

Attention Mechanism and Transformers

This file is meant for personal use by tinlong@iscopedesign.com only.

Sharing or publishing the contents in part or full is liable for legal action.

Proprietary content. © Great Learning. All Rights Reserved. Unauthorized use or distribution prohibited.

Agenda

- Key Questions
- What are Transformer Models?
- How does a Transformer model work?
- What are the components of the Encoder?
- How does the Self-Attention Mechanism work?
- What are the components of the Decoder?
- What are the different transformer architectures?

This file is meant for personal use by tinlong@iscopedesign.com only.

Sharing or publishing the contents in part or full is liable for legal action.

Proprietary content. © Great Learning. All Rights Reserved. Unauthorized use or distribution prohibited.

Key Questions

What are Transformer Models?

How does a Transformer model work?

What are the components of the Encoder?

How does the Self-Attention Mechanism work?

What are the components of the Decoder?

What are the different transformer architectures?

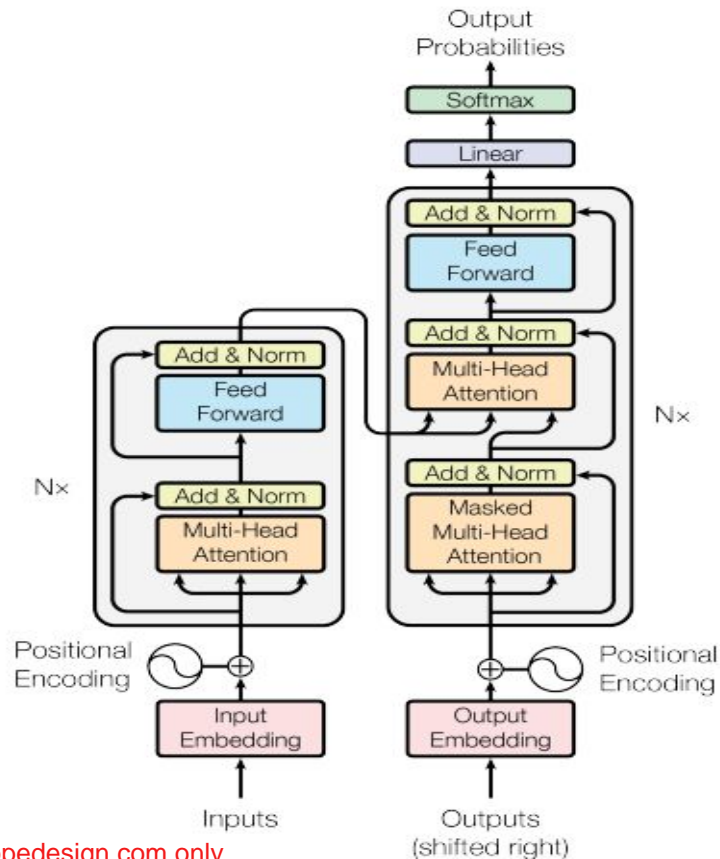
What are Transformer Models?

