

---

# VARIATIONS OF LARGE LANGUAGE MODELS: AN OVERVIEW

---

**Md Shoaibur Rahman**  
Capital One - Quantica  
San Francisco, CA, USA  
shoaibur.rahman@capitalone.com

**Eddy Borera**  
Capital One - Quantica  
McLean, VA, USA  
eddy.borera@capitalone.com

## ABSTRACT

Large Language Models (LLMs) have introduced a paradigm shift in how computers understand and generate language, elevating artificial intelligence to unprecedented levels. These powerful models, anchored in the Transformer architecture, demonstrate unparalleled proficiency in tasks like text interpretation and generation. This review delineates LLMs into three predominant categories: Interpretive AI (iAI) for contextual understanding, Generative AI (gAI) for text generation, and Transformative AI (tAI) that combines both capabilities. For each category, we delve into their mechanisms, diverse applications, and the ethical challenges they introduce. By tracing their evolutionary journey, we provide insights into the rapid progression of linguistic technology and speculate on its future trajectory. As LLMs gain traction, a balanced understanding of their advantages and inherent challenges becomes paramount. This report presents a comprehensive overview of LLMs, highlighting their significance and implications for the future of AI.

## 1 Introduction

The field of Natural Language Processing (NLP) has witnessed a monumental shift in recent years, transitioning from rule-based systems to most recent advanced models known as Large Language Models (LLMs). This transition was made possible by enhanced computational power and the abundance of textual data available today, enabling machines to achieve feats once thought implausible in understanding and generating human language [1].

LLMs, with their vast parameters and intricate architectures, have set new standards in linguistic tasks, ranging from predicting words to crafting well-formed paragraphs. Models such as GPT-4 (Generative Pre-trained Transformer) [2] (and its siblings GPT-3.5 [3] and GPT-3 [4]), BARD [5], BERT (Bidirectional Encoder Representations from Transformers) [6], and many others available on open platforms like Hugging Face [7] represent the capabilities and promise of LLMs in modern-day applications.

As LLMs spread industries like healthcare [8] and finance [9], understanding their capabilities, limitations, and inherent biases is crucial. Their potential applications, from customer service automation to aiding medical diagnoses, emphasize their transformative impact on various sectors.

This review offers a deep dive into the realm of LLMs, highlighting their classifications, the innovations that drive them, and the most recent models that exemplify their evolving capabilities. Through this, we aim to present a holistic view of the present state of LLMs and the future they are set to shape.

## 2 Historical Context

The roots of Natural Language Processing (NLP) extend to the mid-20th century, when the journey began to make machines understand and generate human language. These nascent stages saw rule-based systems dominate, where language tasks were governed by manually-defined rules. The ELIZA program from the 1960s exemplifies this era, simulating a psychotherapist using pattern-matching and substitution methodologies [10].

By the 1990s, the growth in computational resources and the availability of expansive datasets instigated a pivotal shift in NLP. Rule-based models started giving way to statistical ones, with approaches like Hidden Markov Models (HMMs)

Table 1: List of major LLM milestones

Model	Full Name	Developer	Year	Category
Transformer	Transformer	Google	2017	tAI
BERT	Bidirectional Encoder Representations from Transformers	Google	2018	iAI
DistilBERT	Distilled version of BERT	HuggingFace	2019	iAI
T5	Text-to-Text Transfer Transformer	Google	2019	tAI
XLNet	Extended Transformer Network	Google	2019	iAI
RoBERTa	Robustly optimized BERT approach	Meta AI	2019	iAI
BART	Bidirectional and Auto-Regressive Transformers	Meta AI	2019	tAI
GPT-2	Generative Pre-trained Transformer 2	OpenAI	2019	gAI
CTRL	Conditional Transformer Language Model	Salesforce	2019	gAI
GPT-3	Generative Pre-trained Transformer 3	OpenAI	2020	gAI
LaMDA	Language Model for Dialogue Applications	Google	2021	gAI
MT-NLG	Megatron-Turing Natural Language Generation	NVIDIA/Microsoft	2021	gAI
BLOOM	BigScience Large Open-science Open-access Multilingual Language Model	BigScience	2022	gAI
PaLM	Pathways Language Model	Google	2022	gAI
GATO	Generalist Agent	Google	2022	gAI
chatGPT	GPT-3.5 variant for chat applications	OpenAI	2022	gAI
chatGLM	Chat version of General Language Model	Tsinghua	2023	gAI
Claude	Claude	Anthropic	2023	gAI
Falcon	Falcon	TIUAE	2023	gAI
BARD	Bard	Google	2023	gAI
LLaMA	Large Language Model Meta AI	Meta AI	2023	gAI
GPT-4	Generative Pre-trained Transformer 4	OpenAI	2023	gAI

and Conditional Random Fields (CRFs) gaining traction. These models tapped into data to identify linguistic patterns, reducing the dependency on manual rules [11].

The late 2000s and early 2010s signaled the neural network era in NLP. Innovations like Word2Vec [12] and GloVe [[13]] converted words into vectors, capturing their semantic essence. Deep learning structures, such as RNNs (Recurrent Neural Networks), LSTMs (Long Short-Term Memory), and GRUs (Gated Recurrent Unit), powered sequence-based tasks, marking significant progress in areas like language or text translation and sentiment analysis.

A monumental breakthrough came in 2017 with the Transformer architecture, spotlighted in the "Attention is All You Need" paper [14]. Departing from recurrent layers, Transformers employed attention mechanisms, setting the stage for pioneers like BERT and GPT. This introduced the era of Large Language Models (LLMs) that we are familiar with today.

In this context, recent innovations have further elevated NLP’s capabilities. Models like GPT-4 [2] and BARD [5] have seamlessly integrated text and image processing. ChatGPT [3] emphasizes natural conversations, while BLOOM [15] and LaMDA [16] showcase multilingual proficiencies and engaging dialogue capabilities. The range of LLMs continues to expand with models like MT-NLG [17], LLaMA [18], GATO [19], Claude [20], Falcon [21], and ChatGLM [22], each contributing distinct features to the dynamic landscape of NLP.

Refer to Table 1 for an overview of each model, highlighting significant LLM milestones organized by release year, developer, and respective category. A detailed rationale for this categorization can be found in the following sections.

### 3 Methodology

The expansive and diverse world of Large Language Models (LLMs) necessitates a structured approach to understand and categorize them effectively. In this section, we detail the logic underlying our categorization and the research strategies supporting our conclusions.

