# ChatGPT in the Age of Generative AI and Large Language Models: A Concise Survey

Salman Mohamadi[1], Ghulam Mujtaba[1], Ngan Le[2], Gianfranco Doretto[1], Donald A. Adjeroh[1*]

[1]Computer Science & Electrical Engineering, West Virginia University, Morgantown, WV, 26506, West Virginia, USA.
[2]Computer Science & Computer Engineering, University of Arkansas, Fayetteville, AR, 72701, Arkansas, USA.

*Corresponding author(s). E-mail(s): donald.adjeroh@mail.wvu.edu;
Contributing authors: sm0224@mix.wvu.edu; gmujtabakorai@gmail.com;
thile@uark.edu; gianfranco.doretto@mail.wvu.edu;

**Abstract**

ChatGPT is a large language model (LLM) created by OpenAI that has been carefully trained on a large amount of data. It has revolutionized the field of natural language processing (NLP) and has pushed the boundaries of LLM capabilities. ChatGPT has played a pivotal role in enabling widespread public interaction with generative artificial intelligence (GAI) on a large scale. It has also sparked research interest in developing similar technologies and investigating their applications and implications. In this paper, our primary goal is to provide a concise survey on the current lines of research on ChatGPT and its evolution. We considered both the glass box and black box views of ChatGPT, encompassing the components and foundational elements of the technology, as well as its applications, impacts, and implications. The glass box approach focuses on understanding the inner workings of the technology, and the black box approach embraces it as a complex system, and thus examines its inputs, outputs, and effects. This paves the way for a comprehensive exploration of the technology and provides a road map for further research and experimentation. We also lay out essential foundational literature on LLMs and GAI in general and their connection with ChatGPT. This overview sheds light on existing and missing research lines in the emerging field of LLMs, benefiting both public users and developers. Furthermore, the paper delves into the broad spectrum of applications and significant concerns in fields such as education, research, healthcare, finance, etc.

**Keywords:** Generative Artificial Intelligence, Large Language Model, NLP, ChatGPT
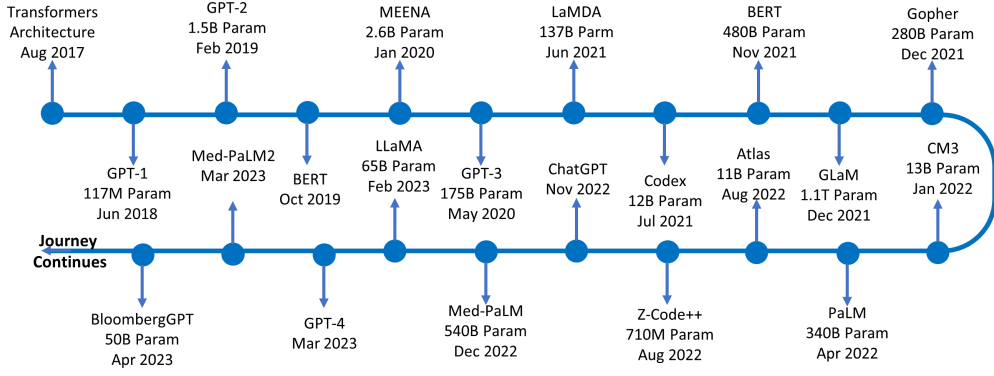
1

# 1 Introduction

The human ability to communicate through language is a remarkable capability that develops during childhood and evolves over a lifetime [1]. In contrast, machines cannot naturally grasp the nuances of human language, their understanding, or communication without powerful artificial intelligence (AI) algorithms. Developing machines that can read, write, and communicate like humans has been a long-standing research challenge [2].

One of the major methods for enhancing the language intelligence of machines involves language modeling (LM). LM is a technical process that seeks to model the probability of word sequences and predict missing or future words (also called tokens). The goal of LM is to enable machines to comprehend and generate human-like language by learning patterns and relationships between words in a given text corpus. This entails the use of advanced algorithms that can analyze vast amounts of data to determine the most likely sequence of words and accurately predict the next word in a sentence. This ability is crucial for natural language interactions with humans that allow machines to generate coherent and contextually appropriate responses. The most popular examples of language interaction with humans include chatbots, language translation, and speech recognition.

The field of LM has been the subject of extensive research, which can be classified into four main stages of development. The first stage is statistical language models (SLMs), which emerged in the 1990s and are based on statistical learning methods [3, 4]. SLMs rely on the Markov assumption to build word prediction models where the most recent context is used to predict the next word. N-gram language models are SLMs with a fixed context length, such as bigram and trigram languages. SLMs have been widely used to improve the performance of tasks in natural language processing (NLP) [5, 6] and information retrieval (IR) [7, 8]. However, SLMs often face the challenge of the curse of dimensionality and limited datasets, making it difficult to estimate high-order language models accurately. Consequently, smoothing techniques such as Good-Turing estimation [9] and backoff estimation [9] have been developed to tackle the issue of data sparsity.

The second stage is the development of NLP which is characterized by the use of neural language models (NLMs) [10]. It uses neural networks to model the probability of word sequences like recurrent neural networks (RNNs) [11]. A significant contribution of this stage was the introduction of the concept of distributed representation, which created word prediction functions based on aggregated context features (i.e., distributed word vectors) [10]. This approach was extended to develop a general neural network solution for various NLP tasks [12], including the development of the popular word2vec embedding method [13, 14]. A simplified shallow neural network is used to learn distributed word representations that are highly effective across different NLP tasks. This stage initiated the use of language models for representation learning which had a significant impact on NLP.

The pre-trained language models (PLMs) have become a major focus in the third stage of LM development/evolution [15]. These models use pre-training to generate context-aware word representations that can be utilized for downstream tasks such as
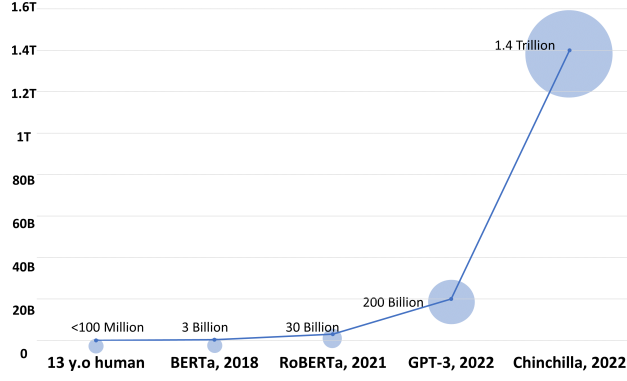
**Fig. 1** Brief timeline of a number of well-known large language models (LLMs).

Embeddings from Language Models (ELMo) and Bidirectional Encoder Representations from Transformers (BERT) [16]. ELMo [15] uses a bidirectional Long Short Term Memory (LSTM) [17] (or biLSTM) network to capture context-aware word representations by analyzing both preceding and following words, which allows for a deeper representation of word meaning. Meanwhile, BERT [16] utilizes a highly parallelizable transformer architecture [18] with self-attention mechanisms to pre-train bidirectional language models on large-scale unlabeled corpora to make them highly effective for many NLP tasks. The pre-trained context-aware word representations have significantly improved the performance of NLP tasks. It enables better performance in areas like question answering, sentiment analysis, and language generation[19]. Furthermore, the success of PLMs has led to the development of the *"pre-training and fine-tuning"* learning paradigm [20]. Here, the PLM is pre-trained on a large unlabeled dataset and then fine-tuned on a small labeled dataset for a specific downstream task. This approach has become a popular method for achieving state-of-the-art results in many NLP tasks [21].

Recent research has shown that scaling up PLMs in terms of model size or data size leads to improved model capacity for downstream tasks by following the scaling law [22]. This has led to a number of studies exploring the upper limits of performance by training even larger PLMs, such as the 175 billion-parameter GPT[1]-3 and the 540 billion-parameter PaLM. These larger PLMs behave differently when compared with the relatively smaller ones such as BERT with 330 million parameters or GPT-2 with 1.5 billion parameters. They have been observed to exhibit surprising emergent abilities that enable them to solve complex tasks that were previously thought to be impossible [23, 24]. For instance, GPT-3 can solve few-shot tasks through in-context learning while GPT-2 cannot [25]. These large-sized PLMs are now referred to as large language models (LLMs) [26]. One prominent application of LLMs is ChatGPT, which adapts LLMs from the GPT series for dialogue and demonstrates a unique ability to converse with humans [27]. Thus, the fourth stage of LM development can be seen as the introduction of these LLMs and subsequent exploration of the upper limits of

---

[1]GPT is the short form for "Generative Pre-trained Transformer."

**Fig. 2** Number of tokens seen during training. Adapted from [33].

their performance through the development of ever more extensive and more powerful models that exhibit emergent abilities and can solve even more challenging tasks. A brief timeline of LLMs is illustrated in Figure 1.
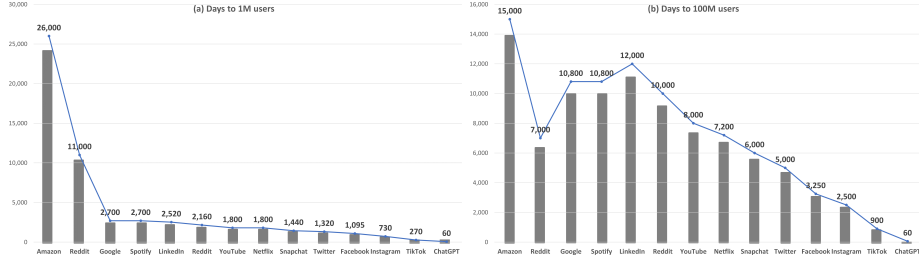
The emergence of such models which exhibit both surprisingly useful and complex capabilities boosted the attention toward a broader area of research, Generative AI (GAI). GAI, a branch of AI, involves the creation of computer models that can generate new and convincing content, such as images, music, or text (more details in Section 2.2). While the concept of GAI has gained significant attention in recent years, its origins can be traced back to the mid-20th century when early attempts at computer-generated art and music emerged [28]. Specifically, the fact emergence of LLMs, such as ChatGPT [29] and GPT-4 [30], has had a considerable effect on the field of AI and has prompted a renewed interest in both the prospects of artificial general intelligence (AGI) [31], and the potential for risks in the rapid development of GAI [32]. To improve performance these models are being trained on enormous datasets. Figure 2 illustrates the amount of training data for LLMs as compared to that seen by human teenagers and other LLMs in terms of the number of tokens [33]. The progress of LLMs has revolutionized research in various areas of AI, particularly in NLP, where LLMs have proven to be a general-purpose language task solver to a certain capacity. Consequently, research efforts are increasingly focusing on the use of LLMs. AI chatbots such as ChatGPT[2] and New Bing[3] are challenging traditional search engines in information retrieval (IR), and researchers are exploring the use of LLMs to improve search results [34]. Additionally, researchers are developing vision-language models similar to ChatGPT to facilitate multimodal dialogues [35]. GPT-4 [30] has also integrated visual information to support multi-modal input, highlighting the potential of this new wave of technology to enable a thriving ecosystem of real-world applications based on LLMs.

ChatGPT can generate responses in multiple languages and perform a range of tasks including question answering, creative writing, problem-solving across a vast range of disciplines, writing code, etc [36]. The popularity and reach of ChatGPT with
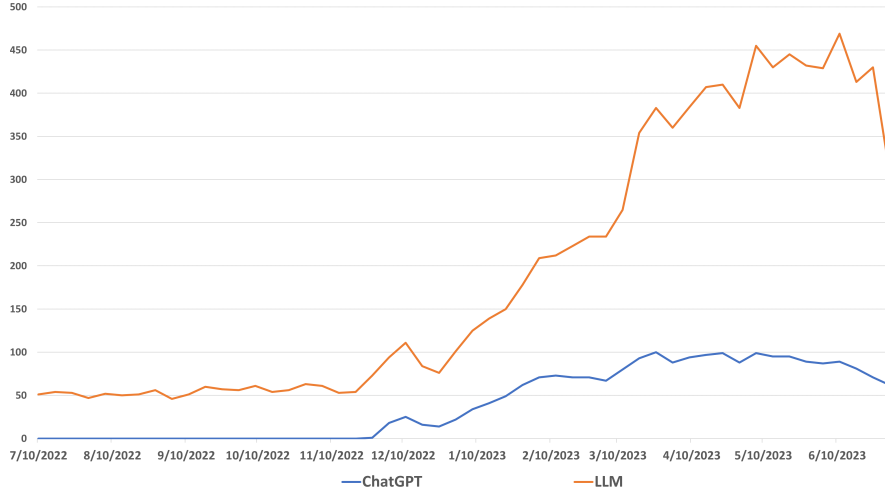
---

[2]https://chat.openai.com/
[3]https://www.bing.com/new

**Fig. 3** Some of the popular applications have reached 1 million (left) and 100 million active users (right) worldwide within days.