Transformers are a **type of neural network architecture**

Transformers were **introduced** in a paper by **Vaswani et al. in 2017**

Transformers are based on the idea of **self-attention**

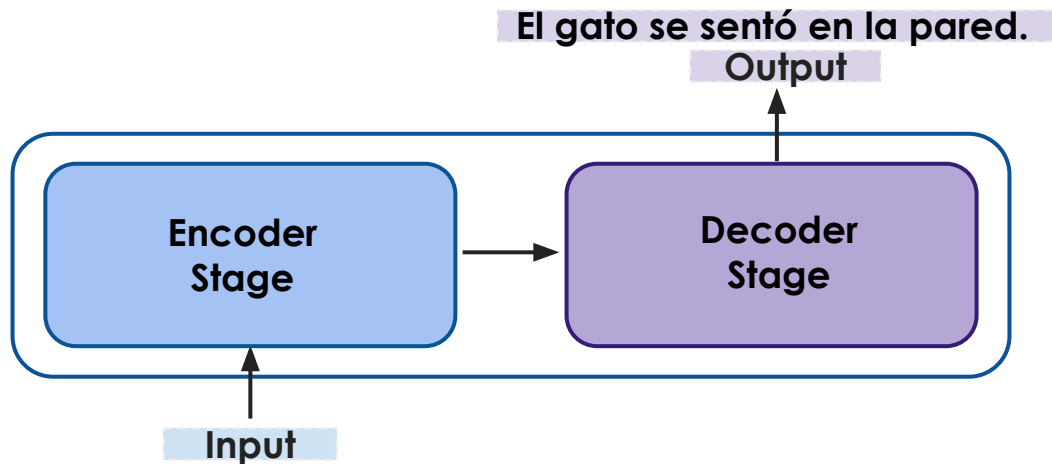
Source: Image from the original research paper [Attention Is All You Need](#)



How does a Transformer model work?

The **encoder** takes in a sequence of tokens (e.g. words or characters) and outputs a **latent representation**

The **decoder** then takes this latent representation as input and outputs a **sequence of tokens**



The cat sat on the wall.

El gato se sentó en la pared.

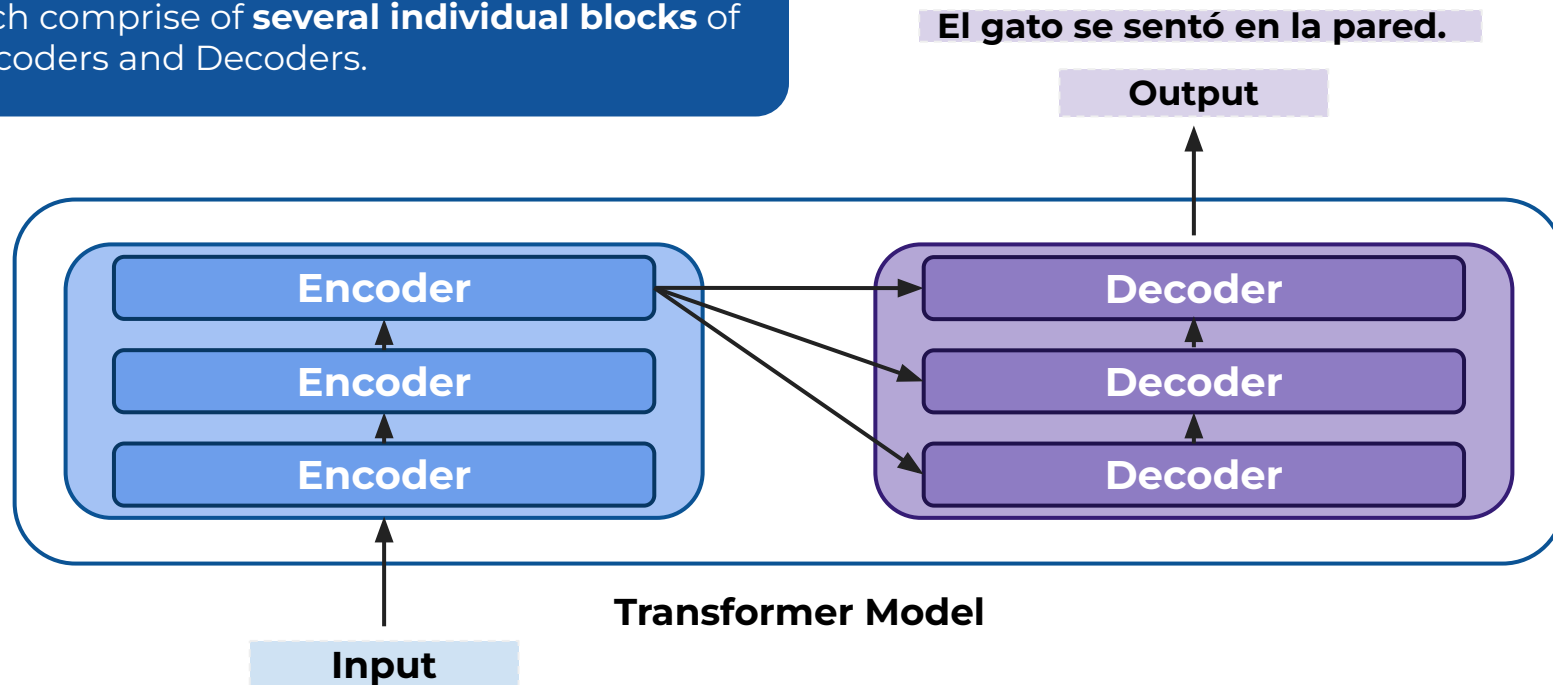
This file is meant for personal use by tinlong@iscopedesign.com only.