#### 3.1 Criteria for LLM Categorization

Our categorization is based on understanding the following functional aspects of LLMs:

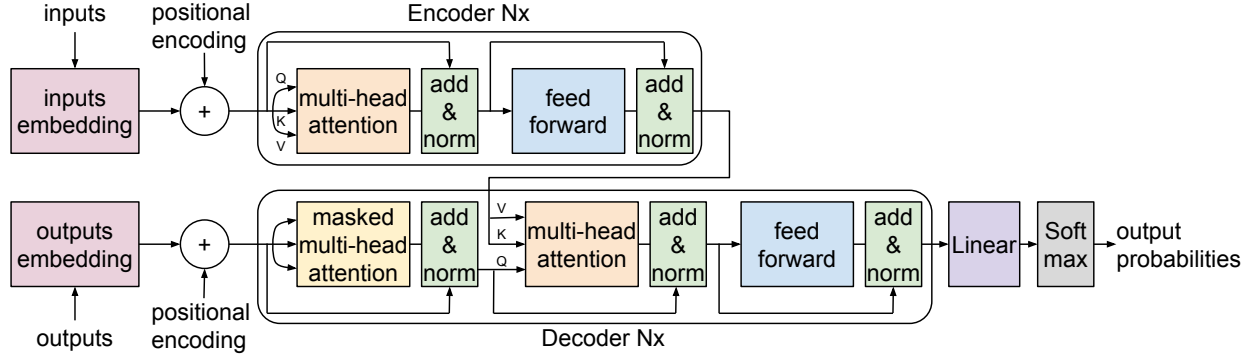


Figure 1: A visual representation of the Transformer architecture, highlighting the interplay of attention mechanisms, feed-forward networks, and positional encodings. Reproduced from [14].

1. **Architectural Composition:** A foundational classification, we dissect models based on their architectural core — whether they are primarily encoder-based, decoder-based, or a hybrid of both.
2. **Functional Role:** Stemming from their architecture, models are categorized by their primary functionality, be it understanding text (Interpretive AI), generating text (Generative AI), or transforming text (Transformative AI).
3. **Application Spectrum:** The variety of tasks LLMs excel at, from sentiment analysis (typically encoder models) to content generation (largely decoder models), offer critical insights for our categorization.

### 3.2 Research Methods

Our insights are the result of a comprehensive research strategy:

1. **Literature Review:** A thorough exploration of existing literature, spanning original papers, articles, and conference contributions, gave us historical context, current trends, and potential future directions.
2. **Model Analysis:** Practical assessments of prominent models like BERT, GPT-3, T5, and others provided direct insights into their capabilities and limitations.
3. **Expert Engagement and Community Feedback:** Interaction with NLP experts and the wider community through seminars, forums, and discussions enriched our theoretical insights with practical experiences and considerations.

Our methodology is both iterative and holistic, ensuring a comprehensive understanding of LLMs while firmly grounding our insights on empirical analysis and community expertise.

## 4 Transformer Architecture: The Backbone of LLMs

The Transformer architecture has been the cornerstone behind the success of modern Large Language Models (LLMs). In this section, we delve into the details of this architecture, emphasizing its components and functionalities. The section concludes with an overview of the Transformer’s pivotal role in reshaping the landscape of NLP.

### 4.1 A Glimpse of the Transformer

The current expertise of LLMs attributes much to the revolutionary Transformer architecture. Introduced in the seminal work by Vaswani et al. [14], this architecture discarded traditional recurrent neural structures, adopting instead the powerful attention mechanisms and marking a transformative era in natural language processing.

Figure 1 illustrates the Transformer’s constituents: the encoder and the decoder, each with attention mechanisms at their core [See Appendix A for how attention mechanisms work]. Let’s delve deeper into these components:

## 4.2 Data Flow and Architectural Details

- **Input embeddings:** Input words are converted into tokens size of  $n$ . Convert these tokens into embeddings using a learned embedding matrix. The resulting matrix has dimensions  $n \times d$ , where  $d$  is the embedding size. For each position in the sequence (from 0 to  $n - 1$ ), generate positional encodings of size  $n \times d$  [See Appendix B for methods of generating position encodings]. Add these encodings to the token embeddings to get input embeddings of size  $n \times d$ .
- **Encoder:** The encoder consists of several identical layers, denoted as  $N_x$  layers, which are stacked one after the other, and the output of one layer is passed to the next as its input. Each encoder layer comprises:
  - **Multi-head Self-Attention Mechanism:** The input embeddings are transformed into Query ( $Q$ ), Key ( $K$ ), and Value ( $V$ ) matrices. The attention scores are computed using the scaled dot product between  $Q$  and  $K$ . These scores signify the importance or emphasis each word places on every other word in the sequence. These raw scores are then passed through the softmax function to yield attention weights. The attention output is derived as the product of these weights and the  $V$  matrix. The process can be represented by:

$$\text{Attention}(Q, K, V) = \text{Softmax} \left( \frac{Q \cdot K^T}{\sqrt{d_k}} \right) V \quad (1)$$

where  $d_k$  denotes the dimension of the key. The attention outputs from all heads are concatenated and linearly transformed, then added to the original input embeddings using a residual connection. This is followed by layer normalization. The size of the resulting matrix is  $n \times d$ .

- **Feed-Forward Networks:** The normalized outputs are sent through feed-forward networks, which can have one or more hidden layers. The output from these networks is then combined with the input through a residual connection, followed by layer normalization. The resultant output size is  $n \times d$ .
- **Target embeddings:** Target words are converted into tokens of size  $m$ . Convert these tokens into embeddings using a learned embedding matrix. The resulting matrix has dimensions  $m \times d$ , where  $d$  is the embedding size. For each position in the sequence (from 0 to  $m - 1$ ), generate positional encodings of size  $m \times d$ . Add these encodings to the token embeddings to produce target embeddings of size  $m \times d$ .
- **Decoder:** The decoder, like the encoder, is made up of several identical layers (denoted as  $N_x$ ). Each decoder layer comprises:
  - **Masked Multi-head Self-Attention:** The self-attention mechanism in the decoder mirrors that of the encoder. However, there are two key differences: the  $Q$ ,  $K$ ,  $V$  matrices are derived from target embeddings, and the attention is masked to ensure that each token can only attend to prior tokens in the sequence. This masking mechanism is represented as:

$$\text{Attention}(Q, K, V) = \text{Softmax} \left( \frac{Q \cdot K^T}{\sqrt{d_k}} + \text{Mask} \right) V \quad (2)$$

In this equation, the Mask matrix contains extremely negative values for positions that should not be attended to, such as future tokens in an autoregressive setting. After passing through the softmax function, these positions approach zero, effectively ignoring those positions in the weighted sum.