**Fig. 4** Comparison of ChatGPT and LLM on Google trends from June 10, 2022 to June 10, 2023.
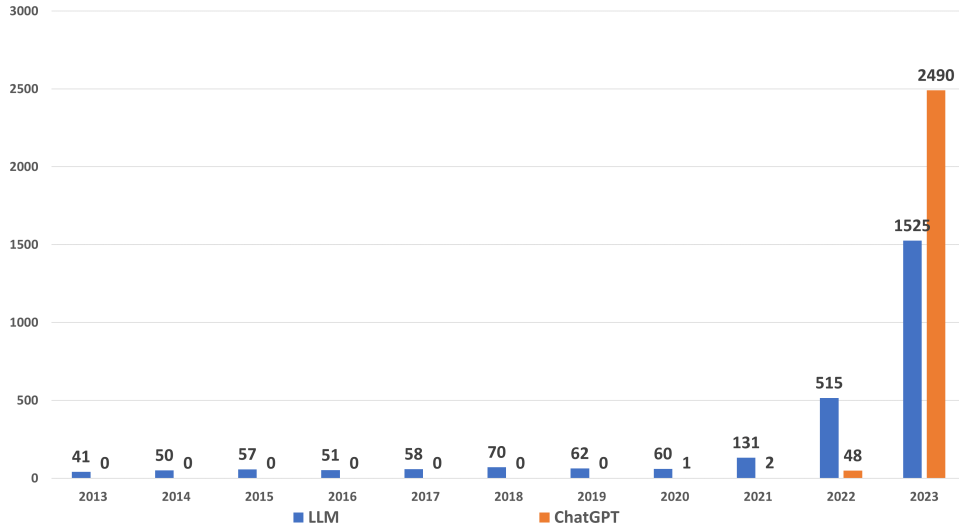
other popular applications are illustrated in Figure 3. The figure shows the time it took different platforms to reach 100 million global active users that are measured in months as reported in [37]. Figure 4 shows the rapid growth of ChatGPT and LLMs as captured by Google Trends from June 2022 to June 2023. Since its inception, ChatGPT has undergone significant evolution in terms of its capabilities, number of users, and in mitigating its limitations. Initially, ChatGPT was limited in its ability to generate coherent and contextually relevant responses. However, through continuous training on large datasets and the development of more sophisticated algorithms, ChatGPT has made significant progress in improving the quality of its responses.

This paper aims to provide a concise servery of the work on ChatGPT from a GAI and LLM perspective. It covers current research, limitations, applications, and concerns related to ChatGPT. Additionally, it highlights significant aspects and issues that are currently overlooked but have substantial implications. The goal of the paper is to promote the safe and effective adoption of ChatGPT while mitigating potential harm. It also discusses responsible AI practices, public responses, regulations, fairness,

privacy, and security. Finally, it outlines required studies and future research directions, including limitations of detectors, LLM-based search engines, societal impact, in-house LLMs training, effects on human languages, and potential future directions.

## 1.1 Related Surveys

The emergence of LLMs has transformed NLP by providing state-of-the-art results in various language tasks. ChatGPT is one of the excellent examples, as it demonstrated impressive capabilities in generating human-like language, understanding contextual information, and performing a wide range of conversational tasks. However, the development of LLMs is accompanied by both opportunities and challenges, requiring careful attention and consideration. In this regard, a range of survey papers have been published. Existing survey papers focus on specific or general aspects of LLMs [19, 26, 38] or ChatGPT [39–42]. In comparison, this paper aims to give an overview of the latest advancements in LLMs and put the main focus on the development of ChatGPT. Through a detailed analysis of the relevant literature, we summarize the most important issues, findings, techniques, and methods related to ChatGPT relying on the preliminary foundation.



**Fig. 5** Number of papers published on LLM and ChatGPT from 2013 to July 2023 was obtained through Google Scholar.

The graph presented in Figure 5 illustrates the number of academic papers that have been published on the subject matters of LLM and ChatGPT, based on Google Scholar records. The data was collected over a period of 10 years from 2013 to July 06, 2023. The focus of the study was on papers that included specific keywords such as "Large Language Model", "Large Language Models", "LLM", "LLMs", and "ChatGPT" in their titles. This approach allowed for a more targeted and accurate analysis

6

of the research data. For those who wish to delve deeper into this subject and retrieve the data, a code repository has been made available for their convenience [4].

## 1.2 Scope

The goal of this work is to present a clear and concise overview of the current research on ChatGPT. The scope of the survey includes foundational knowledge of LLMs and GAI, as well as a review of research across various fields such as education, healthcare, and finance. The paper also on approaches of ChatGPT explores emerging aspects of this technology, highlighting under-explored or missing areas of research. Specifically, the survey covers LLMs in the context of GAI, the evolution of ChatGPT, new sub-fields, research in terms of both glass-box and black-box views of ChatGPT, responsible AI in light of the emergence of ChatGPT, and challenges associated with this technology.

## 1.3 Terminology

We find it necessary to provide some preliminaries including the relevant terminology needed to follow the discussion in this survey. Here are a list of some key terms and their definitions:

- **Machine Learning (ML):** ML is a subset of AI that involves the use of algorithms and statistical models to enable computers to learn from data and make predictions or decisions without being explicitly programmed [43].
- **Deep Learning (DL):** DL is a subset of ML that uses artificial neural networks with multiple layers to learn from large datasets and make complex predictions or decisions. They have achieved state-of-the-art performance in several tasks such as image recognition and NLP [44].
- **Natural Language Processing (NLP):** NLP is a subfield of AI that involves the use of computational techniques to enable computers to understand, interpret, and generate human language. It involves tasks such as text classification, sentiment analysis, language translation, speech recognition, and manipulation of natural language data [45].
- **Tokenization:** The process of breaking up text into individual tokens, which are typically words or subwords, for use in LM [16].
- **Masked Language Modeling:** It is a task in NLP where certain words in a sentence are replaced with a special token. The goal is for the model to predict the original words based on context. This is often used for pre-training language models like BERT to better understand relationships between words and sentences [16].
- **Transformer:** A DL architecture that uses self-attention mechanisms to process sequential input data and capture dependencies between distant positions in the sequence [18].

---

[4]This resource is particularly useful for scholars and researchers who want to examine the quantity and distribution of scholarly articles related to a specific field or topic. You can access the code repository at https://github.com/iamgmujtaba/scholar_search.

- **Attention:** A mechanism used in DL architectures, including transformers, to selectively focus on different parts of the input data based on their relevance to the task at hand [18].
- **Encoder:** The part of a transformer architecture that processes the input sequence and generates a representation that summarizes the information in the sequence.
- **Decoder:** The part of a transformer architecture that generates an output sequence, one token at a time, using the representation generated by the encoder.
- **Embedding:** A representation of a discrete variable as a continuous vector in a high-dimensional space, which allows the model to capture semantic relationships between different tokens [13, 14].
- **Beam Search:** A search algorithm used in sequence generation tasks such as LM and machine translation, which explores multiple possible sequences in parallel and selects the most likely one based on a scoring function [46].
- **Perplexity:** A metric used to evaluate the quality of LMs, which measures how well a model can predict the probability distribution of the next word in a sequence given the previous words [47].
- **Pre-training of LMs:** The process of training a LLM on a vast corpus of text data in an unsupervised manner to learn general language representations that can be fine-tuned for downstream NLP tasks [20].
- **Fine-tuning of LMs:** The process of adapting pre-trained language models (PLMs) to a specific task or domain by training it on a smaller dataset of task-specific examples [48].
- **Conditional Generation**: The task of generating text that meets certain conditions or criteria, such as generating a response to a given input prompt in a chatbot system.
- **Prompt Engineering:** The process of effective communication with large-language models through the refinement of the input query or output recommendations, using input prompts of different types [49].
- **GPT:** A series of large pre-trained transformer-based language models developed by OpenAI, including GPT-1, GPT-2 [23], GPT-3, GPT-3.5, and GPT-4 [30].

## 2 Large Language Models (LLMs)

In this section, we present the characteristics of LLMs and the evaluation of various LLMs based on their primary mappings between input and output data types. Additionally, we address key concerns related to LLMs, including their implications for GAI.

### 2.1 Characterizing LLMs

LLMs like ChatGPT have revolutionized NLP. They excel in tasks like text completion, translation, and question answering. Here we describe work on LLMs in terms of various relevant topics namely data sources, architectures, pre-training, fine-tuning, scaling, optimization, human-in-the-loop interactions, and evaluation [50]. It is essential to comprehend the significance of these aspects to fully understand their potential and direct future research in this field.

- **Data Sources**: To effectively train LLMs, it is imperative to have access to a substantial amount of data sourced from diverse mediums such as books, articles, and websites. The data is carefully selected and processed to ensure quality and avoid biases. Data types used include program codes, however, concerns about data privacy, copyright issues, and the spread of harmful information can arise due to the massive amount of data used [51]. Detailed discussion on LLMs data source and data extraction is addressed in [52].
- **Architectures**: LLMs are created using advanced neural network architectures, namely transformers [18]. These transformers have greatly improved NLP by effectively processing and understanding the context of words in a sentence. These architectures consist of numerous attention and feed-forward neural network layers, which enable the model to recognize intricate patterns and dependencies in the data [18].
- **Pre-training and Pretasks**: To prepare for specific tasks, LLMs undergo pre-training [20] where they are exposed to vast amounts of unlabeled text data. This phase helps models to learn how to predict missing words or masked tokens within sentences, which builds their understanding of language. Consequently, they can generate relevant and coherent responses. This process is crucial in enabling AI-powered assistants to provide users with accurate and useful information.
- **Fine-tuning**: Fine-tuning is a crucial step in developing LLMs. It involves training the model on a smaller, more specific dataset that is tailored to the target task at hand [48]. This could be anything from text completion to translation or question-answering. Adapting the model to the specific task requirements can greatly improve its performance. Essentially, fine-tuning is a way to make the pre-trained model more specialized and efficient in its designated task.
- **Scaling**: Scaling LLMs involves various techniques [53]. One approach is to train them on wider and more diverse datasets. Another method involves increasing the number of parameters in the model, which enables it to recognize intricate patterns. However, this method demands significant computational resources. Parallel computing techniques and specialized hardware like GPUs or TPUs can also aid in computation scaling. Moreover, larger models have demonstrated better performance on language tasks, at times showing emergent characteristics but they may also have limitations and biases that arise during the training process [51].
- **Optimization**: Optimizing parameters and hyperparameters, such as the learning rate, regularization techniques, and loss functions are essential to attain optimal outcomes when training LLMs [54]. Optimizing these parameters can help the model converge to better solutions and improve its overall performance. As researchers continue to develop new optimization algorithms and techniques, the training process of LLMs becomes increasingly efficient and effective.
- **Human-in-the-Loop**: LLMs are often developed with the help of human interactions. This means that human reviewers and annotators are involved in the process of curating data, which is essential to ensure that the training data is of high quality and relevance. Additionally, human interactions are used to provide feedback, evaluate model outputs, and address ethical concerns, including bias, fairness, and

misinformation [55, 56]. Identifying and addressing these issues through interactions is vital for the successful development and deployment of LLMs.

- **Evaluation**: Evaluating the efficacy of extensive LLMs is crucial to gauge their proficiency across diverse tasks [57]. The metrics used to evaluate performance varies depending on the task but typically include accuracy, precision, recall, F1 score, and perplexity. Similar to general ML or NLP methods, it is imperative to conduct comprehensive evaluations to comprehend the constraints, prejudices, and potential hazards that come with extensive LLMs.

It is crucial to consider these when discussing LLMs. Each factor has its own intricacies and ongoing research aimed at overcoming obstacles and enhancing the capabilities of these models. Below, we present a taxonomy for different LLMs in the context of GAI.

## 2.2 Taxonomy for LLMs

Researchers have recently shown a growing interest in developing and improving LLMs using DL architectures. LLMs are trained on large textual datasets to generate coherent, sensible, and natural responses to natural language queries, and are also used in text generation systems for various applications. The success of LLMs has piqued the interest of researchers from both within and outside the field of computer science in creating artificial intelligence generated content (AIGC). This interest has been further fueled by the release of powerful LLMs from major companies such as Google[5], OpenAI[6], Microsoft[7], and Hugging Face[8]. While some of these LLMs are limited to a single modality, such as ChatGPT [58], others, such as GPT-4 [30], can process multi-modal data. AIGC is created using advanced GAI techniques, automated and distinct from human-generated content. For example, ChatGPT developed by OpenAI understands human input and provides meaningful responses through a textual modality. Until the release of GPT-4, ChatGPT was considered the most powerful conversational bot made available to the public [32]. In contrast, OpenAI's DALL-E2 generates high-quality images based on textual descriptions provided by humans.

To facilitate our analysis of various LLMs in the context of GAI, we first organized them into a taxonomy based on the primary mappings between input and output data types shown in Figure 6. It is important to note that the list is not comprehensive as it excludes many competitors, such as DeepMind, Amazon, EleutherAI, BigScience, Aleph Alpha, Huawei, Tsinghua, Together, Baidu, and many others. In the following sections, we provide an overview of the models that belong to each of the categories.

### 2.2.1 Text-to-Text Models

Text-to-Text LLMs are a class of DL models that have gained popularity in recent years due to their ability to transform one text into another. These models have shown excellent performance in various NLP tasks such as question answering, dialogue generation, summarization, and language translation, among others. One of the most
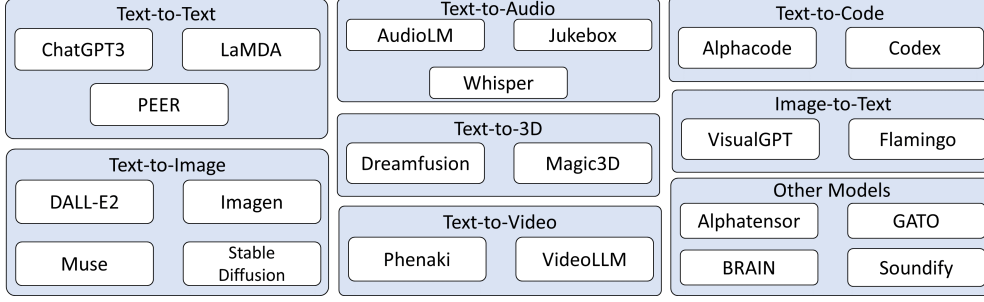
---

[5]https://bard.google.com/
[6]https://openai.com/dall-e-2
[7]https://aka.ms/Aajnnys
[8]https://huggingface.co/blog/bloom

**Fig. 6** A taxonomy of LLMs based on input and output format. Entries in the rectangles are examples of models for each category.

well-known and widely used LLMs is ChatGPT. ChatGPT is designed to converse with humans in natural language and can generate responses to various prompts, including code and math operations [39].

