

A **large language model** is a type of artificial intelligence algorithm that applies neural network techniques with lots of parameters to process and understand human languages or text using self-supervised learning techniques. Tasks like text generation, machine translation, summary writing, image generation from texts, machine coding, chat-bots, or Conversational AI are applications of the Large Language Model.

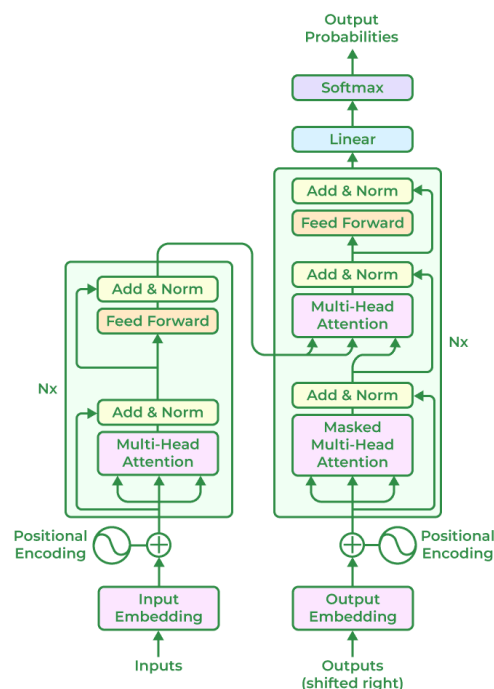
Large Language Model's (LLM) architecture is determined by a number of factors, like the objective of the specific model design, the available computational resources, and the kind of language processing tasks that are to be carried out by the LLM. The general architecture of LLM consists of many layers such as the feed forward layers, embedding layers, attention layers. A text which is embedded inside is collaborated together to generate predictions.

Important components to influence Large Language Model architecture:

- Model Size and Parameter Count
- input representations
- Self-Attention Mechanisms
- Training Objectives
- Computational Efficiency
- Decoding and Output Generation

### Transformer-Based LLM Model Architectures

[Transformer](#)-based models, which have revolutionized natural language processing tasks, typically follow a general architecture that includes the following components:



1. **Input Embeddings:** The input text is tokenized into smaller units, such as words or sub-words, and each token is embedded into a continuous vector representation. This embedding step captures the semantic and syntactic information of the input.
2. **Positional Encoding:** Positional encoding is added to the input embeddings to provide information about the positions of the tokens because transformers do not naturally encode the order of the tokens. This enables the model to process the tokens while taking their sequential order into account.
3. **Encoder:** Based on a neural network technique, the encoder analyses the input text and creates a number of hidden states that protect the context and meaning of text data. Multiple encoder layers make up the core of the transformer architecture. Self-attention mechanism and feed-forward neural network are the two fundamental sub-components of each encoder layer.
  1. **Self-Attention Mechanism:** Self-attention enables the model to weigh the importance of different tokens in the input sequence by computing attention scores. It allows the model to consider the dependencies and relationships between different tokens in a context-aware manner.
  2. **Feed-Forward Neural Network:** After the self-attention step, a feed-forward neural network is applied to each token independently. This network includes fully connected layers with non-linear activation functions, allowing the model to capture complex interactions between tokens.
4. **Decoder Layers:** In some transformer-based models, a decoder component is included in addition to the encoder. The decoder layers enable autoregressive generation, where the model can generate sequential outputs by attending to the previously generated tokens.
5. **Multi-Head Attention:** Transformers often employ multi-head attention, where self-attention is performed simultaneously with different learned attention weights. This allows the model to capture different types of relationships and attend to various parts of the input sequence simultaneously.
6. **Layer Normalization:** Layer normalization is applied after each sub-component or layer in the transformer architecture. It helps stabilize the learning process and improves the model's ability to generalize across different inputs.
7. **Output Layers:** The output layers of the transformer model can vary depending on the specific task. For example, in language modeling, a linear projection followed by SoftMax activation is commonly used to generate the probability distribution over the next token.

It's important to keep in mind that the actual architecture of transformer-based models can change and be enhanced based on particular research and model creations. To fulfill different tasks and objectives, several models like GPT, BERT, and T5 may integrate more components or modifications.

### Popular Large Language Models

Now let's look at some of the famous LLMs which has been developed and are up for inference.

- **GPT-3:** GPT 3 is developed by OpenAI, stands for Generative Pre-trained Transformer 3. This model powers ChatGPT and is widely recognized for its ability to generate human-like text across a variety of applications.
- **BERT:** It is created by Google, is commonly used for natural language processing tasks and generating text embeddings, which can also be utilized for training other models.
- **RoBERTa:** RoBERTa is an advanced version of BERT, stands for Robustly Optimized BERT Pretraining Approach. Developed by Facebook AI Research, it enhances the performance of the transformer architecture.
- **BLOOM:** It is the first multilingual LLM, designed collaboratively by multiple organizations and researchers. It follows an architecture similar to GPT-3, enabling diverse language-based tasks.

For implementation details, these models are available on open-source platforms like Hugging Face and OpenAI for Python-based applications.

### Large Language Models Use Cases

- **Code Generation:** LLMs can generate accurate code based on user instructions for specific tasks.
- **Debugging and Documentation:** They assist in identifying code errors, suggesting fixes, and even automating project documentation.
- **Question Answering:** Users can ask both casual and complex questions, receiving detailed, context-aware responses.
- **Language Translation and Correction:** LLMs can translate text between over 50 languages and correct grammatical errors.
- **Prompt-Based Versatility:** By crafting creative prompts, users can unlock endless possibilities, as LLMs excel in one-shot and zero-shot learning scenarios.

Use cases of LLM are not limited to the above-mentioned one has to be just creative enough to write better prompts and you can make these models do a variety of tasks as they are trained to perform tasks on one-shot learning and zero-shot learning methodologies as well. Due to this only Prompt Engineering is a totally new and hot topic in academics for people who are looking forward to using ChatGPT-type models extensively.