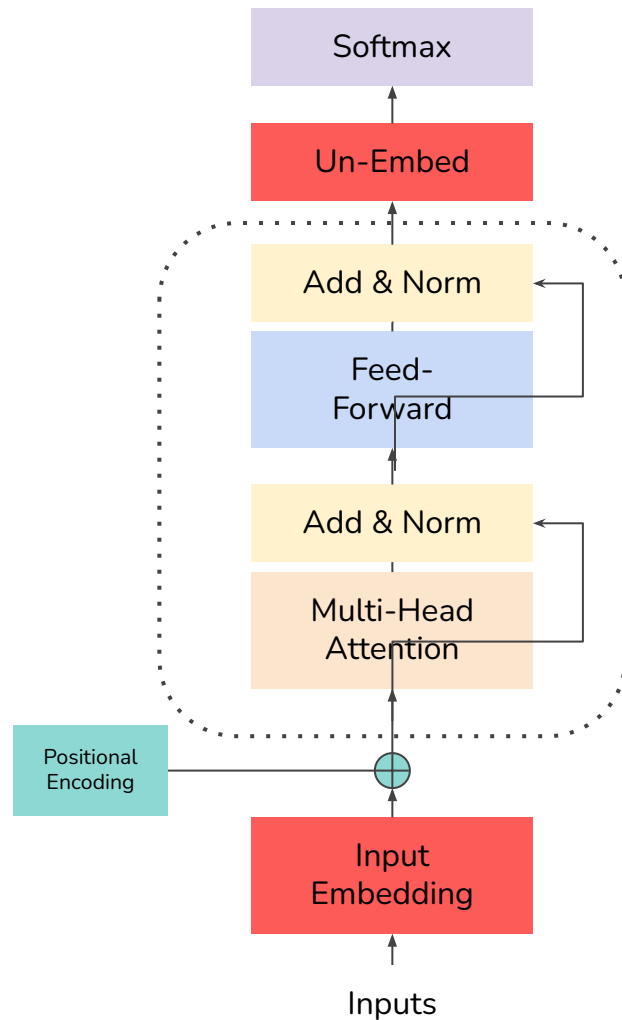
- **Add**, aka **residual connection**:  
 $X \leftarrow X + f(X)$
- $f$  is **Self-Attention** and **MLP** in transformer
- **LayerNorm** normalizes hidden vectors of each token position.





# Un-Embed

- Make a multi-class prediction at each position. This is just softmax regression!
- What is the dimension of the un-embed matrix?
- Typically, this is the transpose of the embedding matrix.



The background is a solid orange color. In the top-left corner, there are three vertical bars of varying heights, each composed of four overlapping circles. In the bottom-right corner, there are four vertical bars of increasing height from left to right, each composed of four overlapping circles.