Sharing or publishing the contents in part or full is liable for legal action.

Proprietary content. © Great Learning. All Rights Reserved. Unauthorized use or distribution prohibited.

The Transformer Model - High-level Flow

In reality, the **Encoder** and **Decoder** stage each comprise of **several individual blocks** of Encoders and Decoders.



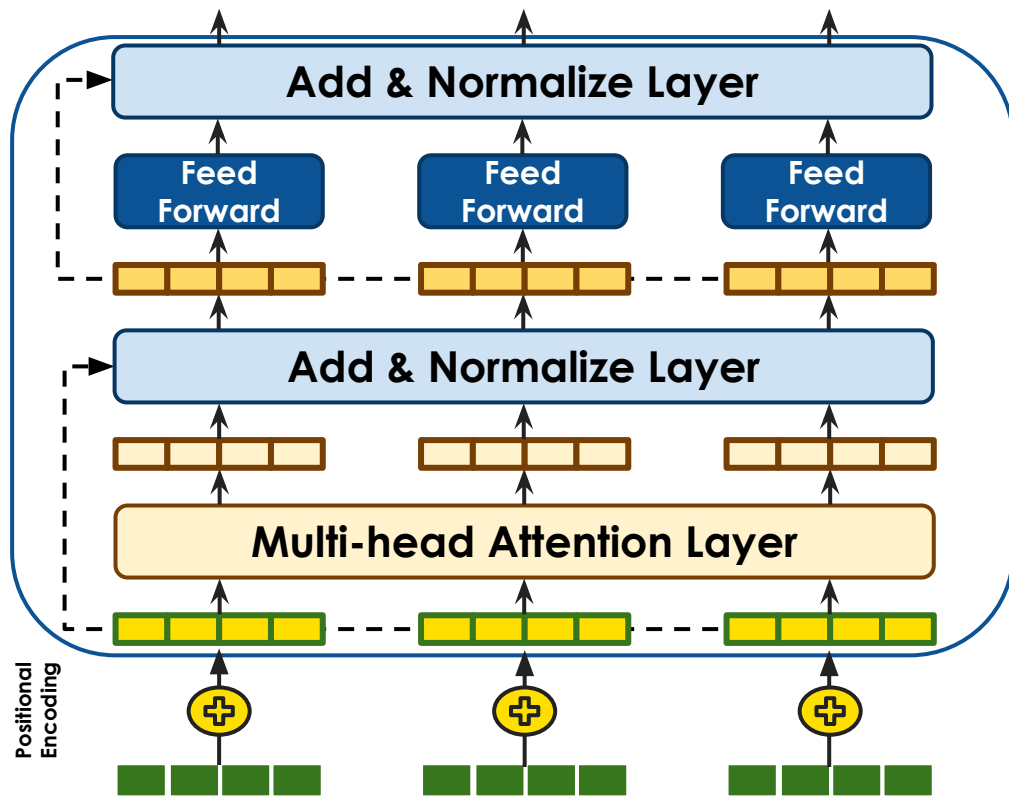
The cat sat on the wall.

El gato se sentó en la pared.

What are the components of the Encoder?