To achieve its high performance, ChatGPT uses reinforcement learning with human feedback and fine-tuning. In the training process, the model is fed with a massive amount of data from various sources to learn different styles and concepts of natural language. During the fine-tuning process, the model is trained on specific tasks, such as question answering or dialogue generation, to further enhance its performance. Additionally, ChatGPT utilizes LM approach that employs a transformer-based neural network architecture to generate responses. This architecture has proven to be highly effective in learning and generating coherent, fluent, and contextually relevant responses. ChatGPT's effectiveness and versatility have made it an essential tool for various NLP applications, including chatbots, language translation systems, and virtual assistants.

LaMDA is another example of an LM for dialog application (LaMDA) developed by Google specifically for dialog applications [59]. It differs from most other LMs in that it was trained on dialogue data, with up to 137B parameters, and pre-trained on 1.56T words of public dialog data and web text. Fine-tuning the model can further enhance its safety and factual grounding. Notably, only 0.001% of the training data is used for fine-tuning, a remarkable achievement of the model. LaMDA takes advantage of Transformers' ability to model long-term dependencies in text, making it well-suited for scaling. With a single model, it can perform multiple tasks, generating multiple responses and filtering them for safety, and grounding them on an external knowledge source. Finally, the model is re-ranked to identify the highest-quality response. These features make LaMDA an ideal candidate for conversational agents that require both a high level of conversational fluency and fact-checking ability.

PEER is a collaborative LM developed by Meta AI research [60]. It is trained on edit histories to cover the complete writing process. The model is based on an iterative process of planning, editing, explaining, and repeating until the text is in a satisfactory state. The model allows users to decompose the task of writing a paper into multiple subtasks, and intervene at any time to steer the model in any direction. It is primarily trained on Wikipedia edit histories, using a self-training approach, where

models are used to infill missing data and then train other models on the synthetic data. However, the model faces a downside due to the noisy comments and lack of citations, which is compensated for by a retrieval system that does not always work. The framework is based on an iterative process, where the entire process of formulating a plan, collecting documents, performing an edit, and explaining it can be repeated multiple times until arriving at a sequence of texts. For the training, the model uses the DeepSpeed transformer [61].
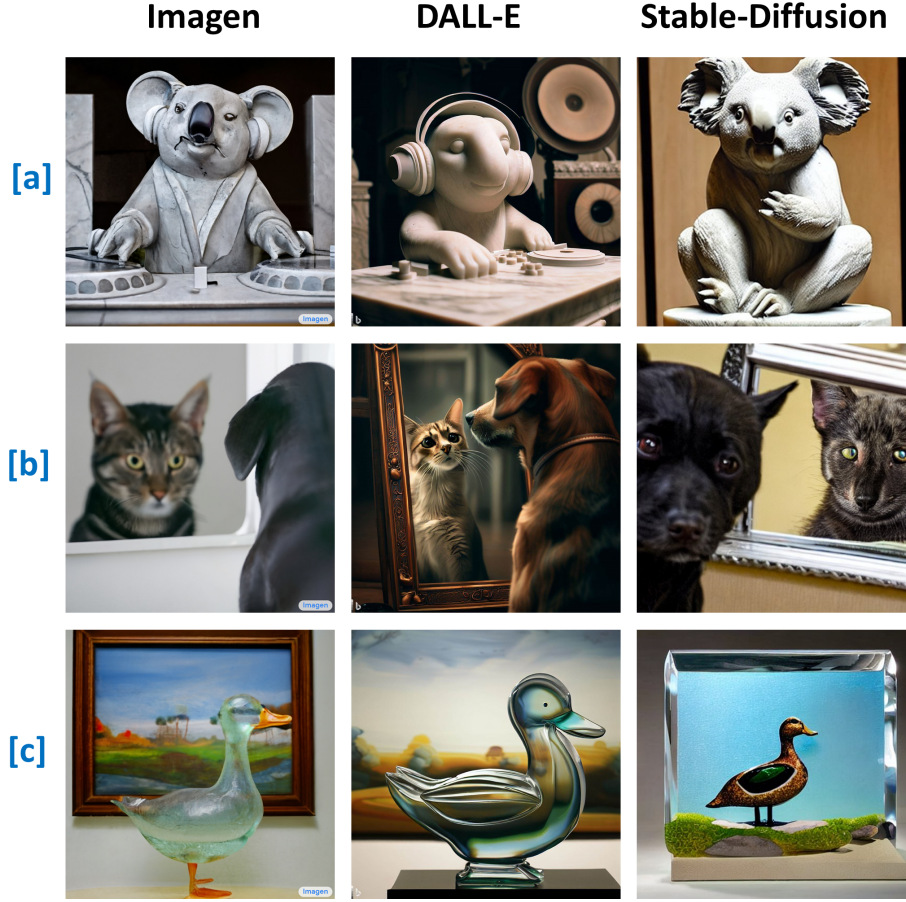
### 2.2.2 Text-to-Image Models

OpenAI has developed DALL-E2 which is a powerful GAI model that can produce highly realistic images based on text prompts [62]. The model leverages CLIP (Contrastive Language-Image Pre-Training), a neural network trained on diverse (image, text) pairs, to combine concepts, attributes, and styles [63]. CLIP enables the model to predict the most relevant text description from an image. The CLIP Image Embeddings Decoder module is integrated with the previous model to generate possible CLIP image embeddings from provided text captions [63]. The resulting DALL-E2 model boasts numerous desirable properties, such as robustness to changes in image distribution, exceptional zero-shot ability, and state-of-the-art performance.

The IMAGEN model is a text-to-image diffusion model comprising a large transformation LM, as described in [64]. According to [65], this model demonstrates that an LLM pre-trained on a text-only corpus can effectively encode text for image synthesis. IMAGEN employs a pre-trained text encoder to map text to a set of word embeddings, and a series of conditional diffusion models to generate high-resolution images from these embeddings. Interestingly, it was observed that the size of the LM is more important for performance than the size of the image diffusion model. To evaluate and compare text-to-image models, Google created Drawbench, a benchmark of 200 prompts [65].

Stable Diffusion is a novel text-to-image LLM developed by the Ludwig Maximilian University of Munich [66]. This model stands out from other existing text-to-image LLMs because of its use of a latent diffusion model, which allows for operations in the latent space for image modification [64]. This approach is much faster and more efficient than previous diffusion models that operate in pixel space. Stable Diffusion consists of three main components: a text encoder, an image generator, and an image information creator. The text encoder takes a textual input and encodes it into a continuous vector representation, which is then used to condition the image generator. The image generator takes the text vector and generates an image that is conditioned on the text input. Finally, the image information creator operates in the latent space and modifies the latent code to perform various image manipulations, such as changing the pose or viewpoint of the generated image. By operating entirely in latent space, it is able to generate high-quality images more efficiently, and it also provides more flexibility for image manipulation.

Muse is a highly efficient text-to-image conversion model that achieves state-of-the-art image generation performance, surpassing diffusion and autoregressive models, as reported in a recent study [67]. This model is trained on a masked modeling task in a distinct token space, which enhances its efficiency as it uses separate tokens

| Imagen | DALL-E | Stable-Diffusion |
|--------|--------|------------------|



**Fig. 7**  Example of text-to-image generation models.The figures are generated based on the text [a] "A marble statue of a Koala DJ in front of a marble statue of a turntable. The Koala has wearing large marble headphones." [b] "A dog looking curiously in the mirror, seeing a cat." [c] "A transparent sculpture of a duck made out of glass. The sculpture is in front of a painting of a landscape."

and requires fewer sampling iterations. Compared to autoregressive models like Parti, Muse is faster due to parallel decoding. During inference, Muse is 10 times faster than Imagen-3B, and three times faster than Stable Diffusion v 1.4, despite the fact that both of these models operate in the latent space of VQGAN [68]. Figure 7 shows the comparison of text-to-image generation models obtained and generated using Imagen[9], DALL-E[10], and stable-diffusion[11] of the given text.

### 2.2.3 Text-to-Audio Models

Google developed AudioLM which takes text as input and generates audio as output generating high-quality audio with long-term consistency [69]. This model uses discrete

---

[9]https://imagen.research.google/
[10]https://www.bing.com/create
[11]https://huggingface.co/spaces/stabilityai/stable-diffusion

tokens to represent input audio and treats audio generation as an LM task in this representation space. By training on vast amounts of raw audio waveforms, AudioLM can generate natural and coherent continuations based on short prompts. Remarkably, the model can even produce coherent piano music continuations without any symbolic representation of music. Generating high-quality audio while maintaining consistency is a challenging task due to the multiple scales of abstraction involved in audio signals [70]. However, this has been achieved by combining recent advances in neural audio compression, self-supervised representation learning, and LM. To determine if the generated audio was of high quality, a sample of 10 seconds was presented to evaluators who were asked to identify whether it was human speech or a synthetic continuation. After gathering 1000 ratings, the success rate was found to be 51.2%, which was not significantly different from random labeling. This suggests that the evaluators were unable to differentiate between synthetic and authentic audio samples, indicating that the generated audio is of high quality.

OpenAI has developed Jukebox, a music generation model that can produce music with singing in the raw audio domain [71]. It is a non-symbolic approach and enables the creation of music that sounds more natural and authentic as it is generated directly as audio. To achieve this, Jukebox uses a hierarchical VQ-VAE (Vector Quantized Variational Autoencoder) [72] architecture that allows audio to be compressed into a discrete space while retaining as much information as possible. This process allows the jukebox to produce long-form music with vocals up to four minutes long while maintaining a high level of consistency and quality. Jukeboxes are not limited to producing a single genre of music. It can create songs in different styles such as pop, rock, country, and classical. In addition, Jukebox can also generate never-before-heard music by combining elements from different genres. This means Jukebox has the potential to be used for music production, sound design, and even music production for film and television. Jukebox's training process is massive, using its massive dataset of 1.2 million songs from LyricWiki. This model has 5 billion parameters and a 9-second audio clip. However, this extensive training process is necessary for Jukebox to be able to produce natural, consistent, high-quality music. The end result is a model that can generate music that is almost indistinguishable from human-generated music.

Whisper is another versatile text-to-audio LLM model that can perform several tasks in the field, including multilingual speech recognition, language identification, and translation [73]. The primary goal of any speech recognition system should be to function accurately in diverse environments without the need for constantly supervised fine-tuning. However, this has been challenging due to the lack of a high-quality pre-trained decoder. To address this, Whisper has been trained on a diverse dataset of 680,000 hours of labeled audio data collected from various sources, which covers a broad distribution of audio from different environments, recording setups, speakers, and languages. The model has been designed for ML generation from the dataset, ensuring that it is only from the human voice. The files are segmented into 30-second intervals paired with the corresponding transcript subset. Whisper employs an encoder-decoder transformer architecture, which has been validated to scale reliably. The model can recognize and transcribe audio in multiple languages, making it highly versatile. Whisper's robustness and accuracy in recognizing speech in different
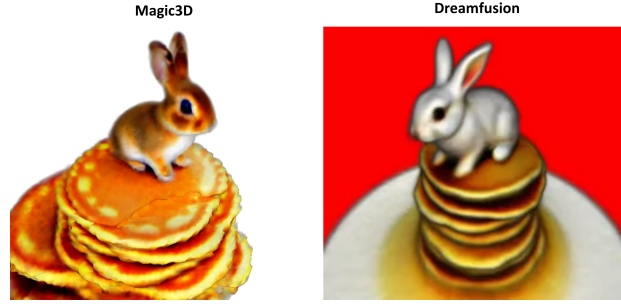
environments and languages make it a valuable tool for various applications, such as automatic transcription, audio indexing, and language learning.

### 2.2.4 Text-to-Video Models

In recent years, researchers have been making progress in generating visual content from textual inputs, specifically images, and videos. While some models can generate still images from text, creating videos from textual input is a more significant challenge. It requires maintaining temporal coherence across frames while ensuring the generated content aligns with the input text. Despite these challenges, research has shown promising results in generating short videos from text descriptions, such as generating short clips of animated characters performing actions described in text prompts. Researchers combine techniques from computer vision and NLP to create coherent sequences of frames that align with the textual input. These advancements have implications for various applications, such as movie and video game production, where generating visuals from text descriptions can streamline the pre-production process.

Phenaki is an advanced video generation model developed by Google Research [74] that uses a C-ViViT encoder, a training transformer, and a video generator to generate coherent and diverse videos from textual prompts. Unlike models that generate still images, Phenaki generates videos from textual input, making it a significant breakthrough in video synthesis. The model is trained on a large dataset of image-text pairs and a smaller dataset of video-text examples, enabling it to produce high-quality videos with a wide range of scenes and styles. Phenaki is capable of generating videos that are several minutes long, even though it is trained on 1.4-second videos. To generate videos, the model compresses the video's representation with the C-ViViT encoder, processes the input text with a temporal transformer and a spatial transformer, and maps the output of the spatial transformer back to pixel space using a linear projection without activation, producing temporally coherent and diverse videos. Phenaki's unique ability to generate videos from open-domain textual input holds great promise for applications such as virtual reality, gaming, and video production. It represents an exciting innovation in video synthesis, with the potential to revolutionize how we create and consume video content.