# Transformer Variants





# Two Flavors of Self-Attention

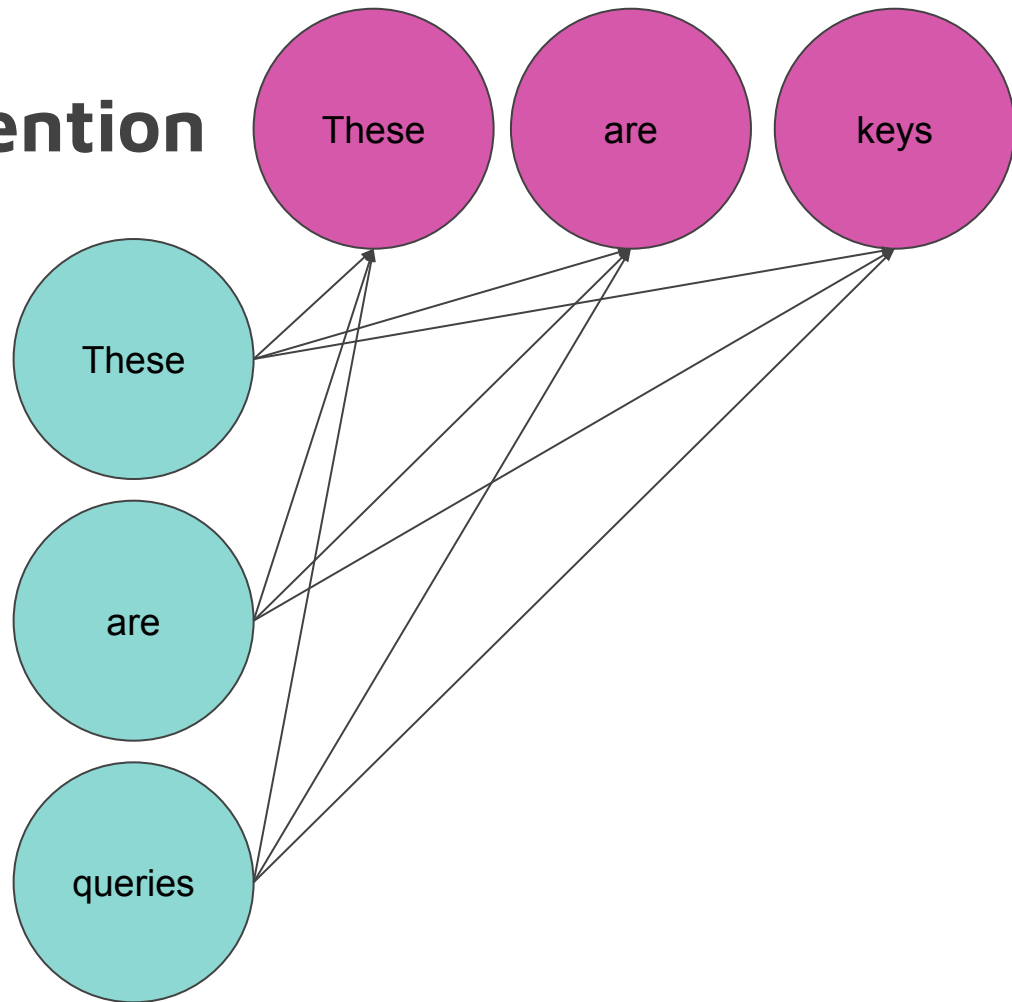
- Multi-Head Attention: Every query can attend to every key
- Masked Multi-Head Attention: Each query can only attend to its associated key and keys earlier in the sequence
  - In practice this done by setting masked positions to a very low number before softmax