- **Encoder-Decoder Attention:** This layer of attention enables the decoder to focus on different parts of the encoder's output. The  $Q$  matrix is derived from the output of the previous masked self-attention layer, whereas  $K$  and  $V$  are taken from the encoder's output. The attention output is computed using equation 1. By allowing the decoder to refer back to the encoder's output, the model can utilize the context provided by the encoder. This mechanism is pivotal for tasks such as machine translation where the decoder needs to generate sequences based on the encoder's context.
- **Feed-Forward Networks:** After the attention mechanisms, the output is passed through feed-forward neural networks. The output from these networks is then combined with the input from the previous layer using a residual connection and then normalized. The final output size remains  $n \times d$ .
- **Final Layers:**
  - **Linear Layer:** The output from the final decoder layer (size  $n \times d$ ) passes through a linear layer that has a weight matrix of size  $d \times W$ , where  $W$  is the vocabulary size. This transforms the output to size  $n \times W$ .
  - **Softmax Activation:** Applied to the above output, producing a probability distribution over the vocabulary for each position in the output sequence. The result remains  $n \times W$ .
  - **Final Output:** A matrix of size  $n \times W$  representing predicted probabilities for each token in the output sequence at each position.

### 4.3 The Transformative Impacts of the Transformer

The Transformer has been pivotal in numerous NLP breakthroughs:

- **Parallelization:** Avoiding recurrent layers, Transformers empower parallel data processing, expediting training and accommodating larger models.
- **Versatility:** Suited for both encoding (as with BERT) and decoding (as with GPT) tasks, and even combined approaches (e.g., T5, BART), its adaptability covers a broad NLP spectrum.
- **Performance Benchmarking:** Transformers have raised the bar in various NLP tasks, from machine translation to sentiment analysis.

In conclusion, the Transformer, with its pioneering attention mechanisms and modular design, has established itself as the cornerstone of modern LLMs, driving NLP research into novel areas.

## 5 LLM Spectrum: A Thorough Categorization

In the rapidly evolving domain of Natural Language Processing (NLP), models can be primarily grouped into three significant categories: Interpretive AI (iAI), Generative AI (gAI), and Transformative AI (tAI). While iAI concentrates on understanding and interpreting input data, gAI is dedicated to generating new, original content. Meanwhile, tAI, as the name suggests, focuses on transforming input data into another form, often combining aspects of both interpretation and generation. This section dives deep into the mechanics, comparative analyses, and insightful evaluations of models across these categories, highlighting their strengths, distinctions, and potential challenges.

### 5.1 Interpretive AI (iAI)

#### 5.1.1 Mechanics of Encoder-based Models

Interpretive AI primarily relies on the encoder component of the Transformer architecture. An encoder processes input sequences, such as sentences, producing a corresponding contextual representation. By operating through self-attention mechanisms, the encoder dynamically assigns attention to parts of the input based on their contextual significance.

#### 5.1.2 Comparative Analysis of iAI Models

Encoder-based models like BERT have gained popularity for their bidirectional processing of inputs to capture context [23]. RoBERTa, an optimized variant of BERT, demonstrates improved performance with larger datasets and extended training [24]. Conversely, DistilBERT provides a more compact solution by retaining only half of BERT's layers, leading to slightly reduced accuracy but faster performance [25]. XLNet, on the other hand, is a permutation-based approach, builds upon the foundational concepts of BERT but aims to overcome its limitations by leveraging all possible permutations of words in a sentence for training [26]. The choice between models often boils down to trade-offs in accuracy, speed, and computational efficiency.

#### 5.1.3 Drawbacks and Criticisms of iAI Models

iAI models, despite their capabilities, face several criticisms:

1. **Computational Intensity:** Particularly the larger variants are computationally intensive, requiring powerful GPUs or cloud infrastructures.
2. **Opaque Decision-making:** iAI models might produce answers without transparent rationale, creating challenges in critical applications.
3. **Bias:** These models can manifest biases from their training data, potentially leading to skewed results [27].
4. **Overfitting:** Their complexity creates risks of overfitting to training data, possibly affecting performance on new data.

Future advancements in NLP might address these current limitations.

## 5.2 Generative AI (gAI)

### 5.2.1 Mechanics of Text Generation

The decoder component of the Transformer architecture is central to generative AI, enabling the generation of coherent and contextually appropriate content. Utilizing attention mechanisms, the decoder weighs the relevance of each token in the input sequence, ensuring generated text remains semantically coherent.

### 5.2.2 Comparative Analysis of gAI Models

Recent LLM advancements include models like GPT-4 [2], ChatGPT [3], BARD [5], BLOOM [15], LaMDA [16], LLaMA [18], InstructGPT [28], Claude [20], Falcon [21], and ChatGLM [22]. These models have surpassed predecessors, with innovations such as GPT’s decoder-centric architecture, and CTRL’s controlled generation [29].

### 5.2.3 Ethical Concerns and Implications of gAI

Generative models, with their capability to produce authentic-sounding content, bring forth a host of ethical challenges:

1. **Misinformation and Fake News:** gAI models can be exploited to produce misleading information or entirely fabricated news articles, posing challenges for identifying truth in the digital age [30].
2. **Impersonation:** These models can be employed to craft realistic sounding emails or messages, potentially being used in phishing attacks or other forms of deception.
3. **Bias and Stereotyping:** If trained on biased data, gAI can perpetuate harmful stereotypes, leading to further misinformation or cultural insensitivity.
4. **Loss of Originality:** As content generation becomes increasingly automated, there’s a philosophical debate about the erosion of human creativity and the value of machine-generated art or literature.

The exponential growth of gAI demands rigorous ethical guidelines and countermeasures to prevent misuse while harnessing its vast potential.

## 5.3 Transformative AI (tAI)

### 5.3.1 Balancing Encoding and Decoding

Transformative AI leverages the strengths of both encoders and decoders within the Transformer architecture. The encoder reads and interprets the input sequence to generate contextualized representations. These representations are then fed into the decoder, which takes the context and produces a relevant and coherent output. This synergy is especially useful in tasks that require understanding of an input sequence and subsequently generating an associated output, such as in translation or summarization.

The key to tAI’s success lies in maintaining a harmonious balance between interpretation (encoding) and generation (decoding). While the encoder must capture the nuances and intent of the input, the decoder should generate outputs that are not only contextually aligned but also fulfill the desired task’s objectives. The interplay between these components, orchestrated through attention mechanisms and shared embeddings, ensures the generated outputs are both meaningful and accurate.

### 5.3.2 Comparative Analysis of tAI Models

The transformative AI paradigm, encapsulating both encoder and decoder mechanisms, has evolved considerably over time. While the foundational Seq2Seq (Sequence-to-Sequence) models were seminal for tasks like translation, capturing input sequences and transforming them into desired outputs [31], and T5 (Text-to-Text Transfer Transformer) showcased its versatility from translation to summarization [32], the real game-changers in the recent era have been models like MT-NLG [17], GATO [19], and PaLM [33]. MT-NLG, being a transformative AI, has redefined how encoder-decoder models function. GATO, with its multimodal capabilities, not only deals with text but also handles varied tasks like image captioning. Similarly, PaLM’s capability in handling a many tasks highlights its unique versatility. Amidst these, BART presents a unique approach, corrupting the input text and optimizing its reconstruction, excelling in tasks such as text generation and summarization [34].