A new framework called VideoLLM is proposed that uses pre-trained LLMs from NLP as input text for video sequence understanding tasks [75]. VideoLLM includes a Modality Encoder and Semantic Translator to convert inputs from various modalities into a unified token sequence. This token sequence is then processed by a decoder-only LLM, providing a unified framework for different video understanding tasks. Multiple LLMs and fine-tuning methods were used in extensive experiments on eight tasks sourced from four different datasets. The results show that LLMs' comprehension and reasoning abilities can be effectively transferred to video understanding tasks, achieving state-of-the-art or comparable performance with fewer trainable parameters compared to task-specific models. The proposed VideoLLM framework can be used for multiple video sequence understanding tasks, establishing LLM as an effective video reasoner.

**Fig. 8** Example of text-to-3D generation models using "A baby bunny sitting on top of a stack of pancakes" text.

### 2.2.5 Text-to-3D Models

Dreamfusion [76] and Magic3D [77] are two popular text-to-3D models used in the gaming industry to generate 3D images. Dreamfusion is a text-to-3D model developed by Google Research that uses a combination of an image-based approach and a language-based approach. It employs a pretrained text-to-image diffusion model to generate 2D images from textual prompts. The 2D diffusion model is then used as a loss function for a continuous optimization problem that generates 3D models. Dreamfusion addresses the challenge of generating 3D models that look good from any angle when rendered by using a differentiable generator that can create 3D models that look like good images from random angles. Despite achieving remarkable results, Dreamfusion has two problems: long processing time and low quality of the generated images.

Magic3D, on the other hand, is a text-to-3D model developed by NVIDIA Corporation that uses a two-stage optimization framework to generate high-quality 3D models from the text [77]. In the first stage, it builds a low-resolution diffusion prior and accelerates it with a sparse 3D hash grid structure. In the second stage, a textured 3D mesh model is further optimized with an efficient differentiable render. Magic3D uses a generative adversarial network (GAN) to generate 3D models from textual descriptions. It consists of two networks: a generator and a discriminator. The generator generates 3D models from textual descriptions, while the discriminator tries to distinguish between real and generated 3D models. The two networks are trained together in an adversarial manner, where the generator tries to generate 3D models that fool the discriminator, and the discriminator tries to distinguish between real and generated 3D models. Magic3D has been shown to be effective in generating realistic 3D models from textual descriptions, such as cars, airplanes, and chairs. According to human evaluation, Magic3D generates better 3D models than Dreamfusion, with 61.7% of participants preferring Magic3D [77]. Figure 8 shows an example of text-to-3D generation models using Dreamfusion and Magic3D with the text query "a baby bunny sitting on top of a stack of pancakes".

### 2.2.6  Text-to-Code Models

While text-to-text models have been discussed, it's worth noting that not all text follows the same syntax. One such example is code, where it's crucial to translate an idea into a relevant programming language. Codex [78] and Alphacode [79] models are particularly helpful for this purpose. These are innovative AI systems that utilize natural language descriptions to generate functional code.

Codex is developed by OpenAI which is a versatile programming model that can be applied to various programming tasks [78]. The model uses the technique of *program synthesis*, which involves breaking down complex problems into simpler subproblems and mapping those subproblems to existing code libraries, APIs, or functions. Codex is trained on a massive dataset of 179GB of unique Python files under 1 MB, which were collected from public software repositories hosted on GitHub. This extensive dataset allows Codex to learn from a diverse range of programming practices and programming styles and helps it to generate accurate and efficient code. The model's ability to use natural language descriptions to generate code makes it more accessible to non-programmers and can help reduce the time and effort required to write code for complex programming tasks. Additionally, Codex's ability to leverage existing code libraries, APIs, and functions can help reduce errors and improve code quality by ensuring that best coding practices are followed. These innovative features of Codex hold great promise for various applications, including web development, automation, and AI.

On the other hand, Alphacode is an advanced AI system specifically designed to generate functional code for complex and unseen problems that require deeper reasoning [79]. It utilizes transformer-based architectures, including a shallow encoder and a deep encoder, to optimize its efficiency. Alphacode leverages an encoder-decoder transformer architecture, which allows for bidirectional description and provides greater flexibility compared to decoder-only architectures typically used in other code generation models. Multi-query attention is another unique feature of Alphacode, which helps to minimize the cost of sampling. The dataset used for training and evaluation in Alphacode is significantly larger than Codex's dataset. It includes 715.1 GB of code from GitHub repositories and a fine-tuning dataset sourced from the Codeforces platform. The vastness of the dataset ensures that Alphacode is trained on diverse coding styles and scenarios, making it better equipped to handle various coding tasks. The shallow encoder of Alphacode functions by extracting high-level information from the input, while the deep encoder generates detailed information required for the decoding process. This process enables Alphacode to break down complex problems into simpler problems and map them to existing code libraries, APIs, or functions, just like Codex. However, Alphacode's approach is more comprehensive and can handle more complex and intricate programming tasks.

Both models are capable of generating functional code efficiently and accurately, but Alphacode is particularly helpful for generating code for problems that require deeper reasoning, such as those encountered in research or data analysis. Alphacode has been shown to outperform other LMs in generating code for complex, unseen problems. While Codex has a demo and API available for general use, Alphacode is still in the research stage and is not yet available for widespread use.

17

### 2.2.7 Image-to-Text Models

An image-to-text model is a type of computer vision model that aims to generate a natural language description of an image. These models are important for a variety of applications such as image search, image captioning, and assisting the visually impaired. Flamingo [80] and VisualGPT [81] are two LLMs that have received a lot of attention in recent years.

DeepMind developed a visual LM called Flamingo that utilizes few-shot learning techniques on various open-ended tasks involving vision and language [25, 80]. Its unique feature is its visually conditioned autoregressive text generation models, which take in a sequence of text tokens along with images and/or videos as inputs and produce text as output. Flamingo's vision model analyzes visual scenes, while its LM is capable of performing basic forms of reasoning. The LM is trained on large amounts of text data, making Flamingo a powerful tool in situations where only a few labeled examples are available. Flamingo has achieved state-of-the-art performance on various benchmark datasets, including image and video captioning, image question answering, and visual dialog tasks [80].

VisualGPT is an image captioning model developed by OpenAI [81] that builds on the PLM GPT-2 [23] to generate text descriptions for images. It features an innovative encoder-decoder attention mechanism [18] with an unsaturated rectified gating function, which helps to bridge the semantic gap between different modalities. VisualGPT's attention mechanism focuses on the most relevant visual features while generating the text description, resulting in more accurate and relevant captions. Despite being trained on a relatively small amount of labeled data compared to other image captioning models, VisualGPT has shown impressive performance on several benchmark datasets. Its API is available on GitHub, making it accessible for researchers and developers to use in their own projects.

### 2.2.8 Others Models

This section discusses some models that do not easily fit into the above-mentioned categories. One of these models is Alphatensor [82] which is a groundbreaking model created by DeepMind that uses deep reinforcement learning to discover more efficient algorithms for computations such as matrix multiplication. This model is based on a game called TensorGame, where the agent (Alphatensor) is trained to find tensor decompositions within a finite factor space. Alphatensor leverages a specialized neural network architecture to exploit symmetries and improve algorithm efficiency. This model's achievement is significant because it has a widespread impact on computations, from scientific computing routines to neural networks.

Another noteworthy model created by DeepMind is GATO [83], designed as a multi-modal, multi-task, multi-embodiment generalist policy. GATO is a single generalist agent that performs various tasks using the same network and weights. With approximately 1.2 billion parameters, the model is trained at the operating point of the model scale that enables real-time control of real-world robots. GATO's adaptability to different tasks with little extra data is due to the use of a single neural sequence
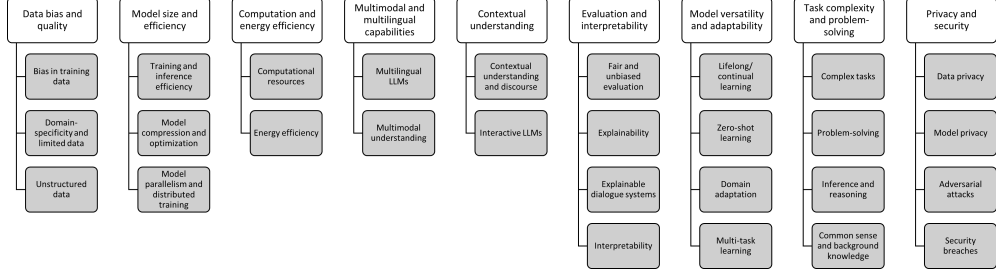
model across all tasks, reducing the need for hand-crafted policy models with their own inductive biases and increasing the amount and diversity of training data [84].

Meta AI Speech Brain is a new model that can directly decode language from noninvasive brain recordings [85]. This is a safer and more scalable alternative to traditional techniques that rely on invasive brain-recording techniques. The model uses a combination of electroencephalography and magnetoencephalography to measure neuronal activity and a DL model with contrastive learning to align brain recordings and speech sounds. The model was trained on 150 hours of recordings from 169 volunteers listening to audiobooks. Results show that the algorithm improves as the number of EEG and MEG recordings increases, indicating that self-supervised trained AI can decode perceived speech regardless of noise and variability in the data. However, the model is limited to speech perception, and future work is needed to expand this to speech production.

Soundify is a video editing system developed by Runway that aims to simplify the process of finding the right sound and matching it to your video [86]. This is achieved by using a high-quality sound effects library and CLIP, a neural network with zero-shot image classification capabilities. Soundify is divided into three main components: categorize, sync, and mix. In the classification component, Soundify matches sound effects to the video by identifying sound emitters in the video. The video is split into segments based on absolute color histogram distance to reduce the number of individual sound emitters. The system then matches sound effects to these segments based on the type of sound the user wants. In the Sync component, Soundify identifies intervals in the video that require sound effects by comparing effect labels to each frame and identifying consecutive matches that exceed a certain threshold. This ensures that the sound effects are in sync with the video visuals. In the mix component, Soundify splits the sound effects into 1-second chunks and seamlessly stitches them together via crossfades. This results in a mix that is both consistent and diverse. Soundify's innovative approach to video editing streamlines the process of adding the right sound to the given videos while saving time and effort. Using a high-quality sound effects library and advanced neural network technology, sounds are well-matched and precisely synchronized with the video for a better overall viewing experience.

In recent years, there has been an influx of LLMs that have been developed and published. These models have a diverse range of capabilities, including the ability to generate human motion [87], perform language translation, and even create presentations. One notable example is ChatBCG, which utilizes ChatGPT as a surrogate model to generate slides based on the input text. The model is trained on a large dataset of slides and their corresponding textual descriptions. ChatBCG employs a sequence-to-sequence architecture that consists of an encoder and a decoder, similar to ChatGPT. The encoder processes the input text, while the decoder generates the slide content based on the encoded text. These models demonstrate the versatility and potential of LLMs and GAI in various applications and industries.

The following section describes the major challenges in LLMs that are still the focus of active research.

19

**Fig. 9** Key issues hierarchy in current LLMs.

## 2.3 Key Challenges in LLMs

Several LLMs based on their inputs and outputs are described in the aforementioned section. However, LLMs also pose several important challenges that are currently the focus of active research [57]. Figure 9 shows the key issues hierarchy for the current LLMs. These issues can impact the usability and applicability of LLMs across various domains. Some of the major issues are described briefly below.
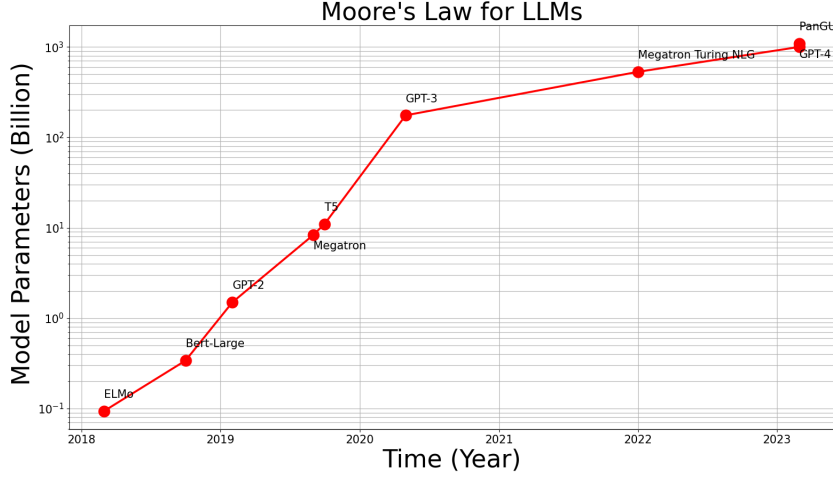
### 2.3.1 Data bias and quality

LLMs require extensive amounts of training data to learn statistical relationships between words and phrases. The effectiveness of LLMs is significantly dependent on the quality of the training data. Acquiring high-quality data can be a costly and time-consuming process. However, training data can be biased, leading to unfair outcomes in predictions [84, 88]. Moreover, since LLMs are trained on data generated by humans, they may reflect the biases present in the real world [84, 89], producing text that is discriminatory, harmful, or offensive [55, 56]. Another challenge facing LLMs is the limited availability of data in specific domains. This can result in poor predictions and hinder the model's performance in these domains. LLMs are typically trained on structured data such as text or speech. It may be challenging to perform well on unstructured data such as images or video. Additionally, LLMs can be used to generate fake news and other forms of misinformation, which can have a negative impact on society. They can also be used to generate creative content, such as articles, code, and scripts. This raises the question of who owns the intellectual property rights to this content, potentially creating new ethical and legal challenges.