The Encoder block of a Transformer architecture consists of the following components:

1. **Multi-head Attention:** A stack of self-attention layers that allows the Encoder to attend to different parts of the input sequence simultaneously.
2. **Feedforward Neural Network:** Processes the outputs of the Multi-head Attention layer using a standard fully connected neural network with activations like ReLU.
3. **Skip Connections and Layer Normalization:** Improves the flow of information through the Encoder and avoids the vanishing gradient problem. These are added after each sub-layer.
4. **Positional Encoding:** Typically added to the input embeddings of the Encoder to provide positional information for words, using a set of learned sinusoidal functions.



This file is meant for personal use by tinlong@iscopedesign.com only.

Sharing or publishing the contents in part or full is liable for legal action.

Proprietary content. © Great Learning. All Rights Reserved. Unauthorized use or distribution prohibited.

How does the Self-Attention Mechanism work?

The **self-attention mechanism** lies at the **core** of **transformer models**

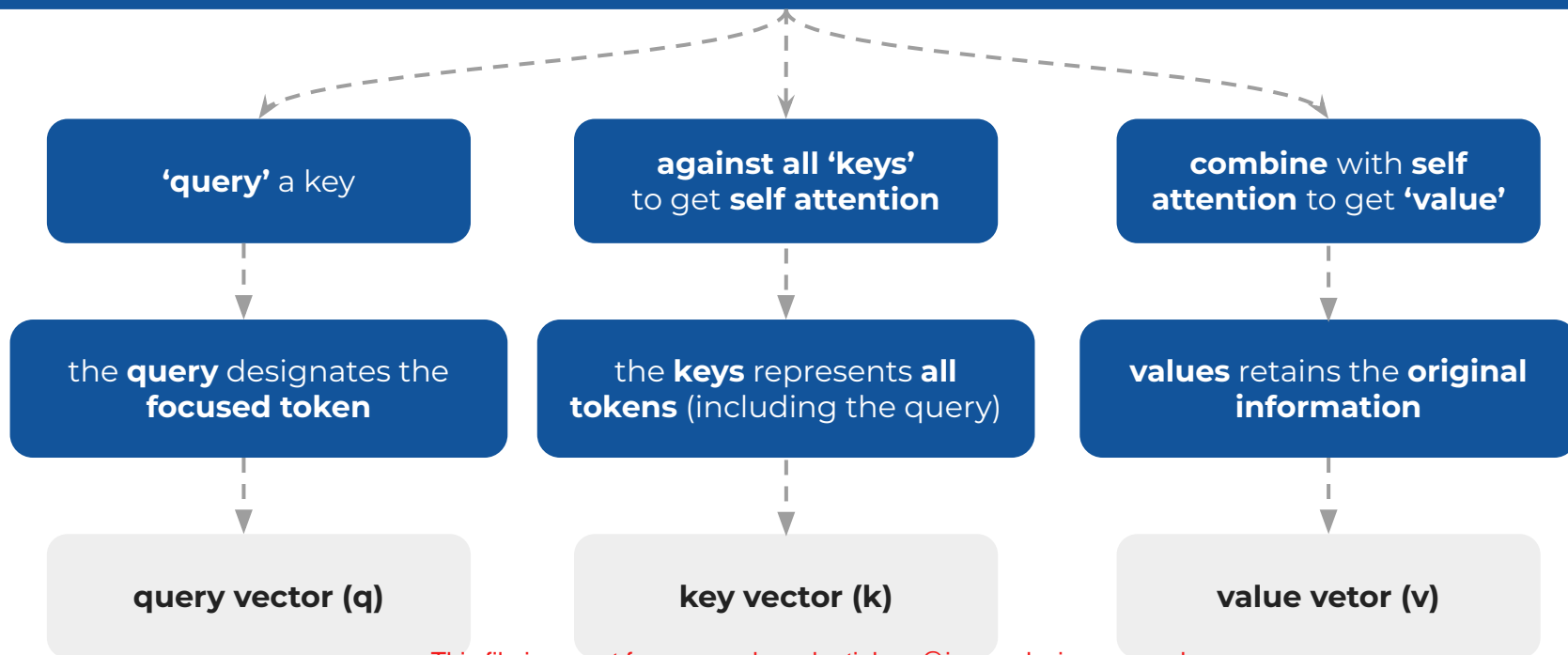
Self attention allows us to generate a **context-aware representation** of **each token** in the input