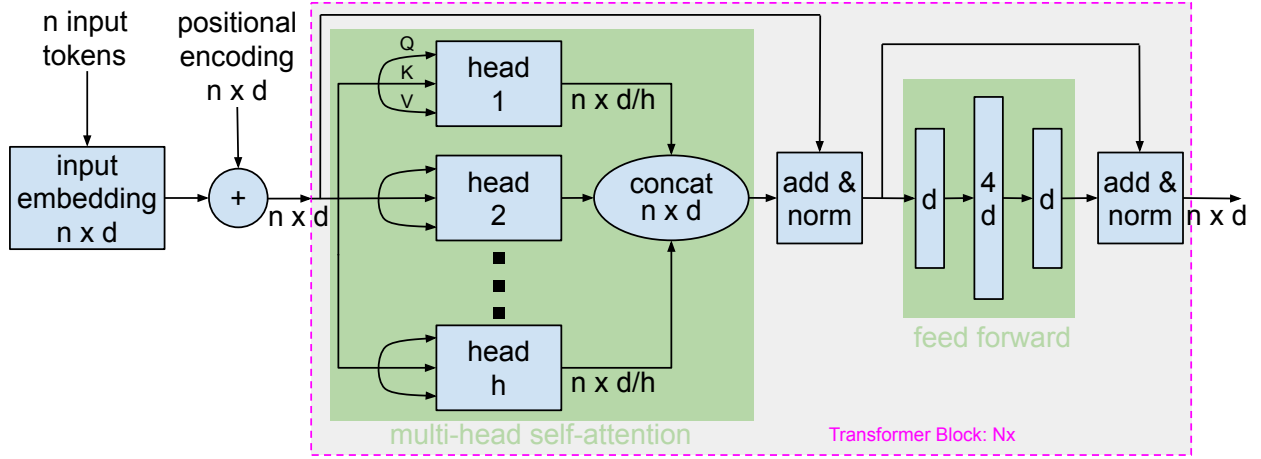


Figure 2: A visual representation of an encoder-based model, showcasing its layers, attention mechanisms, and bidirectional nature. Parameterized by  $n$ : number of input tokens,  $d$ : embedding size,  $h$ : number of heads in the attention mechanisms,  $N_x$ : number of encoding transformer blocks. Both input and output layers of the feed-forward network are of size  $d$ , while the hidden layer size is typically  $4d$ . In the context of the base-BERT model,  $d = 768$ ,  $h = 12$ , and  $N_x = 12$ .

### 5.3.3 Current Perspectives on tAI Limitations

As of this writing, the discourse around Transformative AI is less saturated with criticisms compared to its iAI and gAI counterparts. This can be attributed to its nascent stage in the research community, the specific applications it addresses, or a limited exploration of its potential pitfalls. Future research may shed more light on the broader implications and challenges associated with tAI.

## 6 In-Depth Analysis of Foundational Models

In this section, we deep dive into some foundational models in each category: BERT from interpretive AI (iAI), GPT-3 from generative AI (gAI) and T5 from transformative AI (tAI). We carefully selected the representative model from each category to provide a in-depth understanding of what happens behind the scenes when using a model from any of the three categories.

### 6.1 BERT for iAI

BERT, a hallmark of Interpretive AI models, primarily attributes its success to its bidirectional processing capability, which comprehends context from both left and right sides of a token in a sentence. This bidirectional ability set new benchmarks, such as in text classification on the AG News dataset [35], and showcased robust performance in sentiment analysis tasks like the Stanford Sentiment Treebank [36].

#### 6.1.1 BERT Architecture and Data Flow

The base BERT model architecture is depicted in Figure 2. Key details include:

- **Model Size:** BERT's base architecture contains 12 layers (transformer blocks), 768 hidden units, and 12 attention heads, summing up to approximately 110 million parameters.
- **Input Representations:** BERT's inputs are token sequences from texts, further complemented by special tokens: [CLS] (for classification tasks) and [SEP] (to demarcate segments). These embeddings combine token, segment, and position embeddings.
- **Transformer Blocks:** BERT's core consists of transformer blocks, with each featuring a multi-head self-attention mechanism and a feed-forward neural network. Surrounding each main component are residual connections followed by layer normalization.
- **Attention Mechanism:** BERT's attention mechanism, bidirectional in nature, focuses variably on different tokens. With 12 attention heads in the base model, each word has 12 unique attention weight sets.

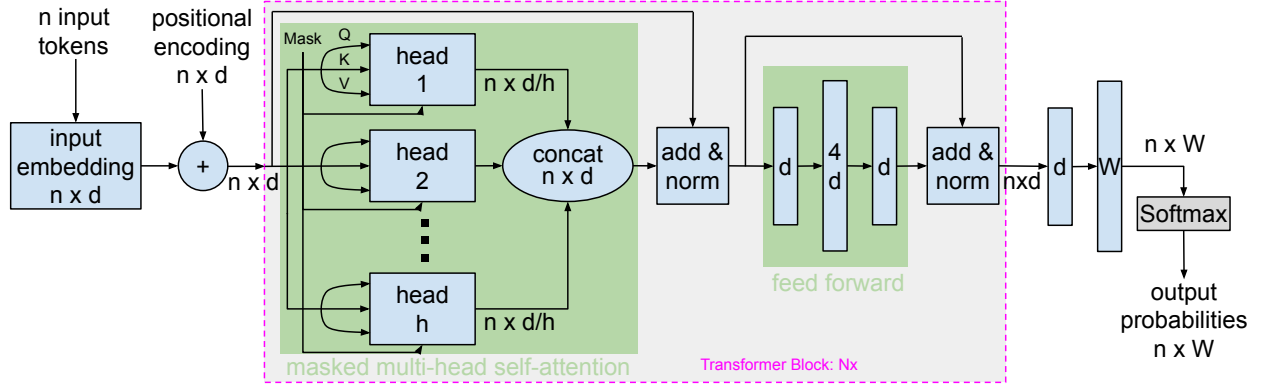


Figure 3: A visual representation of a decoder-based model, showcasing its layers, attention mechanisms, and generative capabilities. Parameterized by  $n$ : number of input tokens,  $d$ : embedding size,  $h$ : number of heads in the attention mechanisms,  $W$ : vocabulary size, and  $N_x$ : number of decoding transformer blocks. Both input and output layers of the feed-forward network are of size  $d$ , while the hidden layer size is typically  $4d$ . Vocabulary size can vary between 32,000 to 50,000. In the context of the GPT-3 model,  $d = 2048$ ,  $h = 16$ , and  $N_x = 96$ .

- **Position-wise Feed-forward Networks:** Present in each transformer block, these networks treat each position separately. Both the input and output layers match the input embedding size, while the hidden layer is  $4x$  larger. GELU (Gaussian Error Linear Unit) serves as the activation function.
- **Output:** The [CLS] token's output is suitable for classification, while outputs related to specific tokens can be used for tasks like named entity recognition.