### 2.3.2 Model size and efficiency

LLMs are computationally intensive, require a lot of memory, and can be expensive and slow to deploy due to their significant memory requirements [57]. Figure 10 shows the trend of the number of parameters LLMs in comparison to Moore's law [90]. Therefore, current research efforts are focused on developing more efficient training and inference techniques to make LLMs more scalable and accessible [91]. This involves exploring methods to reduce the memory and computational requirements of LLMs while maintaining their performance. Additionally, model parallelism and distributed

**Fig. 10** LLMs: A New Moore's Law, as presented, the size of models grows exponentially with time, even faster than the rate of the original Moore's law regarding the semiconductor industry.

training techniques, such as data parallelism, are being investigated to enable the training of larger LLMs [91]. This approach involves splitting the LLM into multiple smaller models that can be trained simultaneously, reducing the overall training time.

### 2.3.3 Computation and energy efficiency

LLMs are known for their high computational demands both during training and during deployment [92], which can be a challenge for organizations with limited resources [93]. To address this, model optimization techniques are under investigation that can reduce the computational requirements of LLM without sacrificing performance [54]. LLM's large memory and processing power requirements remain a potential barrier for many organizations.

### 2.3.4 Multimodal and multilingual capabilities

Developing LLMs that perform well across multiple languages and modalities is an active research area. Multilingual LLMs are being developed to enable effective linguistic processing in a variety of languages [94]. On the one hand, techniques are being developed that can address the limitations of LLM's single-modality processing, such as text and speech, by allowing information from other modalities to be effectively integrated [95]. These advancements are intended to enhance LLMs' multimodal and multilingual capabilities and facilitate their application in various domains.

### 2.3.5 Contextual understanding

Currently, LLMs evaluations are limited to individual sentences or documents without considering the broader context or discourse. To address this, researchers are working on methods to enable LLMs to understand and incorporate the broader context and

discourse [96]. Another research direction is to create interactive LLMs that can engage in multi-turn dialogue or respond to user feedback in a contextually appropriate manner [94]. Interactive LLMs can provide more personalized and natural responses by understanding and responding to the intent of the user.

### 2.3.6 Evaluation and interpretability

Several directions are being explored by researchers to improve the performance of LLMs. Firstly, developing techniques to evaluate LLM performance fairly and without bias is critical since factors such as the quality and representativeness of training and evaluation data can affect LLM performance [84]. Secondly, it is challenging to understand how LLMs make decisions, as they are trained on vast amounts of data and complex algorithms. This lack of transparency makes it hard to trust LLMs and use them in safety-critical applications [95]. Thirdly, LLM-based dialogue systems need to provide transparent and interpretable responses. Researchers are working on techniques to enable LLM-based dialogue systems to explain their predictions [97]. Finally, developing techniques to enable LLMs to provide interpretable explanations for their predictions is essential for high-stakes tasks such as medical diagnosis [98] or legal decision-making [99]. These research directions are critical for the development of transparent and trustworthy LLMs, enabling their use in diverse applications.

### 2.3.7 Model versatility and adaptability

Making LLMs more versatile and adaptable for new tasks and data is another important but challenging problem. To address this, several techniques are under development, including lifelong learning, which allows LLMs to continue learning by adapting to new tasks and data [100]. Another technique being explored is zero-shot learning [101], which enables LLMs to generalize tasks and domains not seen during training. Domain adaptation techniques [102] are also being developed to allow LLMs to adapt to new domains. Additionally, researchers are exploring multi-task learning techniques that enable LLMs to learn multiple related tasks simultaneously. These research directions are critical for developing LLMs that can effectively handle a wide range of tasks and adapt to new scenarios, making them more useful in real-world applications.

### 2.3.8 Task complexity and problem-solving

Developing LLMs that can handle complex tasks such as long-term planning and decision-making is challenging. These tasks require LLMs to integrate multiple sources of information, reason about uncertainty and ambiguity, and perform problem-solving [56]. While LLMs are good at recognizing patterns and making predictions based on previous examples, they may struggle with tasks that require more complex problem-solving skills, such as decision-making or planning. Additionally, LLMs have limited ability to reason about complex causal relationships, which can hinder their performance on tasks that require causal inference or counterfactual reasoning [95]. LLMs often lack common sense and background knowledge, making it difficult for them to make accurate predictions in certain contexts.

**Fig. 11** Word cloud of LLM (left) and ChatGPT (right).

### 2.3.9 Privacy and security

Data privacy is a critical concern in the development and use of LLMs. These models require large amounts of data, which often contain sensitive information about individuals [32]. It is essential to ensure the privacy and security of the data to prevent unauthorized access or misuse. Adversarial attacks can be used to manipulate the input data intentionally and cause the model to make incorrect predictions or reveal sensitive information. Additionally, the trained models themselves may contain sensitive information, making them vulnerable to attacks that compromise the privacy and integrity of the data. Security breaches can also occur, resulting in the leakage of sensitive data or the manipulation of the model's predictions. Therefore, it is essential to develop techniques to ensure data privacy and security in the development and deployment of LLMs [32].

## 3 ChatGPT: Foundation, Emergence, and Evolution

This section provides an overview of the foundation, emergence, and evolution of ChatGPT. ChatGPT is a state-of-the-art LM that has revolutionized human-machine interactions and significantly contributed to the progress and increased interest in generative AI. Moreover, its emergence has laid the groundwork for the development of increasingly sophisticated and intelligent AI systems. Figure 11 illustrates a comparison between the word clouds for LLM (left) and ChatGPT (right) respectively, effectively highlighting the distinctions in scope and capabilities between the two.

### 3.1 Foundations

From the foregoing, the foundations and origin of ChatGPT can be traced to the developments of LLMs, and hence to the key concepts that underlie LLMs, such as transformers, attention mechanisms, transfer learning/domain adaptation, pretraining, fine-tuning, and generative models. ChatGPT uses a transformer neural network to process and generate text. It is fine-tuned for the specific task of generating human-like responses to natural language prompts using a combination of supervised and unsupervised learning techniques. ChatGPT's foundation also includes various optimization techniques such as weight initialization and gradient clipping to improve its performance and reduce the risk of overfitting to the training data. It is worth noting

that, besides the significant increase in model capacity and architectural modifications, the carefully engineered training process actively contributes to the remarkable performance of ChatGPT. ChatGPT and other LLMs have excelled in various tasks primarily due to two key factors. Firstly, they leverage the transformer architecture, incorporating self-attention mechanisms by effectively teasing out the relationships between input elements. Secondly, they adopt a two-stage training pipeline that comprises self-supervised pre-training and subsequent fine-tuning using small, annotated datasets. This approach enabled efficient learning from unannotated data and achieved high accuracy on specific tasks [103].

## 3.2 Emergence and Evolution

ChatGPT emerged as a large-scale generative language model. The first version, GPT-1, was released in 2018 with 117 million parameters. It achieved state-of-the-art results on several LM tasks.

GPT-2 was released in 2019 with 1.5 billion parameters and showed impressive performance in tasks such as language translation, question answering, and text completion. However, due to concerns about its potential misuse, OpenAI released only a smaller version of GPT-2 to the public. GPT-3, released in 2020 with up to 175 billion parameters, marked a significant leap in LM performance. It achieved state-of-the-art results in a variety of NLP tasks and demonstrated remarkable capabilities in generating coherent, human-like text. Table 1 shows the proportional weighting of GPT-3's dataset sources.

**Table 1** Proportional Weighting of GPT-3's Dataset Sources

| Source | Proportional Weighting |
|---|---|
| Common Crawl | 60% |
| WebText2 (Reddit posts) | 22% |
| Internet-based book collections | 16% |
| English-language version of Wikipedia | 3% |
| Other sources (non-English) | 7% |
| **Total** | **100%** |

The GPT-3.5 architecture is an advanced iteration of the GPT-3 released in 2020 [104]. The GPT-3.5 architecture has been trained on an extensive corpus of data, surpassing 570GB. This extensive training has endowed it with an exceptional ability to learn intricate language patterns and comprehend the subtleties of human communication. With its wide-ranging language tasks, including translation, text completion, and question answering, the GPT-3.5 architecture is at the forefront of natural language processing advancements. Its potential to transform human-machine interactions is remarkable, paving the way for further breakthroughs in the future as AI continues to progress.

OpenAI released the GPT-4 technology on March 14th, 2023 [30]. It accepts text and image inputs, delivers outputs in both text and images and has a context length of 8,192 tokens. GPT-4 is proficient in 26 languages, with superior performance to GPT-3.5 in 24 of those languages.

**Table 2** Comparison of Exam Percentiles between GPT-4 and GPT-3.5

| Category | Exams | GPT-4 % | GPT-3.5 % |
|---|---|---|---|
| Law | Uniform Bar Exam | 90 | 10 |
| | LSAT | 88 | 40 |
| SAT | Evidence-based Reading & Writing | 93 | 87 |
| | Math | 89 | 70 |
| Graduate Record Examination (GRE) | Quantitative | 80 | 25 |
| | Verbal | 99 | 63 |
| | Writing | 54 | 54 |
| Advanced Placement (AP) | Biology | 85 | 62 |
| | Calculus | 43 | 0 |
| | Chemistry | 71 | 22 |
| | Physics 2 | 66 | 30 |
| | Psychology | 83 | 83 |
| | Statistics | 85 | 40 |
| | English Language | 14 | 14 |
| | English Literature | 8 | 8 |
| Medical | Knowledge Self-assessment | 75 | 53 |
| Competitive Programming | Codeforces Rating | < 5 | < 5 |

GPT-4 is an improved version of GPT-3.5 with a maximum token limit of 32,000 a significant increase from GPT-3.5's 4,000 tokens [30]. This latest iteration offers numerous enhancements, including linguistic finesse, information synthesis, creativity, coherence, complex problem-solving, programming capabilities, image and graphics understanding, and a reduction in inappropriate or biased responses. As a result of these qualities, GPT-4 demonstrates versatility across various types of exams. Table 2 compares exam percentiles between GPT-4 and GPT-3.5 [30]. In the Law category, GPT-4 outperformed GPT-3.5 in the LSAT (88th percentile vs. 40th percentile), while both models performed well on the Uniform Bar Exam. GPT-4 achieved higher percentiles in the SAT exams, excelling in Evidence-based Reading & Writing (93rd percentile) and Math (89th percentile) compared to GPT-3.5. For the Graduate Record Examination (GRE), GPT-4 scored exceptionally well in Verbal (99th percentile), while both models achieved the same percentile in Writing (54th percentile). In the Advanced Placement (AP) exams, GPT-4 consistently achieved higher percentiles across various subjects. Notably, GPT-4 scored well in Biology (85th percentile), Chemistry (71st percentile), and Psychology (83rd percentile). In the medical knowledge self-assessment exam GPT-4 (75 percentile) scored higher than GPT-3.5 (53 percentile). In Competitive Programming, both models achieved a rating of less than 5 in Codeforces. It's important to note that the percentile scores provided are relative to their respective exams, allowing for a comparison between the two models' performance in different test scenarios [30]. GPT-4 is a significant advancement over GPT-3.5, with more parameters compared to GPT-3.5. It features multi-modal capabilities, extended short-term memory of 64,000 words, improved multilingual support for 25 languages, and enhanced steerability for more control over responses. These enhancements result in better contextual understanding, the ability to process both text and image data, accurate responses, and significant versatile performance in different exams.

**Table 3** Timeline and evolution of ChatGPT

| Model | Date | Key Capabilities Introduced | Model Size | Data Sources |
|---|---|---|---|---|
| Transformers [18] | Jun 2017 | Attention | 165M Par | Eng-Ger and Eng-Fre datasets |
| GPT1 [105] | Jun 2018 | Generative pre-training | 117M Par | Common Crawl, BookCorpus |
| GPT2 [21] | Feb 2019 | Human Preference and Feedback | 1.5B Par | Common Crawl, BookCorpus, WebText |
| GPT3 [25] | May 2020 | In-context Learning (ICL) | 175B Par | Common Crawl, BookCorpus, Wikipedia, Books, Articles |
| Codex [78] | Jul 2021 | Fundamental on Code Data | 12B Par | Public software repositories hosted on GitHub |
| InstructGPT [95] | Mar 2022 | RLHF-Safety, Instruction following | 1.3B Par | Labelers, SFT, RM, and PPO datasets |
| GPT3.5 [104] | Mar 2022 | Enhanced language understanding | 175B Par | Common Crawl, BookCorpus, Wikipedia, Books, Articles |
| ChatGPT [58] | Nov 2022 | Optimized for Dialog | $\sim$>1T Par (Estimated) | Unknown |
| GPT4 [30] | Mar 2023 | Multimodal Input, Read Teaming | $\sim$>1T Par (Estimated) | Unknown |

ChatGPT has emerged as one of the most powerful and versatile LMs in existence, revolutionizing human-machine interactions. Its development showcases significant advancements in NLP and GAI. Comparing its adoption rates, Google search trends, and mentions in recent literature with other technological innovations demonstrates ChatGPT's rapid emergence and evolution(refer to Fig. 3, Fig. 4, and Fig. 5).

Table 3 provides a timeline and summary of the evolution of ChatGPT, including key capabilities introduced during the evolution.