The **context-aware representation** of **each token** is generated with respect to all other tokens in the input

The context-aware representation **focuses** on the **relevant parts of the input** for a given task

Self Attention - Computation

The **steps** to get the **context-aware representation** for **each** of the **token** in the **input** are



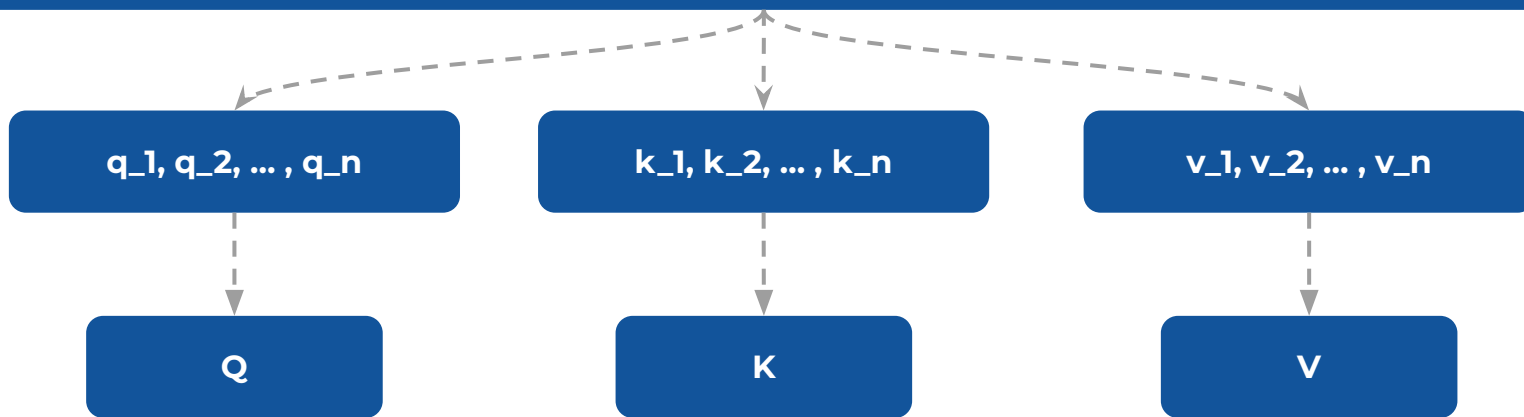
This file is meant for personal use by tinlong@iscopedesign.com only.

Sharing or publishing the contents in part or full is liable for legal action.

Proprietary content. © Great Learning. All Rights Reserved. Unauthorized use or distribution prohibited.

Self Attention - Computation

We **stack** these query, key, and value **vectors** into **matrices**



d_k here refers to the dimension of the vectors used for representing the input

$$\text{softmax} \left(\frac{Q * K^T}{\sqrt{d_k}} \right) * V$$

context-aware representations of all tokens

What are the components of the Decoder?

Most of these operations in the Decoder are identical to the Encoder.