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



# Multi-Head Attention

- Queries are in blue circles
- Keys are in pink circles

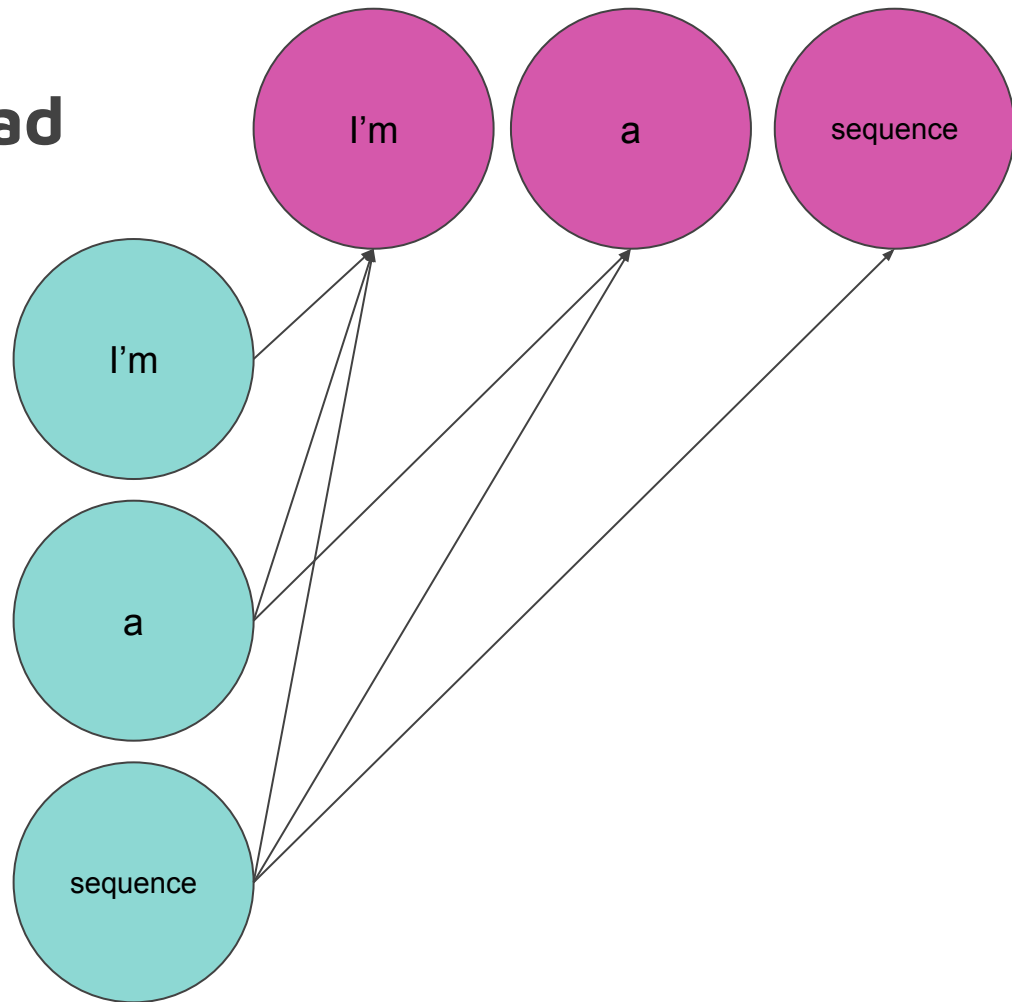


$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



# Masked Multi-Head Attention

- Queries are in blue circles
- Keys are in pink circles
- The spots with no link here would be set to a very low number before softmax



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



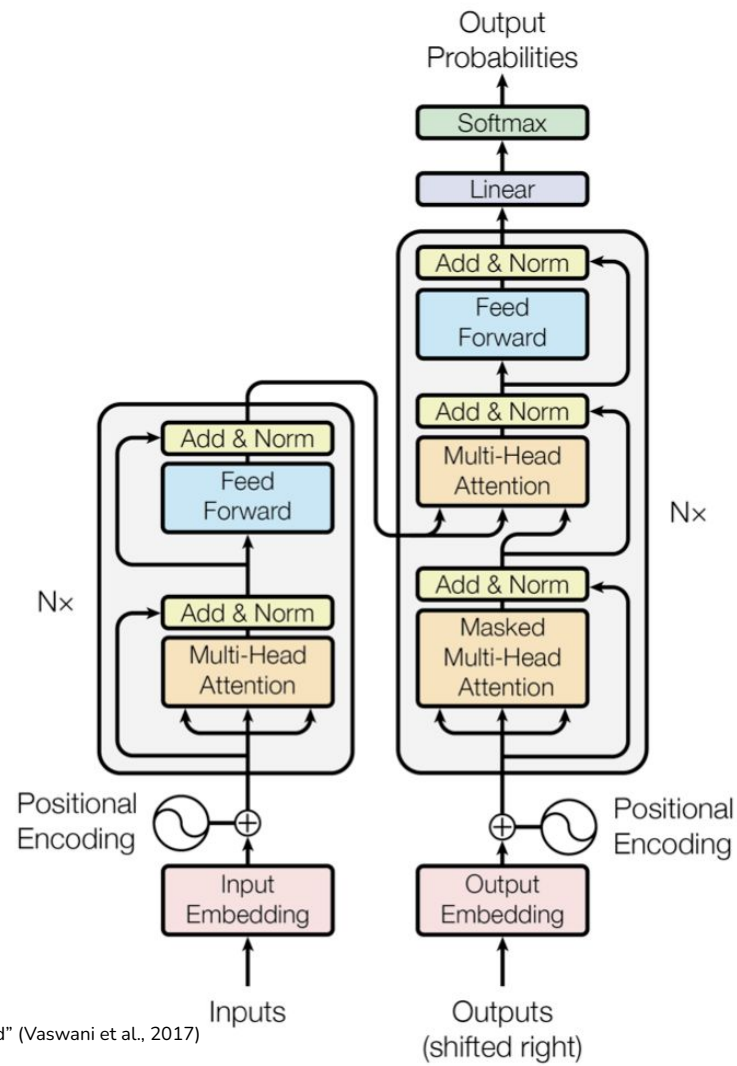
# The Transformer Family

- Encoder-decoder models
  - Original Transformer
  - T5
- Encoder-only models
  - BERT
- Decoder-only models
  - GPT-2
  - GPT-3
  - ChatGPT