# 4 ChatGPT: Spawning New Subfields

ChatGPT's popularity and usage have surged dramatically, giving rise to various emerging trends. Member profiles worldwide are increasingly incorporating terms related to ChatGPT and GAI. This trend showcases the significant impact of ChatGPT and its role in shaping new subfields within the discipline. Prompt engineering is one of the subfields, which involves designing high-quality prompts to guide the model's responses and has emerged as a crucial area of research due to ChatGPT's capabilities. Additionally, ChatGPT's exceptional performance in question answering has paved the way for new research in this area, with the potential to create breakthroughs in fields like education [106–110] and medicine [99, 111–113]. The emergence of these new subfields showcases ChatGPT's remarkable impact on the field of GAI

and its potential to drive progress in various domains. Below we present a number of such subfields.

## 4.1 Prompt Engineering

Prompt engineering is a crucial subfield of NLP that plays a vital role. It focuses on crafting high-quality prompts to guide AI models like ChatGPT in generating responses. Given ChatGPT's remarkable ability to produce human-like outputs, prompt engineering has gained significance in facilitating effective and accurate communication with machines. ChatGPT is actively driving advancements in this field by employing innovative techniques that minimize human intervention and enable more efficient and personalized user interactions. This is evident in the model's remarkable performance in tasks such as language translation, conversational agents, and text completion, which are all areas where prompt engineering is critical. There has been a surge of interest among both job seekers and current employees in acquiring the skill of working with ChatGPT and adding it to their resumes. With the model's widespread adoption and impact on various industries, the ability to work with ChatGPT has become a highly sought-after skill that can give individuals a competitive edge in the job market. Zhong et al. [114] investigate ChatGPT's language understanding abilities as well as advanced prompting. They provide some useful strategies for advanced prompting suggested that result in significant improvements in its understanding abilities.

## 4.2 Question Answering

Conversational AI and Question-Answering systems (QASs) have emerged as promising approaches to NLP in recent years [115]. These systems provide users with a more natural way of interacting with machines and allow for more efficient and effective information extraction. One area in which these systems have been particularly successful is in the realm of knowledge graphs (KGs). KGs are structured databases of knowledge that can be used to provide accurate and relevant answers to user queries [116]. While both conversational AI and QASs can be used for information retrieval in KGs, there are some key differences between the two approaches. Conversational AI systems are designed to mimic human conversation, generating responses in natural language that are more engaging and user-friendly. QASs, on the other hand, focus specifically on retrieving information from KGs and use a structured query language to search for information within the KG [117].

ChatGPT as a popular conversational AI model has gained widespread adoption in recent months. This model is capable of generating human-like responses by analyzing large datasets of text, including online chat logs and other forms of unstructured data. KG chatbots, which use KGs as their primary source of knowledge, can benefit from ChatGPT's ability to generate natural language responses [118]. However, KG chatbots also require a QAS mechanism to update their responses with the latest information from the KG. This mechanism allows the chatbot to identify when new information has been added to the KG and update its responses accordingly. In addition, user feedback is an important component of KG chatbots [115]. Users can provide

feedback to help improve the accuracy and relevance of the chatbot's responses, which can be used to refine the chatbot's algorithms and improve the overall user experience.

One challenge faced by QA systems is identifying unanswerable questions [116]. These are questions for which the KG does not contain enough relevant information or the system is unable to retrieve a relevant answer. KG chatbots need to be able to identify unanswerable questions and provide feedback to the user about where to find the information they are looking for [117]. This is an important consideration for the future development of KG chatbots, as users will expect accurate and reliable responses to their queries [115].

## 4.3 Enhanced Search Engine

New enhanced search engines such as Microsoft's Bing are being developed, powered by LLM models like ChatGPT [119]. The engine can automate the manual process of decomposing complex questions into multiple simple questions, retrieving answers for each question, and merging the questions to provide a comprehensive response. Unlike traditional search engines that display lists of links, the LLM-powered search engine can provide an interactive chat-like experience for complex searches, allowing users to participate in dynamic conversations to get useful information. It can also help to provide answers to complex queries and responses [120].

One of the key features of LLM-powered search engine is the ability to summarize results from multiple documents. When a user asks a question, the search engine retrieves relevant information from various sources and provides a concise summary of the results [119]. This will allow users to quickly get to the point and easily get the information they need without having to read multiple documents. Another notable feature of the LLM-powered search engine is the ability to answer questions about recent events [121]. Retrieving information from the web in real-time can enable search engines to provide up-to-date answers to questions about current events, news, and other time-sensitive topics, making it a popular choice for staying informed about the latest developments. This capability is important in mitigating the challenge of recency for LLMs inherently imposed by their training data.

Overall, LLM-powered search engines can represent a major advancement in search technology that leverages the capabilities to provide users with a more sophisticated and interactive search experience. The search engine can answer complex queries, summarize results, provide interactive chat for complex searches, provide quotes for answers, retrieve information from the web in real-time, and more [120]. This can support various features and provides a powerful and efficient tool for search seekers to find accurate and comprehensive information on a wide range of topics [119].

## 4.4 Detection of AI-Generated Content

LLMs like ChatGPT are capable of generating grammatically flawless and seemingly-human replies to different types of questions on various domains [122]. Educators are integrating such technology into classroom instruction, while GAI tools are increasingly becoming publicly available, leading to greater exposure to synthetic content during web browsing. The academic community is also actively exploring methods to

discern whether a given text was produced by a machine or a human, leveraging the fact that there exist systematic differences between human-generated and machine-generated text. Expectedly some tools have quickly emerged to address this problem. Examples of existing tools for such detection include GPTZero [123], GPTRadar [124], Turnitin [125], and Originally.AI [126], However, the problem is far from being resolved [127, 128].

One objective of research in AI-generated content is to enhance the quality of machine-generated text so that it closely approximates human-produced text. In recent research [129] watermarking techniques are proposed that can be incorporated into LLMs used to generate AI-generated text, although these methods are not infallible. During the generation process, LLMs predict the next likely word in a sentence by comparing various alternatives. By designating certain word patterns as off-limits for the AI text generator, a watermark may enable the detection of text produced by a human when the watermark rules are violated multiple times during text scanning.

In the recent study [130], the classification results clearly show that the disambiguation between ChatGPT-generated and human-generated reviews is more challenging for the ML model when the reviews have been rephrased from existing human-generated texts and are not generated by custom queries. ChatGPT can now generate text that is difficult to distinguish from that written by a human to the point [128]. The study suggested that existing techniques for the detection of AI-generated text are not reliable in practice. This paradigm shift challenges our conventional notions of fluency in the language, prompting the need for novel techniques to distinguish between human and machine-generated content. While new tools may become necessary to discern synthetic media in the future, the advice against writing like a robot remains relevant.

# 5 ChatGPT: A Glassbox View

ChatGPT, an LLM, possesses the ability to execute various NLP tasks, including language translation, text summarization, text completion, question answering, and conversation generation. Research on ChatGPT can be categorized into two groups: internal research conducted by the company and other developers of similar LLMs, (we call this the glassbox view), and external research conducted by numerous researchers who evaluate ChatGPT as a blackbox by investigating its responses to specific queries and analyzing multiple aspects of this technology across various domains. In this section, we consider this galassbox view, and we provide a brief overview of key technical components of ChatGPT. Ongoing research efforts focus on enhancing different aspects of ChatGPT, and several important branches are being explored.

1. **Improving language understanding:** Researchers are working on improving ChatGPT's ability to understand the nuances of human language, including sarcasm, irony, and other forms of figurative language [131].
2. **Developing better conversational models:** ChatGPT is being trained to have more engaging and natural conversations with humans, and researchers are working on developing better dialogue systems to achieve this [132].

3. **Enhancing knowledge representation:** Researchers are exploring ways to improve ChatGPT's ability to represent knowledge and build more accurate and effective models for various applications, such as information retrieval, question-answering, and recommendation systems [115, 118].
4. **Fine-tuning for specific applications:** ChatGPT is being fine-tuned for specific applications, such as customer service [133], mental health support [111], and education [107, 113], to improve its performance in these domains [101, 134].
5. **Ethical considerations:** With the increasing use of AI in language models like ChatGPT in various domains, there is a growing interest in researching ethical considerations around their use, such as bias, fairness, and privacy [32]. Here, the focus is on research on reducing bias and improving algorithmic fairness in ChatGPT and LLMs in general [55, 57, 84].

Overall, the research work on ChatGPT is diverse and multifaceted, spanning a wide range of topics related to NLP, ML, and AI ethics.

# 6 ChatGPT: A Blackbox View

Researchers are actively conducting various forms of external research about Chat-GPT, from its capabilities to applications to its limitations. The subsequent subsections describe some of the research covered in the existing literature, while also addressing some underexplored or overlooked aspects.

## 6.1 Capabilities and Abilities

OpenAI introduced ChatGPT, based on the GPT-3.5 architecture, which offers improved features such as user input retention, follow-up corrections, and the ability to decline inappropriate requests. Progressing from GPT-1 to GPT-3, each version exhibited enhanced language generation and expanded capabilities for various NLP tasks. ChatGPT Plus, powered by GPT-4 technology, excels in accepting text and image inputs, providing outputs in both formats and demonstrating proficiency in multiple languages. Notable experimental works have showcased the impressive capabilities of ChatGPT and its evolution in language understanding and generation (refer Table 2).

Deng and Lin [135] conducted a study on earlier versions of ChatGPT (excluding GPT-4), confirming its contextual understanding and ability to provide coherent answers across various languages and tones. In a study by Zhong et al. [114], the language understanding abilities of ChatGPT were explored. Results revealed strong performance on inference tasks, but challenges were observed in paraphrasing and similarity tasks. Advanced prompting strategies were investigated, leading to notable improvements in understanding abilities. Hassani et al. [136] examined the impact of ChatGPT on technical tasks in data science, proposing automation of data pre-processing, framework training, and inference on unstructured data. They also emphasized fine-tuning for specific language tasks. Additionally, Zhai [137] highlighted how ChatGPT can enhance science learning at various stages, offering significant automation potential and increased productivity for trainers and trainees.

A quick summary of the current capabilities of GPT-3.5 and its successor, GPT-4 on various standard examinations can be gleaned from Table 2. The table displays their performance on high-school subject exams such as AP exams and SAT, graduate school exams like GRE, and professional exams in fields like law and medicine. The OpenAI technical report on GPT-4 [30] provides further examples. Examination of the capabilities of ChatGPT in various other fields is an ongoing endeavor.

## 6.2 Current Limitations and Failures

Since its initial release, ChatGPT has undergone significant improvements leading up to the latest version, GPT-4. However, it still possesses inherent limitations, both fundamental and non-fundamental in nature. OpenAI's recent updates acknowledge these limitations, including occasional generation of incorrect or harmful content and limited knowledge updated only until 2021. Researchers from various communities have also examined and identified the current limitations and failures of ChatGPT. In this subsection, we provide a comprehensive review of this line of research, presenting a summary of their findings.

In a series of comprehensive studies, Zhong et al. [114] quantitatively examined the language understanding abilities of ChatGPT. The findings revealed strong performance in inference problems but difficulties in paraphrasing and similarity problems. The study also explored advanced prompting strategies, leading to significant improvements in understanding capabilities. Borji [36] discussed eleven categories of failures experienced by ChatGPT, addressing reliability, trustworthiness, ethical perspectives, and suggestions for efficient utilization. Failures included common sense reasoning, math, factual errors, and coding, emphasizing the need for improved explainability. Kocon et al. [138] conducted an early large-scale automated study analyzing ChatGPT's responses to 38,000 queries across 25 NLP tasks. They observed performance correlation with state-of-the-art benchmarks, indicating lower performance on challenging tasks and vice versa. The study also highlighted bias in ChatGPT's responses linked to its training procedure, emphasizing the importance of establishing standard evaluation protocols for a better understanding of its usefulness.

In a comprehensive investigation by Hariri [139], various aspects of ChatGPT were explored, including applications, limitations, and ethical concerns. The study offered insights on improving utilization and adapting to the technology more effectively. Haleem et al. [140] conducted an early-stage analysis, highlighting the current and future benefits of ChatGPT across different domains of aspects. The study examined the challenges while emphasizing the system's understanding of nuanced aspects of human language. Azaria [141] analyzed ChatGPT's limitations in computing long mathematical expressions, its bias towards certain digits, and its ability to self-correct. The study also revealed the potential for contradictory answers resulting from minor query changes. Yang et al. [142] quantitatively evaluated ChatGPT's text summarization capability (up to GPT-3), demonstrating performance on par with humans and emphasizing the need for further investigation into aspect or query-based summaries. Ray et al. [57] provided a comprehensive survey covering the background, applications, challenges, and future directions of ChatGPT as an advanced chatbot technology. The paper addressed critical challenges such as ethical considerations, biased responses, and

safety concerns, proposing strategies to mitigate them and highlighting the importance of maintaining a balance between AI-assisted innovation and human expertise.

In a query-answer analysis conducted by Shahriar et al. [143], several limitations of ChatGPT were identified, including issues with factual information, reasoning, logic, mathematics, and certain application-based abilities. Aljanabi [144] explored the current capabilities of ChatGPT, along with future possibilities and limitations, to facilitate better adaptation of the technology. Huang et al. [145] conducted an analytical study focusing on the potential and limitations of ChatGPT in implicit hateful speech detection, revealing current limitations in hate speech recognition. Shen et al. [103] extensively investigated the foundations and limitations of ChatGPT, highlighting concerns such as the "hallucination effect" where incorrect responses seem correct, as well as its tendency to act more as an instructions-follower rather than engaging in genuine conversation. The study also discussed limitations in clinical settings and journalism. ChatGPT is known to make up stuff or stories e.g., citing non-existent scientific publications [146]. Guo et al. [147] analyzed and compared responses generated by ChatGPT and human experts across different domains, identified differences and gaps, and explored future directions for language models. The study also sought to develop effective detection systems for distinguishing between human-generated and ChatGPT-generated content.