### 6.1.2 BERT Pre-training and Fine-tuning

- BERT's pre-training capitalizes on the Masked Language Model (MLM) and Next Sentence Prediction (NSP). In MLM, BERT masks and then predicts roughly 15% of sentence words. For NSP, BERT detects if a second sentence naturally succeeds the first. Pre-trained BERT models are available from open source platforms such as Hugging Face.
- These models can be further refined for specific tasks by introducing task-specific layers to the pre-trained model.
- Constructing a BERT equivalent from scratch requires considerable resources, expertise, and time.

BERT's bidirectional context comprehension laid the groundwork for achieving unparalleled results across diverse NLP tasks. This innovation paved the way for subsequent BERT variants like DistilBERT, RoBERTa, and ALBERT.

## 6.2 GPT for gAI

Generative AI, particularly models like OpenAI's GPT-4 [2] and Google's BARD [5], has experienced significant advancements in recent years. These models have garnered considerable attention due to their remarkable generative abilities. GPT-3, a precursor in the series, is known for generating near-human text and has been implemented in platforms like GitHub Copilot for suggesting code snippets. A variant, ChatGPT or GPT-3.5, has been fine-tuned for conversational AI, driving the development of interactive bots and virtual assistants with human-like interactivity.

### 6.2.1 GPT Architecture and Data Flow

Figure 3 showcases the architecture of one of GPT-3's versions. The model's additional details include:

- **Model Size:** This version of GPT-3 consists of 96 layers (transformer blocks), 2048 hidden units, and 16 attention heads per block, totaling approximately 175 billion parameters.
- **Input Representations:** GPT-3 employs token and positional embeddings for input data representation. These embeddings help capture the sequence and context of the input.
- **Transformer Blocks:** GPT-3's core lies in its transformer blocks, which encompass a multi-head self-attention mechanism and feed-forward neural networks. Each component features residual connections followed by layer normalization.



- **Attention Mechanism:** The masked multi-head self-attention mechanism enables GPT-3 to attend to preceding words for next-word prediction. This attention processes text uni-directionally, capturing past context. This is a clear distinction from the BERT model, where attention is bidirectional, enabling it to understand the context from both the preceding and following words in a sentence.
- **Position-wise Feed-forward Networks:** Each block has feed-forward networks operating per position. These networks use the GELU activation function and maintain a hidden layer size four times that of the input embedding.
- **Output:** The last transformer block's outputs are directed through a linear layer, matching the input embedding size. This output is subsequently passed through a softmax function to yield a vocabulary-spanning probability distribution, determining the next predicted word.

### 6.2.2 GPT Pre-training and Fine-tuning

- Initially, GPT is trained on a vast corpus, aiming to predict the next word in a sequence. This process enables the model to grasp the language's inherent patterns and structures.
- Unlike many models which require task-specific fine-tuning post pre-training, GPT uniquely generalizes across multiple tasks. By understanding detailed prompts, GPT can efficiently execute diverse tasks without additional fine-tuning.
- GPT can operate in zero-shot, one-shot, or few-shot scenarios by incorporating examples within the input prompt to perform tasks. In this context, "zero-shot" refers to performing tasks without any provided examples, "one-shot" involves using a single example to understand a task, and "few-shot" leverages several examples to understand and execute a task.
- Given GPT's complexity, creating an equivalent model from scratch would require substantial resources.

GPT's architectural design and vast parameter count have propelled it to achieve unparalleled performance across diverse tasks, often rivaling or outperforming task-specific models. It has sparked interest in the possibilities and limits of large AI models.

## 6.3 T5 for tAI

Blending interpretation and generation, tAI models cater to diverse application domains. For instance, T5, by redefining translation as a text-to-text challenge, has marked its prominence in tasks requiring instant multilingual translations. The primary idea behind T5 is to treat every NLP problem as a text-to-text problem. Such a framing allows even classification tasks to be perceived as predicting a label in textual form. This versatility extends to applications like document summarization where the aim is to retain the essence of the content in a condensed format. Figure 4 depicts the architecture of the T5-11B model, and the subsequent sections provide a detailed exploration of its functionalities.

### 6.3.1 T5 Architecture and Data Flow

The illustration in Figure 4 showcases the architecture of the most expansive version of T5, named T5-11B. Key features and components of this model include:

- **Model Size:** T5 is available in multiple sizes, but the largest, T5-11B, boasts 11 billion parameters, enabling it to set benchmarks in numerous NLP tasks.
- **Input Representations:** T5 employs a tokenization strategy that breaks text into subwords, which are then translated into embeddings. Positional embeddings are combined with these token embeddings to ensure sequence coherence.
- **Encoder:** Comprising 24 transformer blocks, the encoder captures intricate context from the input. The multi-head self-attention mechanism within these blocks allows for the careful assessment of token significance when encoding. Additionally, each transformer block embeds a feed-forward network for processing token representations.
- **Decoder:** Initiated with a <start> token (e.g., for text summarising tasks) or target output tokens (e.g., for translation tasks), the decoder leverages embeddings analogous to the encoder. The decoder houses its own suite of transformer blocks, each equipped with self-attention and cross-attention mechanisms to ensure contextually relevant decoding. Feed-forward networks are also embedded within these blocks to further process token data.