# Encoder-Decoder Model

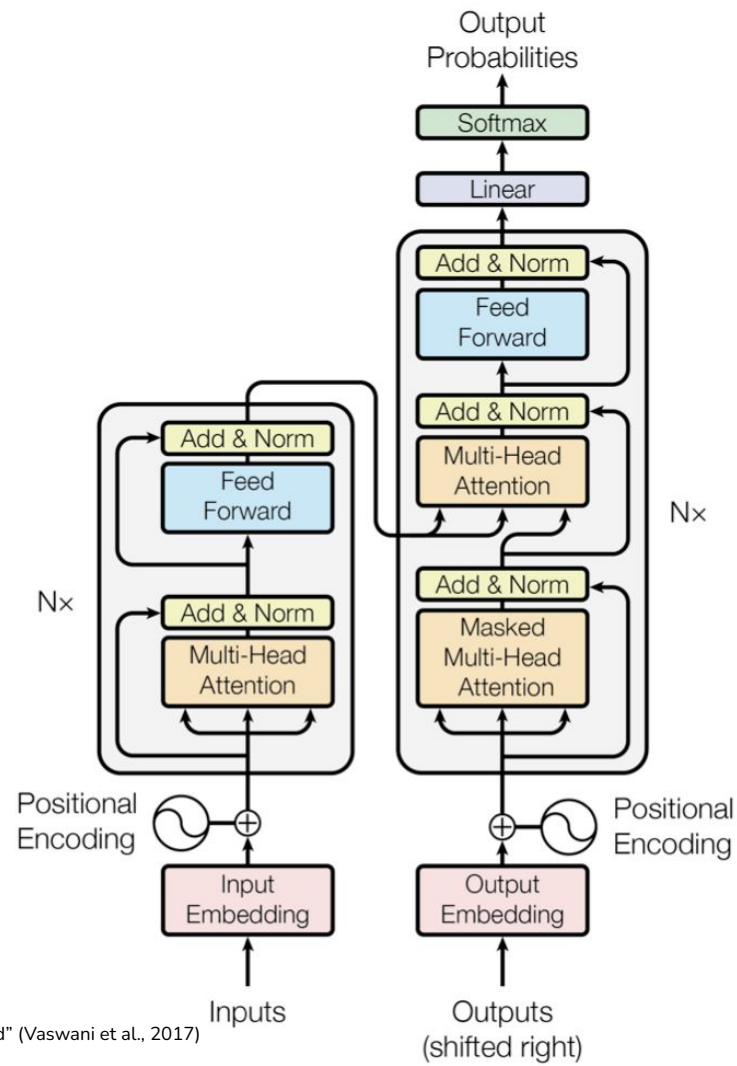
- Machine translation
- Summarization





# The Original Transformer

- [Learn more](#)
- Input like “How are you?”
- Outputs like “<s> Cómo estás?”
- Targets like “Cómo estás? <e>”
- Minimize cross entropy between predictions and labels





# T5

- [Learn more](#)
- “Span corruption” objective
- Minimize cross entropy between predictions and labels

Original text

Thank you ~~for~~ ~~inviting~~ me to your party ~~last~~ week.

Inputs

Thank you <X> me to your party <Y> week.