## 6.3 Applications

The widespread popularity and rapid adoption of ChatGPT have sparked significant interest in its applications across diverse fields. This section, we highlight some of the key applications, recognizing that the list is not exhaustive and that numerous additional applications are anticipated in the near future.

### 6.3.1 Education

ChatGPT can be helpful as a general educational tool or for personalized education. Some examples of these applications are given below.

- Personalized Learning: ChatGPT can be used as a personalized learning tool that adapts to the student's level of knowledge, pacing the learning as needed to provide an interactive learning experience. It can provide personalized recommendations for learning resources and practice exercises based on the student's needs and performance.
- 24/7 Availability: ChatGPT can be accessed 24/7, making it convenient for students to access learning resources and support whenever they need it.
- Answering Questions: ChatGPT can be used to answer students' questions, clarify concepts, and provide additional information on topics they are studying.
- Language Learning: ChatGPT can be used as a language learning tool, allowing students to practice their speaking and writing skills by conversing with the model.
- Accessibility: ChatGPT can be used to make education more accessible for students with disabilities or those who are unable to attend traditional classroom settings.

Since its emergence in late 2022, ChatGPT has gained significant attention as an educational tool. A notable study focused on its performance on the United States Medical Licensing Examination (USMLE) [111], where the model demonstrated impressive knowledge and achieved near or passing scores without specific training. This suggests potential applications in medical education and decision-making. Early research in economics and finance [134] explored ChatGPT's ability to analyze data, propose scenarios, and present results while highlighting limitations and considerations for its effective use. A case study by Susnjak [148] raised concerns about the integrity of online academic exams and proposed alternative evaluation methods and detection techniques. Furthermore, studies in [108, 109] explored how ChatGPT can enhance teaching and learning, addressing opportunities and challenges such as ethical concerns. Kohne et al. [149] examined the use of ChatGPT in language teaching, emphasizing adaptation strategies for teachers.

In another study [110], researchers examined the ethical and responsible use of ChatGPT in education, proposing strategies to enhance privacy, fairness, and sustainability. Farrokhnia et al. proposed strength, weakness, opportunity, and threat (SWOT) analysis in [150], where strengths were identified in generating coherent answers and offering personalized training, while weaknesses included a lack of human-level understanding and challenges in assessing response quality. Banerjee et al. [106] explored the implications of ChatGPT in computer science and engineering education, suggesting practical ways to improve educational quality. Khan et al. [112] investigated the usage and impact of ChatGPT in medical training and management. Bishop [151] examined the early version of ChatGPT in terms of its potential for education and scientific writing, considering theoretical aspects in computing science and philosophy. Additionally, Atlas [152] provided a comprehensive guide for better utilization of ChatGPT in education and professional settings. Numerous other studies explored various aspects of ChatGPT in education, covering academia adaptation [106, 153, 154], examination abilities [155], future implications [156], open education [157], medical education and examination [158–161], bioinformatics education [99], physics education [162], engineering [163], responsible use [164], opportunities and challenges [165, 166], perspectives and concerns [167], proper assessment [168], and tourism research [169].

### 6.3.2 Research

ChatGPT proved to be valuable for various research activities in academia, encompassing literature review, data analysis, virtual research assistance, collaborative platforms, and personalized learning. Ongoing research has been conducted in this domain, examining the potential of ChatGPT for research purposes. In "ChatGPT is fun, but not an author" [146], an early editorial published in Science, the author expressed concerns about the use of ChatGPT in scientific writing, leading to updated editorial policies to prevent scientific misconduct. King et al. [170] evaluated the hypothesis development and testing capabilities of ChatGPT Plus, powered by the GPT-4 algorithm, showcasing promising results in generating novel hypotheses and performing numerical testing. Quintans et al. [171] investigated the implications of ChatGPT for academia, addressing advantages, concerns related to chatbots as authors, disinformation, and ethical issues, and emphasizing the need for regulations and awareness in the research

33

community. MacDonald et at. [172] conducted a comprehensive study on ChatGPT's ability in data analysis and research paper drafting, developing a framework based on simulated data on vaccine effectiveness and highlighting both the advantages and concerns for academia and the research community. Other notable research explored topics such as plagiarism in academia [107], the impact of ChatGPT's hallucination [173], and ChatGPT as an author [174]. Beyond the difficult issue of ethics and plagiarism and connected with the problem of hallucination, the use of ChatGPT and (GAI tools in general) in research could mean fake and incorrect information. Generated by AI can now be easily introduced as acceptable peer-review publications, which can have dire consequences on the credibility and trust in research enquiry in every field.

### 6.3.3 Healthcare

In this section, we summarize the research literature on ChatGPT for healthcare, focusing primarily on its application in healthcare education and exploring the opportunities and challenges associated with this technology. The recent advancements in AI have paved the way for the utilization of conversational AI platforms like ChatGPT, presenting the potential to revolutionize healthcare delivery. By leveraging its natural language understanding and intelligent response capabilities, ChatGPT can assist healthcare professionals in various tasks, including diagnosis, treatment planning, and patient communication. Furthermore, ChatGPT can serve as a valuable resource for patients, providing them with on-demand healthcare information and support, ultimately enhancing health literacy and improving patient outcomes. However, it is crucial to prioritize the protection of patient data and adhere to regulatory guidelines, such as Health Insurance Portability and Accountability Act (HIPAA), to ensure the ethical and secure use of this technology.

In early research, Aydin et al. [175] explored ChatGPT's application in preparing literature review articles, demonstrating its potential to expedite the compilation and presentation of research literature in the healthcare field. Hosseini et al. [113] conducted a comprehensive study analyzing the attitudes of 844 participants towards ChatGPT in education, healthcare, and research. Their findings highlighted the potential advantages and disadvantages in each sector, emphasizing the importance of a measured approach to adoption. In dental medicine and oral biology research [176, 177], Eggmann et al. [178] studied the implications of LLMs like ChatGPT. Kung et al. [111] suggested that ChatGPT and similar models could be valuable in medical education and decision-making. Lyu et al. [179] investigated ChatGPT's capability to translate radiology reports into understandable language for patients and healthcare providers, while Lecler et al. [180] focused on unlocking the technology's potential for revolutionizing radiology. Sallam et al. [159, 167] conducted extensive investigations into how ChatGPT can enhance various aspects of healthcare education and services.

Computational biology and specifically genetic [181–183] is also another area that is benefiting from ChatGPT. Biswas [184] discusses ChatGPT's potential uses in public health, highlighting its abilities, advantages, limitations, and concerns. Hisan et al. [160] present a study on the usefulness, advantages, limitations, and disadvantages of ChatGPT in medical training and education, emphasizing the need for careful considerations, see also [99]. Kung et al. [111] and Gilson et al. [161] report interesting

findings on ChatGPT's performance in the United States Medical Licensing Examination (USMLE) and its future implications for medical education. Luo et al. [185] introduce BioGPT, a domain-specific language model pre-trained on biomedical literature, establishing its state-of-the-art performance in biomedical tasks. Li et al. [186] investigate the ethical aspects of language models in medicine and medical research, highlighting concerns such as biases, trust, authorship, accessibility, and privacy. Jin et al. [187] develop GeneGPT, a model pre-trained on genomic data, establishing its state-of-the-art performance in bioinformatics and genomics.

Multiple other research works presented their explorations and findings on ChatGPT's accuracy and usefulness for medical research [188], clinical and medical applications [189, 190], medical journalism [190], ChatGPT's general knowledge about and usefulness for healthcare [191], implications for discharge summaries [192], the evolution of AI in medicine [193], systematic reviews of research on ChatGPT in healthcare [186], medical writing [194, 195], healthcare for mariners [196], GPT-3 as a virtual doctor [197], discussions on the role of AI in translational medicine [198], information access for cancer patients [199], protection of medical information [200], and more.

### 6.3.4 Finance

In the following section, we provide a summary of the research conducted in the field of business and finance, as numerous businesses and organizations embraced this technology and successfully incorporated the associated changes. ChatGPT proved to be advantageous in finance, offering various benefits, including:

- Customer support: ChatGPT can be used as a virtual assistant to provide customer support and answer queries related to banking, investment, and insurance. This can reduce the workload on human agents and improve customer satisfaction.
- Risk management: ChatGPT can analyze large volumes of financial data and identify potential risks and frauds. This can help financial institutions to make better decisions and reduce losses.
- Investment management: ChatGPT can provide personalized investment advice based on individual financial goals and risk preferences. This can help individuals to make informed investment decisions and achieve their financial objectives.
- Trading: ChatGPT can assist traders in making real-time trading decisions by analyzing market data and providing insights on market trends and conditions.
- Financial education: ChatGPT can provide financial education to individuals by answering their queries and providing information on financial planning, budgeting, and investment strategies. This can help individuals to improve their financial literacy and make better financial decisions.

In the business context, George and George [201] highlighted ChatGPT's potential in customer service applications and its versatility across sectors such as education, finance, health, news, and productivity. Additionally, Downling and Lucy [202] emphasized ChatGPT's assistance in finance research and data identification while acknowledging limitations in literature preparation and ethical concerns. Zaremba and Demir [203] emphasized the need for further exploration of ethical issues and

interpretable AI in financial applications. Finally, Wenzlaff et at. [204] evaluated ChatGPT's accuracy in answering finance-related questions, compared it to human scholars' responses, and discussed its implications for institutions, particularly in crowdfunding, alternative finance, and community finance.

In [205] Yue et al. explored the potential of ChatGPT combined with Explainable AI (XAI) to enhance financial literacy and informed investment decisions by explaining complex financial concepts. The study suggested the revolutionary impact of this technology in the finance field. In the domain of accounting, Alshurafat [206] highlighted the benefits and challenges of using ChatGPT, emphasized the need for human expert intervention, and addressed privacy and security concerns. Ali et al. [207] investigated the advantages of ChatGPT in finance and banking systems, particularly as a chatbot for customer service, while other studies focused on financial decision-making [208], the textile industry [209], the future of businesses [210], and the transformation of organizations [211].

### 6.3.5 Others

ChatGPT has been utilized in various applications, including law [212–215], where it has aided legal professionals in different tasks such as research, contract drafting, and automated legal advice. Additionally, some studies reported the direct use of ChatGPT in specific research areas such as global warming [216] and ChatGPT as a translator [217]. Furthermore, ChatGPT has played a crucial role in combating disinformation [218] through the analysis of multiple sources to identify false content, thereby supporting the dissemination of accurate information. In the field of policy and politics [219], ChatGPT has provided insights by analyzing documents, speeches, and public sentiment [220], contributing to policy summaries, impact evaluation, and simulated discussions to enhance evidence-based decision-making and foster public discourse. In a recent study [221], it was revealed ChatGPT's political bias favoring liberal politicians, its ability to create short stories, and its capacity to rank U.S. presidents based on historical perspectives.

## 7 Responsible AI in the Age of ChatGPT

In this section, we first examine the public reaction to this technology, considering the various perspectives and attitudes toward its use. We then delve into the important aspects of regulations, fairness, privacy, and security, addressing the critical concerns surrounding the implementation and use of ChatGPT.

### 7.1 Public Response to ChatGPT

Various news outlets and experts have provided their reflections on ChatGPT and its impact. The New York Times [222] discusses the potential degradation of science and ethics due to the flawed conception of language and knowledge in machine learning. The Wall Street Journal has published an article on the introduction of ChatGPT, its applications for daily productivity at work, and the government's perspective [223]. Bloomberg has covered different aspects of ChatGPT, including its efficient usage and

the challenges faced by regulators in keeping up with its advances [224]. CBS Morning interviewed Geoffrey Hinton, addressing the potential of AI and its impact on humanity [225]. The Guardian has explored various topics related to ChatGPT, including its impact on health, happiness, productivity, and its ability to generate fake news articles, as well as the contrasting effects of human touch and AI newsreaders [226]. The Economist has raised concerns about the potential automation of tasks traditionally performed by humans, such as telemarketing, teaching, and trading, due to the advancement of ChatGPT [227]. More recently, the ChatGPT delivered a sermon in Germany's church with a large crowd [228]. Various other news organizations have expressed their perspectives on the rapid advancement of ChatGPT, offering a mix of optimism and caution (see for more examples [229–231]).

## 7.2 Regulations

The use of LLMs for commercial purposes has raised concerns regarding ethics, privacy, copyright infringement, and regulatory compliance. ChatGPT as a widely used LLM with over 100 million active users to date, has been banned in several countries due to privacy breaches. For example, ChatGPT is banned in Italy where a data breach involving payment information and user conversations led to its prohibition [232]. Regulatory compliance with the General Data Protection Regulation (GDPR) has been a major issue for OpenAI, as compliance in some countries remains uncertain.

Legislation for AI-based systems is underway in Europe. The lawmakers in the USA want to create an AI regulator body to protect people and control AI [233]. However, it may take years to have a significant impact. Compliance with regulatory bodies is crucial to safeguard individuals' privacy, and the recent ban of ChatGPT in Italy highlights the importance of such compliance. OpenAI's approach of using scraped data for training ChatGPT without compensating individuals has also sparked debate regarding copyright infringement. A recent court case by a writers' organization has accused OpenAI of using copyrighted materials in books and novels without appropriate permission [234]. Taking all the above mentioned together, it is necessary that governments or other entities in charge, develop and apply appropriate regulations that prevent undesired outcomes of such technologies and smooth out their public adoption while being careful to avoid halting their advancement. It is important to note that such regulations need long-term ongoing research as they need to evolve at the same pace that technology evolves and could provide appropriate guidelines both for the developers and users of such technologies.