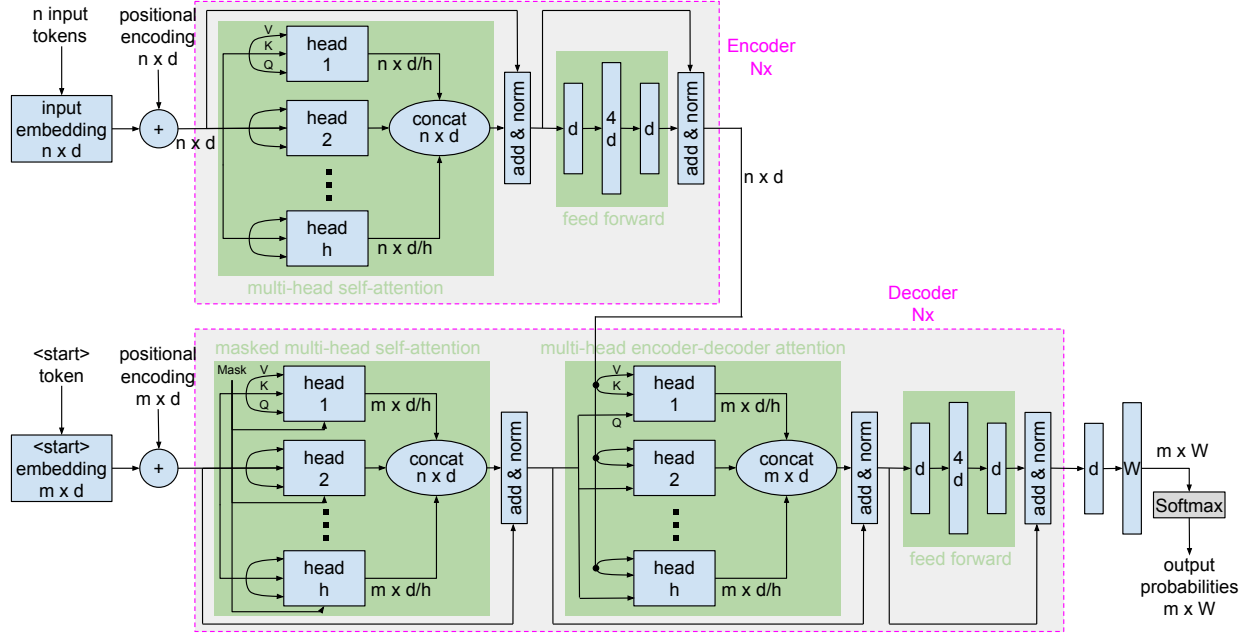


Figure 4: A visual representation of an encoder-decoder-based model, showcasing its layers, attention mechanisms, and text-to-text transformation capabilities. Parameterized by  $n$ : number of input tokens,  $m$ : number of target output tokens,  $d$ : embedding size,  $h$ : number of heads in the attention mechanisms,  $W$ : vocabulary size, and  $N_x$ : number of transformer blocks for both encoder and decoder. Both input and output layers of the feed-forward networks are of size  $d$ , while the hidden layer size is typically  $4d$ . Vocabulary size can vary between 32,000 to 50,000. In the context of the T5 model,  $d = 4096$ ,  $h = 128$ , and  $N_x = 24$ .

- **Output:** Post decoding, the output undergoes transformation through linear layers and a softmax activation to derive word probabilities. Depending on these probabilities, word selections are made, with options to introduce variability for more dynamic outputs like in generative AI.

### 6.3.2 T5 Pre-training and Fine-tuning

- Initially, T5 is trained on a vast text corpus. By reimagining every NLP task as text-to-text during this phase, T5 learns to transform input text sequences into fitting output sequences, thereby gaining a broad and versatile linguistic understanding.
- Post pre-training, while most models undergo task-specific fine-tuning, T5's text-to-text paradigm simplifies this process. Every task, be it sentiment analysis or translation, is treated as a transformation of one text form to another, making T5 an adaptable model for a wide range of tasks.
- T5's capabilities extend to zero-shot, one-shot, and few-shot learning. Such capability, coupled with its foundational training and the text-to-text model, empowers it to interpret and execute tasks based on minimal examples or prompts.
- Building a model of T5's stature from the ground up demands significant resources, both computational and expertise-wise, attesting to the strides made in NLP and the investment in crafting such advanced models.

## 7 Quality Assessment of Generated Texts

This section focuses on the technical evaluation of LLM outputs, and outlines the methods and metrics to determine how good or reliable an LLM's generation is. While these methods help in achieving better quality outputs, they are also crucial in addressing some of the ethical concerns outlined in the next section.

- **Human Evaluation:** A subjective but essential method where human raters assess the quality, relevance, and coherence of generated content.

- **Entity and Topic Verification:** Ensure the presence of desired themes or entities using methods like topic modeling, named entity recognition, and keyword matching.
- **Avoidance of Undesired Outputs:** Curtail harmful or biased content through techniques such as fine-tuning, rule-based filtering, and user feedback. Automate the detection process with tools like sentiment analysis or other dedicated models.
- **Consistency and Coherence:** Ensure logical structure and non-contradictory content. Check embedding similarity across sentences or analyze connected components in entity-relationship graphs.
- **Diversity Measurement:** Prevent monotonous AI outputs by evaluating the ratio of unique words to total words or by checking embedding similarity and clustering.
- **Quantitative Metrics:** These metrics provide an objective and standardized assessment of text quality, distinguishing them from the more subjective or heuristic evaluations.
  - **BLEU (Bilingual Evaluation Understudy):** BLEU checks how many n-grams in the generated content match those in the reference text. It favors precision and is more about "out of all the n-grams in my generated output, how many are correct?" The score also incorporates a brevity penalty, discouraging overly short or incomplete translations.
  - **ROUGE (Recall-Oriented Understudy for Gisting Evaluation):** ROUGE calculates the overlap of n-grams between generated and reference content. It favors recall and is more about "out of all the n-grams in the reference, how many did I capture in my generated output?" There are several variants of ROUGE, including ROUGE-N (overlap of n-grams), ROUGE-L (longest common subsequence), and ROUGE-S (skip-bigram).
  - **METEOR (Metric for Evaluation of Translation with Explicit Ordering):** More sophisticated than BLEU, METEOR considers precision and recall, synonymy, stemming, and even word order, and therefore, offers a more holistic evaluation.
  - **BERTScore:** It uses the power of BERT's contextual embeddings. By comparing the embeddings of the generated text and reference text, BERTScore captures semantic similarities that other metrics might miss. It provides precision, recall, and F1-scores. Typically, the F1-score is used as the representative BERTScore value if not stated otherwise.
  - **Perplexity:** It is a measure for probabilistic models like language models, and does quantify how well the model's probability distribution aligns with the true distribution of the words in the text. A lower perplexity indicates a better fit, meaning the model is less "surprised" or "perplexed" by the given text.

## 8 Ethical Considerations and Criticisms

The emergence of Large Language Models (LLMs) across various sectors highlights not just their technological breakthroughs but also introduces numerous ethical dilemmas. Given the considerable influence these models hold, a thorough examination of their ethical concerns and the feedback they receive is crucial.

### 8.1 Biases in LLMs

Among the most pronounced apprehensions tied to LLMs are the intrinsic biases they might encase. Addressing and detecting these biases requires robust quality assessment, as discussed in the previous section.

1. **Data-Originated Biases:** As LLMs undergo training on massive datasets typically derived from the internet, they are susceptible to inheriting societal biases present in this data. These biases range from racial and gender stereotypes to socio-economic disparities [27].
2. **Amplification of Biases:** LLMs, in some contexts, might not only mirror but intensify pre-existing biases. This can result in outputs that show a higher degree of bias than the training data [37].
3. **Presumption of Objectivity:** The algorithmic foundations of LLMs can mislead users into assuming their outputs are impartial, overlooking over the potential for biased information.

### 8.2 Potential Misuse and Mitigation Strategies

The capabilities of LLMs also give rise to scenarios for potential misuse:

1. **Fake News and Disinformation:** The generative ability of AI models can be misemployed to create convincing yet wholly spurious news articles, driving disinformation campaigns.

2. **Impersonation and Phishing:** Cutting-edge models can author communications that adeptly imitate genuine human-crafted content, thereby elevating risks in domains such as email phishing.
3. **Dependency and Over-reliance:** An unchecked reliance on LLMs, especially within decision-making fields, could erode human judgment and supervision. Critical sectors like healthcare, finance, and law increasingly utilize LLMs. However, relying solely on model outputs without human judgment can result in severe consequences. For instance, in medical diagnostics, LLM recommendations should complement, not replace, expert opinions. Similarly, financial or legal decisions based solely on model outputs might overlook contextual nuances, leading to potentially flawed outcomes. It's imperative for users to interpret LLM results with a thorough understanding, integrating them with domain-specific knowledge.