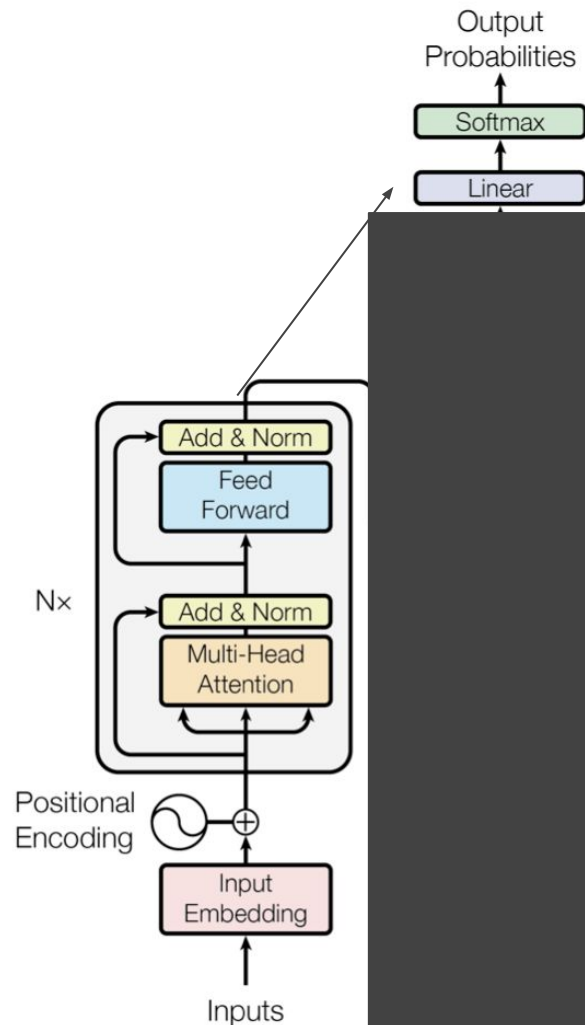
Targets

<X> for inviting <Y> last <Z>



# Encoder-Only Model

- Text classification
- Sentence representation







# BERT

- [Learn more](#)
- Objective is masked language modeling (MLM)
  - For 15% of tokens
    - Replace 80% with [MASK]
    - Replace 10% with a random token
    - Leave 10% unchanged
- Minimize cross entropy between predictions and labels
- Input is something like “I saw a [MASK] on my way to class tennis morning.”
- Use [CLS] token for classification
- So many related models: RoBERTa, DeBERTa, ALBERT, ELECTRA, etc.



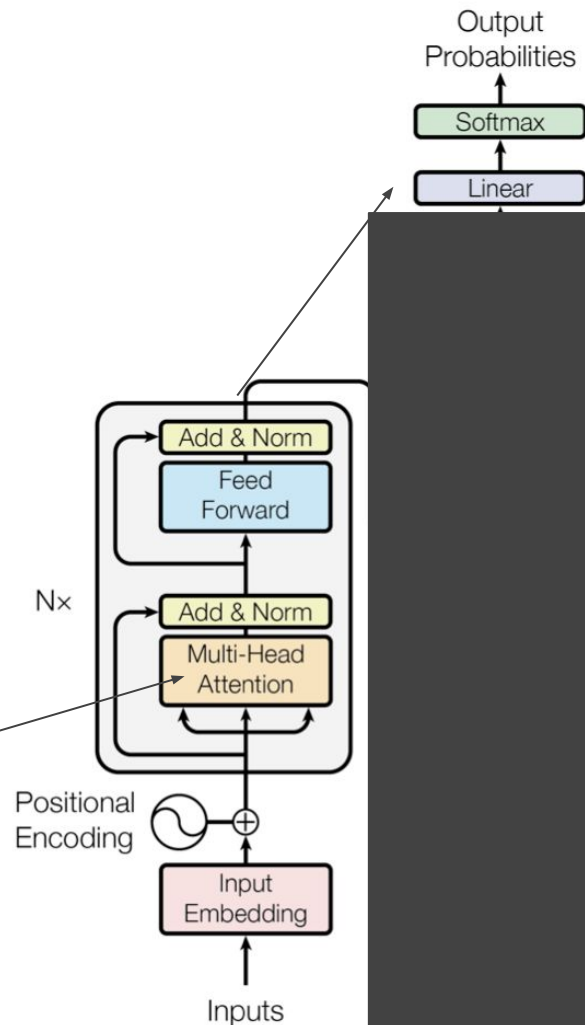
Image is from <https://muppet.fandom.com/wiki/Bert>