Self-Attention

Add & Normalize Layer

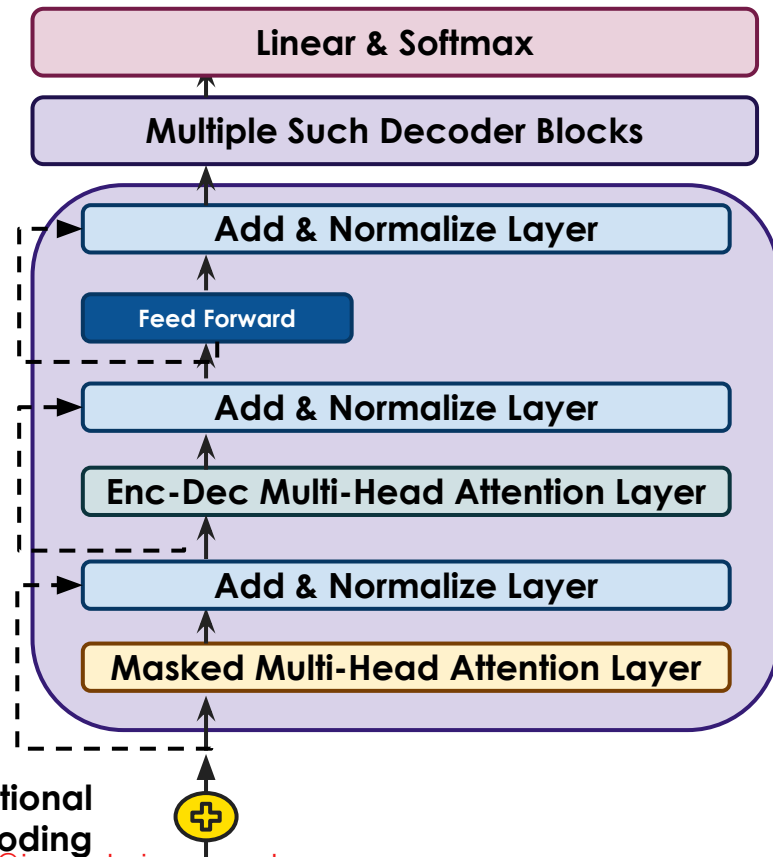
Feed Forward

But there are a few other operations **unique to the Decoder**.

Masking

Encoder-Decoder Attention Layer

Linear & Softmax



Components of the Decoder

Masking

Involves **hiding** (masking) **information** from the **future** to keep the model focused on the present and past during each run

Encoder-Decoder Attention

Aligns the **decoder's output** with the **context** provided by the **encoder** by enabling selective attention to different segments of the input sequence

Linear & Softmax

Linear layer **converts** contextual information into a **format suitable** for subsequent computations

Softmax produces **probability scores**, facilitating the selection of the most likely output token

What are the different transformer architectures?

There are broadly three-types of transformer models today, based on their usage of Encoder and Decoder blocks

Encoder-Decoder

Encoder-only

Decoder-only

This file is meant for personal use by tinlong@iscopedesign.com only.

Sharing or publishing the contents in part or full is liable for legal action.

Proprietary content. © Great Learning. All Rights Reserved. Unauthorized use or distribution prohibited.

Encoder-Decoder Transformer

There are broadly three-types of transformer models today, based on their usage of Encoder and Decoder blocks

Encoder-Decoder

Utilize the Encoder and Decoder blocks in tandem, similar to the original transformer architecture

Encoder-only

Typically used in **tasks** where the **output heavily relies on the input**, like **Machine Translation** and **Text Summarization**

Decoder-only

Examples: **T5** and **FLAN-T5**

This file is meant for personal use by tinlong@iscopedesign.com only.

Sharing or publishing the contents in part or full is liable for legal action.

Proprietary content. © Great Learning. All Rights Reserved. Unauthorized use or distribution prohibited.

Encoder-only Transformer

There are broadly three-types of transformer models today, based on their usage of Encoder and Decoder blocks

Encoder-Decoder

Encoder-only

Decoder-only

Utilize only Encoder blocks to generate continuous embeddings from the input

Typically used in **discriminative tasks** that require embeddings, like for **Text Classification** and **Semantic Search**

Examples: **BERT** and **DistilBERT**

Decoder-only Transformer

There are broadly three-types of transformer models today, based on their usage of Encoder and Decoder blocks

Encoder-Decoder

Encoder-only

Decoder-only

Utilize only Decoder blocks to auto-regressively predict* the next token based on the input

Typically used in **generative tasks** like **Sentence Completion** and **Question-Answering**

Examples: **GPT** and **Llama**

* Autoregressive prediction involves predicting future values based on past values.



Happy Learning !