Below are some strategies that can be used to mitigate these potential misuse:

- **Fine-tuning on Curated Data:** Models trained on meticulously curated datasets can reduce biases, aligning outputs more faithfully with societal ideals and ethics [38].
- **Transparent Model Training:** Open-sourcing both the model training protocols and the datasets introduces transparency, facilitating the pinpointing of bias and inaccuracies. Furthermore, understanding attention mechanisms helps make the model more transparent by showing how it weighs different words. This allows us to see into its decision-making and find any bias or mistakes.
- **User Education:** By providing end-users with knowledge about LLMs' confines and inherent biases, a more thoughtful engagement with their outputs can be encouraged.
- **Regulation and Oversight:** Regulatory frameworks can demarcate boundaries for LLM utilization, establishing ethical commitments and limiting misuse.
- **Reinforcement Learning from Human Feedback (RLHF):** Incorporating RLHF into the mechanics of text generation enriches the output's quality and relevance. Once an initial model is trained, human evaluators can refine it by providing explicit feedback or even ranking multiple generated text outputs. This feedback fine-tunes the text generation capabilities of the model. Moreover, RLHF addresses challenges by mitigating biases, upholding contextual appropriateness, and ensuring safety measures in the generated content.

In summary, although LLMs show transformative prospects, they are not without ethical challenges. It is incumbent upon both the NLP community and the broader society to confront and navigate these challenges, aiming for responsible development, enhancement, and deployment of these models.

## 9 Future Prospects

The rapidly expanding domain of Large Language Models (LLMs) has revolutionized our approach to textual data processing and interpretation, marking an exciting course for the evolution of natural language processing. Though LLMs bear intrinsic challenges, they equally carry the promise of novel advancements. This section explores the present challenges faced by LLMs, potential enhancements, and the foreseeable trends in their development and utility.

### 9.1 Current Limitations and Areas for Improvement

**Model Size and Computational Overheads:** The descriptor "Large" in LLMs is indicative of their scale. Models, exemplified by GPT-3, encompass billions of parameters, inducing significant computational burdens. Such size constraints not only impede deployment in environments with limited resources but also instigate concerns regarding energy usage and environmental sustainability [39].

**Fine-tuning and Domain Adaptation:** While LLMs demonstrate proficiency in broad language comprehension, tasks rooted in specific domains typically necessitate targeted fine-tuning. Present fine-tuning methodologies occasionally fall short in efficiency, calling for enhanced strategies to adapt models to specialized areas [4].

**Bias and Ethical Challenges:** The biases inherent to training datasets can skew LLM outputs, leading to ethical dilemmas and potential misinterpretations. Rectifying these biases is vital for upholding the principles of fairness and trustworthiness [40].

### 9.2 Anticipated Future Developments

**Leaner and More Efficient LLMs:** To counter the hefty computational requisites of contemporary LLMs, there's an expected trajectory towards devising models that reconcile robust performance with fewer parameters and minimized

training datasets [25]. NVIDIA’s FasterTransformer offers optimized transformer-based encoder and decoder components for NLP tasks [41]. Meanwhile, SmoothQuant introduces a post-training quantization method for LLMs, enabling efficient 8-bit weight and activation quantization while preserving accuracy [42]. Additionally, enhancing the quality of training data can lead to more accurate model predictions, reduced biases, and a better generalization to real-world scenarios.

**Enhanced Domain Adaptation:** With industries seeking customized solutions, LLMs are projected to incorporate superior domain adaptation mechanisms, providing specialized insights without the necessity of exhaustive fine-tuning.

**Expansive Multi-modal Applications:** Beyond the domain of text, the fusion of LLMs with visual and auditory information holds immense promise. Prototypes like OpenAI’s CLIP have already unveiled the capabilities of such integrative multi-modal applications [43].

**Ethical and Fairness Guidelines:** As LLM adoption becomes pervasive, we anticipate the emergence of rigorous ethical standards, possibly shaped by both communal consensus and regulatory oversight, to safeguard responsible utilization.

## 10 Conclusion

Large Language Models (LLMs) have played a key role in advancing Natural Language Processing (NLP) and the wider area of artificial intelligence. This review explored the details of LLMs, from their beginnings with the Transformer architecture to their use in many fields. We divided them into three groups: Interpretive AI (iAI), Generative AI (gAI), and Transformative AI (tAI), showing the unique strengths of each.

We highlighted how LLMs are now used in many areas involving context understanding and content creation. But, it is clear that powerful tools come with big responsibilities. We need to use LLMs carefully, making sure we do not misuse them or let them promote biases.

LLMs have a lot of impacts. Economically, they can change industries, improve operations, and create new business types. For society, they can make information easier to access, help with education, and bridge language gaps. At the same time, they push us to think differently about data, understanding, and AI ethics.

To wrap up, LLMs bring a lot of exciting opportunities for AI’s future. It is up to experts, businesses, and policymakers to make sure we move forward in the right way. The future for LLMs is not just about better algorithms or more data; it is about combining technology, ethics, and society in a balanced way.

## Appendix

### A Attention Mechanism in Transformers

#### A.1 Preliminaries

Given a set of input vectors  $X$ , the attention mechanism allows a model to focus on different parts of the input with varying intensity. This ‘attention’ is computed based on three main components:

- **Query ( $Q$ ):** It corresponds to the current word or token being considered. The Query matrix seeks to understand which parts (or words/tokens) of the input sequence are relevant in the context of the current word.
- **Key ( $K$ ):** The Key matrix is used to represent every word or token in the input sequence. It essentially provides a representation against which the Query matrix can be compared to determine the relevance of each word in the sequence.
- **Value ( $V$ ):** Once the relevance (or attention scores) of each word in the sequence is computed by comparing the Query with all Keys, the Value matrix provides the actual values (or word representations) that are used in the weighted sum to produce the final attention output for the current word.

These are derived from the input  $X$ .

## A.2 Computation Steps

- **Step 1: Obtain Q, K, V matrices:** First, we derive the  $Q$ ,  $K$ , and  $V$  matrices using learned weight matrices  $W_Q$ ,  $W_K$ , and  $W_V$  respectively.

$$Q = X \times W_Q \quad (3)$$

$$K = X \times W_K \quad (4)$$

$$V = X \times W_V \quad (5)$$

- **Step 2: Calculate Attention Scores:** To determine the attention scores, compute the dot product between the Query and Key matrices. For self-attention, both the query and key come from the same previous layer output.

$$\text{Scores} = Q \times K^T$$

- **Step 3: Scale the Attention Scores:** The scores are scaled down by a factor of the square root of the depth of the key vectors (often denoted as  $d_k$ ).

$$\text{Scaled Scores} = \frac{\text{Scores}}{\sqrt{d_k}}$$

This scaling helps in stabilizing the gradients, especially for larger values of  $d_k$ .