# Decoder-Only Model

- Language modeling

**MASKED!**





# GPT-2

- [Learn more](#)
- Objective is next-token prediction
- Minimize cross entropy between predictions and labels

“

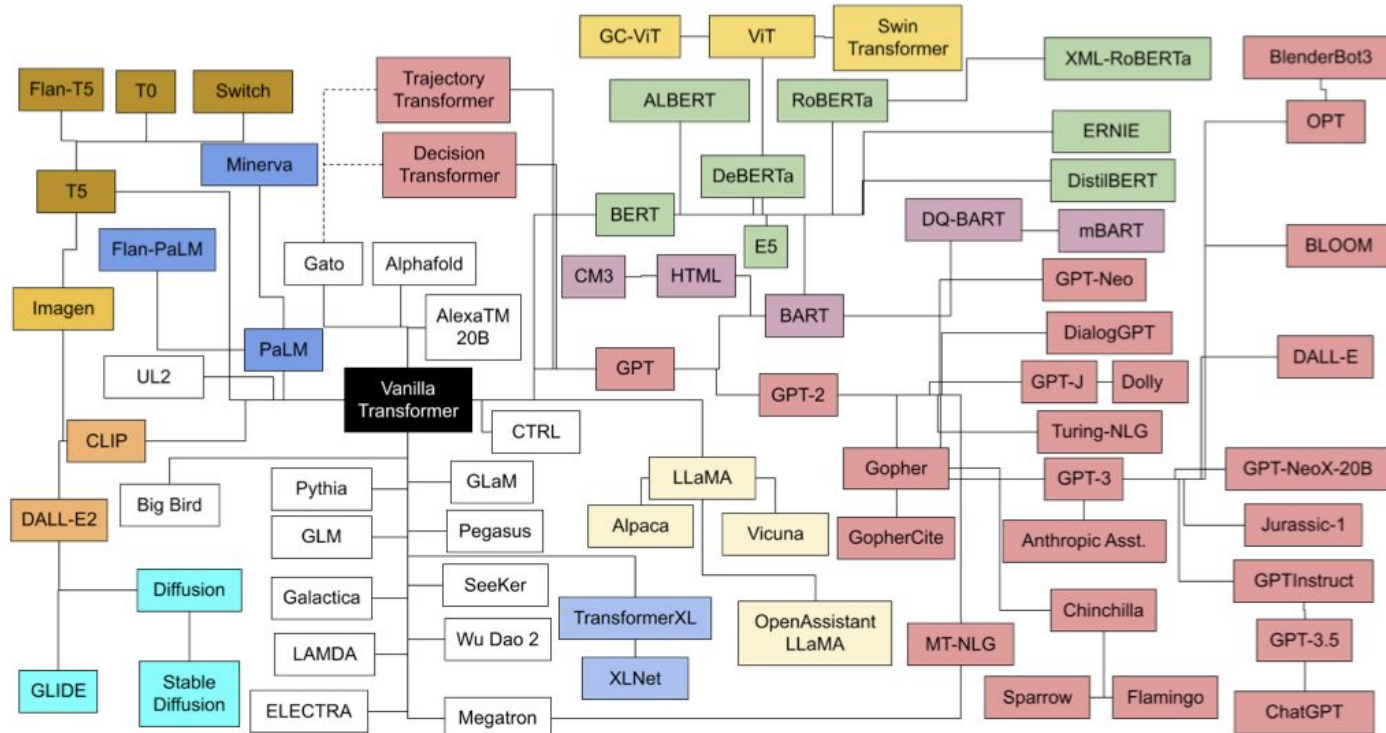
*My heart, why come you here alone?  
The wild thing of my heart is grown  
To be a thing,  
Fairy, and wild, and fair, and whole*

GPT-2

”

Image is from <https://gwern.net/gpt-2>

# Tip of the Iceberg



The figure is from "Transformer Models: an introduction and catalog" (Amatriain et al., 2023)