- **Step 4: Apply Softmax:** The scaled scores are then passed through a softmax operation to obtain the attention weights. This ensures the weights are between 0 and 1, and they sum up to 1.

$$\text{Attention Weights} = \text{softmax}(\text{Scaled Scores})$$

- **Step 5: Multiply by Value matrix:** To get the final attention output, multiply the attention weights with the Value matrix:

$$\text{Attention Output} = \text{Attention Weights} \times V$$

The attention mechanism provides a way for models to weigh input components differently, enabling them to focus or "attend" to more relevant parts of the input when processing data. It has been instrumental in the success of the Transformer architecture across a range of NLP tasks.

## B Positional Encoding in Transformers

The Transformer model, introduced by Vaswani et al. in "Attention is All You Need", lacks any inherent notion of the order or position of words. To allow the model to take into account the position of words in a sequence, the concept of *positional encoding* was introduced. This encoding provides additional information about the relative position of words in a sequence.

### B.1 Formulation

The positional encodings have the same dimension,  $d$ , as the embeddings, so that the two can be summed. They are defined as:

$$PE_{(pos, 2i)} = \sin\left(\frac{pos}{10000^{2i/d}}\right)$$

$$PE_{(pos, 2i+1)} = \cos\left(\frac{pos}{10000^{2i/d}}\right)$$

Where  $pos$  is the position and  $i$  is the dimension. This formulation provides a way for the model to easily learn to attend by relative positions since for any fixed offset  $k$ ,  $PE_{pos+k}$  can be represented as a linear function of  $PE_{pos}$ .

### B.2 Intuition

The intuition behind the sinusoidal choice is that it allows the model to easily determine the relative positions of words, even if they are far apart in a sequence. This is due to the unique properties of sine and cosine functions. The oscillating nature of these functions and the way they're formulated for positional encoding ensures that each position in a sequence gets a unique encoding.

### B.3 Usage in Transformers

The Transformer adds the positional encodings to the input embeddings at the bottom of the encoder and decoder stacks. The sum of the positional encodings and the embeddings are fed into the encoder and decoder stacks, allowing the model to consider the position of words in its computations.

### References

- [1] R. Schwartz et al. Green ai. *arXiv preprint arXiv:1907.10597*, 2019.
- [2] OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [3] OpenAI. Introducing chatgpt. <https://openai.com/blog/chatgpt>, 2022.
- [4] T. B. Brown et al. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.
- [5] Google. An overview of bard: an early experiment with generative ai. <https://ai.google/static/documents/google-about-bard.pdf>, 2023.
- [6] J. Devlin et al. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2019.
- [7] T. Wolf et al. Huggingface’s transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*, 2020.
- [8] E. Callaway. ‘it will change everything’: Deepmind’s ai makes gigantic leap in solving protein structures. *Nature* 588, 203–204, 2020.
- [9] J. M. Tshimula et al. Characterizing financial market coverage using artificial intelligence. *arXiv preprint arXiv:2302.03694*, 2023.
- [10] J. Weizenbaum. Eliza—a computer program for the study of natural language communication between man and machine. *Communications of the ACM*, pp.36–45, 1966.
- [11] M. Collins. Log-linear models, memms, and crfs. <https://api.semanticscholar.org/CorpusID:18582888>, 2011.
- [12] T. Mikolov et al. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.
- [13] J. Pennington et al. Glove: Global vectors for word representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages: 1532–1543, 2014.
- [14] A. Vaswani et al. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.
- [15] T. Le Scao et al. Bloom: A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*, 2022.
- [16] R. Thoppilan et al. Lamda: Language models for dialog applications. *arXiv preprint arXiv:2201.08239*, 2022.
- [17] S. Smith et al. Using deepspeed and megatron to train megatron-turing nlg 530b, a large-scale generative language model. *arXiv preprint arXiv:2201.11990*, 2022.
- [18] H. Touvron et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [19] S. Reed et al. A generalist agent. *arXiv preprint arXiv:2205.06175*, 2022.
- [20] Anthropic. Introducing claude. <https://www.anthropic.com/index/introducing-claude>, 2023.
- [21] G. Penedo et al. The refinedweb dataset for falcon llm: Outperforming curated corpora with web data, and web data only. *arXiv preprint arXiv:2306.01116*, 2023.
- [22] THUDM. Chatglm-6b. <https://github.com/THUDM/ChatGLM-6B>, 2023.
- [23] J. Devlin et al. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [24] Y. Liu et al. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- [25] V. Sanh et al. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019.
- [26] Z. Zhang et al. Xlnet for named entity recognition. *arXiv preprint arXiv:2002.00260*, 2020.

- [27] N. Mehrabi et al. A survey on bias and fairness in machine learning. *arXiv preprint arXiv:1908.09635*, 2019.
- [28] L. Ouyang et al. Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155*, 2022.
- [29] N. Keskar et al. Ctrl: A conditional transformer language model for controllable generation. *arXiv preprint arXiv:1909.05858*, 2019.
- [30] R. Zellers et al. Defending against neural fake news. *arXiv preprint arXiv:1905.12616*, 2019.
- [31] I. Sutskever et al. Sequence to sequence learning with neural networks. *arXiv preprint arXiv:1409.3215*, 2014.
- [32] C. Raffel et al. Exploring the limits of transfer learning with a unified text-to-text transformer. *arXiv preprint arXiv:1910.10683*, 2019.
- [33] A. Chowdhery et al. Palm: scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.
- [34] M. Lewis et al. Bart: denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461*, 2019.
- [35] X. Zhang et al. Character-level convolutional networks for text classification. In *Advances in Neural Information Processing Systems*, 2015.
- [36] R. Socher et al. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 conference on empirical methods in natural language processing*, 2013.
- [37] S. Paun et al. Comparing bayesian models of annotation. *Transactions of the Association for Computational Linguistics*, 2018.
- [38] T. Gebru et al. Datasheets for datasets. *arXiv preprint arXiv:1803.09010*, 2020.
- [39] N. Shazeer. Glu variants improve transformer. *arXiv preprint arXiv:2002.05202*, 2020.
- [40] I. Turc et al. Well-read students learn better: On the importance of pre-training compact models. *arXiv preprint arXiv:1908.08962*, 2019.
- [41] NVIDIA. Fastertransformer. <https://github.com/NVIDIA/FasterTransformer>, 2023.
- [42] G. Xiao et al. Smoothquant: Accurate and efficient post-training quantization for large language models. *arXiv preprint arXiv:2211.10438*, 2022.
- [43] A. Radford et al. Language models are unsupervised multitask learners. *Openai Blog*, 1(8):9, 2019.