The European Parliament is currently negotiating new regulations for AI [235]. The proposed regulations would ban some AI applications, such as predictive policing, and require increased transparency for high-risk AI systems used in border control that are used to make decisions about people's lives [236]. The negotiations are also considering the role of technology companies in the regulatory process. The goal of the negotiations is to establish a global standard for AI regulation that balances the interests of different stakeholders.

## 7.3 Fairness

The LLMs are often referred to as "black boxes" as it is hard to understand how they provide predictions or recommendations. The use of LLMs can create unintended biases and unfairness in their outputs [84]. For example, if the data used to train the model is biased or incomplete, the model may perpetuate these biases in its outputs [57, 84]. Additionally, the algorithm used to generate the outputs may not consider all relevant factors, leading to unfair or discriminatory results. We have explored ongoing research on such cases in the section on limitations and failures.

To address these issues, it is important to ensure that LLMs are developed and deployed with fairness in mind. This includes careful consideration of the data used for training, as well as the design of the algorithm itself. It may also involve testing and validation to ensure that the model does not produce biased or discriminatory outputs. It is also imperative to have policies in place to address and remediate unintended biases and injustices that may arise over time.

## 7.4 Privacy and Security

The growing use of LLMs has given rise to concerns about privacy and security violations, as personal data is a critical component in the commercial use of these models. ChatGPT is one of the fastest-growing consumer applications, boasting over 100 million active users to this date, and has been the foundation for numerous startups. Despite its widespread use, however, the issue of privacy violations has received insufficient attention. Articles have highlighted the potential for privacy violations with ChatGPT's training data, which is scraped from various sources, including posts, websites, articles, books, and personal information, often without proper consent [32].

The absence of consent in the data scraping process may raise privacy and security concerns, especially when personal information or identification is involved. Moreover, even the use of publicly available data can breach contextual integrity when used in ways that were not intended. Furthermore, the storage of personal information by OpenAI, the creator of ChatGPT, raises questions about compliance with the GDPR in some countries such as Italy [232]. The watchdog overseeing the ban raised concerns about verifying users' ages and the appropriateness of some responses generated by ChatGPT for underage users. These events highlight the importance of complying with regulatory bodies to safeguard individuals' privacy and security.

# 8  Emergent Required Studies, and Future Directions

In this section, we explore the need for additional lines of research that demand long-term studies. The following will help researchers and relevant entities to undertake comprehensive investigations into generative AI, with a specific focus on ChatGPT, ensuring a more inclusive understanding of its capabilities and implications.

## 8.1 Detector Limitations

One of the limitations is that detectors require access to the inner workings of the LMs, which may not always be feasible or practical. Additionally, detectors are model-specific, meaning that a detector trained to recognize plagiarism with ChatGPT may not be effective in identifying text generated by other LMs, such as LaMDA. Moreover, the accuracy of detectors may be compromised, as these models are often inaccurate and do not provide explanations for their decisions. False positives can occur, making it ethically challenging to solely rely on them for detecting plagiarism. Furthermore, watermarking detectors that rely on specific word lists to identify human-written text can be easily manipulated, as people can post-edit AI-generated text to include "forbidden" words, thus evading detection. These limitations highlight the challenges in using detectors with LMs and the need for further research and development to address their shortcomings.

## 8.2 LLM-based Search Engines

The LLM-powered search engine may have several limitations. Firstly, the model may generate inaccurate or fabricated information with confidence, as it relies on patterns learned from data and may not always produce reliable responses (training data issue, biased data, fine-tuned data) [84]. Additionally, when it comes to dealing with numbers, particularly in the context of financial reports, the LLM model may struggle and may not always provide accurate summaries or interpretations of numerical data. Moreover, running the model with every search query can be computationally expensive, potentially affecting the speed and efficiency of search results during peak usage times. The chat-like interactions of the search engine, while designed to enhance user experience, may also come across as "creepy" or unnatural when the model generates responses that mimic human-like language. Lastly, companies that train the models (e.g. OpenAI, Microsoft) can implement limitations on the number of utterances in conversations to address concerns, which means that the length and depth of conversations may be limited, potentially affecting the comprehensiveness of information obtained through the chat-like interface.

## 8.3 Training In-house LLMs

Data ownership is a crucial factor in decision-making about whether or not to use in-house LLMs. Sharing data with third-party APIs like OpenAI can raise privacy concerns. Although OpenAI's current policy states that they won't use user data to improve their models without explicit opt-in, many people still worry about data leakage. Thus, owning and controlling data is one of the primary reasons why companies opt for in-house LLMs.

Control is another key consideration for companies. Since third-party APIs are out of their control, companies are at the mercy of the API providers. Companies may face issues if API providers decide to stop supporting their country, industry, or company. Such an eventuality could cause operational disruptions, and companies would have to switch to another API, which can be challenging.

Inference cost is a considerable factor in the decision to use in-house LLMs. Although APIs provide quick and easy access to LLMs, the cost of inference may be high. With increasing API usage, inference costs may continue to rise. On the other hand, the cost of training an LLM like GPT-4 is relatively high, but once it is done, the cost of inference is lower.

Performance is a crucial factor in an organization's decision-making to use in-house LLMs. Although open-source models are available, their performance often falls short of GPT-4. Training a model like GPT-4 is an expensive and complex process. Therefore, it can be risky and time-consuming for companies to invest in creating in-house LLMs that may not perform as well as the industry standard.

Safety is another important consideration for organizations. When using in-house LLMs, companies are responsible for any output generated, including toxic content. Third-party API providers like OpenAI, Anthropic, and DeepMind have spent years researching ways to make their LLMs less harmful. Companies would need to invest a significant amount of time and money in research to ensure that their in-house LLMs do not generate any harmful content. Moreover, with the rapid advancements in AI, keeping up with safety protocols can be challenging and expensive.

## 8.4 Impact on Human Languages

ChatGPT has brought about significant changes in the way humans use and interact with language. Its ability to generate human-like responses in multiple languages has made communication more efficient and accessible in the short term. However, its long-term impact on the evolution of language remains uncertain, with potential risks of language simplification and increased language barriers. Moreover, concerns about bias and fairness arising from ChatGPT's reliance on data and algorithms need to be addressed to prevent the perpetuation of existing linguistic biases and inequalities [84].

In the short term, ChatGPT has undoubtedly made communication more efficient and accessible, especially for people who are not proficient in a particular language. By understanding and generating human-like responses, ChatGPT has facilitated cross-cultural communication and increased the connectivity of the world. However, its long-term impact on language evolution is complex and multifaceted. While ChatGPT has the potential to preserve and document endangered languages, its reliance on preprogrammed language patterns and structures could lead to a loss of nuance and complexity in human communication. This may result in a gradual shift in the way humans use language, which could have profound long-term effects on the way we express ourselves.

Another significant issue with ChatGPT is the potential for language barriers to become more pronounced. While ChatGPT can help people communicate across different languages and cultures, it may also create a "standard" language that only some people can understand. This could reinforce linguistic divides and exacerbate existing language barriers. Moreover, the displacement of human translators and interpreters due to ChatGPT's widespread adoption could have negative economic impacts on certain communities.

ChatGPT's reliance on data and algorithms raises concerns about ethics, bias, and fairness [84]. If the data used to train ChatGPT is not diverse or inclusive, it may

perpetuate existing biases and prejudices in language use. This could lead to further marginalization of certain groups and exacerbate social inequalities. Therefore, it is crucial to address these concerns and work towards a more inclusive and equitable use of ChatGPT.

## 8.5 Effect on Society at Large

The widespread use of LLMs will have significant implications for society at large [32]. Firstly, it will lead to the creation of misinformation at an unprecedented scale. With AI-generated content becoming more prevalent [237], it will become increasingly challenging to discern between genuine and fabricated information, leading to confusion and potentially harmful consequences. Secondly, the issue of accountability will arise, as the responsibility for bad AI-generated legal, medical, financial, or other types of advice becomes a major issue. Determining who will be liable for the outcomes of such advice becomes complex, posing legal and ethical challenges. Additionally, people may develop imaginary relationships with AI entities, as depicted in the movie "Her" (2013), blurring the lines between human and machine interactions. Lastly, with increasing automation in content creation there is the risk that society can become trapped in a perpetual state of intellectual stagnation, resulting in limited creative evolution, stifled innovation, and a general decline in societal progress.

## 8.6 Other Areas

ChatGPT, as an LLM, has continually improved and expanded its capabilities. The future of ChatGPT appears promising, with potential avenues for growth and development, including enhancements in NLP, integration with other technologies, personalization, and expansion into domains like healthcare, education, and finance. Researchers from various fields have actively engaged in ongoing research, examining the utility, concerns, and future directions of this technology [147]. Furthermore, in terms of future directions, it is crucial to address concerns and regulations associated with the rapid evolution of ChatGPT. In their study, Guo et al. [147] analyze and compare the responses generated by ChatGPT and human experts across diverse domains, identify disparities and gaps, and propose future directions for LLMs. The study also aims to develop effective detection systems to distinguish between content generated by humans and ChatGPT. The ethical issue of training LLMs with freely available data owned by other entities (e.g. books, novels, papers, reports, etc). This must be seriously considered the resolution of which could have far-reaching impacts on the future of LLMs and generative AI tools such as ChatGPT, and by extension on the evolution of human languages and society at large.

# 9 Conclusion

The emergence of LLMs such as ChatGPT has revolutionized the field of NLP and opened up new avenues for research and development in generative AI. This paper has provided a concise overview of the current lines of research on ChatGPT and its different versions as a black box, as well as established a roadmap for further experimental

research and studies. While significant progress has been made, there are still gaps in our understanding of these models and their impact on various fields. As the field continues to evolve, it is essential to continue exploring the various applications of LLMs and addressing any concerns that may arise. Ultimately, this will lead to the development of more robust and useful tools that can benefit both developers and end-users in many areas, including education, research, healthcare, law, finance, and beyond.

# 10 Acknowledgment

# References

[1] Hauser, M.D., Chomsky, N., Fitch, W.T.: The faculty of language: What is it, who has it, and how did it evolve? Science **298**(5598), 1569–1579 (2002)

[2] Turing, A.M.: Computing machinery and intelligence. In: Parsing the Turing Test: Philosophical and Methodological Issues in the Quest for the Thinking Computer, pp. 23–65 (2009)

[3] Gao, J., Lin, C.-Y.: Introduction to the special issue on statistical language modeling. ACM Transactions on Asian Language Information Processing (TALIP) **3**(2), 87–93 (2004)

[4] Rosenfeld, R.: Two decades of statistical language modeling: Where do we go from here? Proceedings of the IEEE **88**(8), 1270–1278 (2000)

[5] Bahl, L.R., Brown, P.F., Souza, P.V., Mercer, R.L.: A tree-based statistical language model for natural language speech recognition. IEEE Transactions on Acoustics, Speech, and Signal Processing **37**(7), 1001–1008 (1989)

[6] Brants, T., Popat, A.C., Xu, P., Och, F.J., Dean, J.: Large language models in machine translation. arXiv:2304.02210 (2007)

[7] Croft, B., Lafferty, J.: Language Modeling for Information Retrieval vol. 13, (2003)

[8] Zhai, C.: Statistical language models for information retrieval. Synthesis Lectures on Human Language Technologies, 1–141 (2008)

[9] Katz, S.: Estimation of probabilities from sparse data for the language model component of a speech recognizer. IEEE transactions on Acoustics, Speech, and Signal Processing **35**(3), 400–401 (1987)

[10] Bengio, Y., Ducharme, R., Vincent, P.: A neural probabilistic language model. Advances in Neural Information Processing Systems **13** (2000)

[11] Mikolov, T., Karafiát, M., Burget, L., Cernockỳ, J., Khudanpur, S.: Recurrent neural network based language model. In: Interspeech, vol. 2, pp. 1045–1048 (2010)

[12] Collobert, R., Weston, J., Bottou, L., Karlen, M., Kavukcuoglu, K., Kuksa, P.: Natural language processing (almost) from scratch. Journal of Machine Learning Research **12**, 2493–2537 (2011)

[13] Mikolov, T., Sutskever, I., Chen, K., Corrado, G.S., Dean, J.: Distributed representations of words and phrases and their compositionality. Advances in Neural Information Processing Systems **26** (2013)

[14] Mikolov, T., Chen, K., Corrado, G., Dean, J.: Efficient estimation of word representations in vector space. arXiv:1301.3781 (2013)

[15] Sarzynska-Wawer, J., Wawer, A., Pawlak, A., Szymanowska, J., Stefaniak, I., Jarkiewicz, M., Okruszek, L.: Detecting formal thought disorder by deep contextualized word representations. Psychiatry Research **304**, 114135 (2021)

[16] Devlin, J., Chang, M.-W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv:1810.04805 (2018)

[17] Hochreiter, S., Schmidhuber, J.: Long short-term memory. Neural Computation **9**(8), 1735–1780 (1997)

[18] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. Advances in Neural Information Processing Systems **30** (2017)

[19] Zhao, W.X., Zhou, K., Li, J., Tang, T., Wang, X., Hou, Y., Min, Y., Zhang, B., Zhang, J., Dong, Z., et al.: A survey of large language models. arXiv:2303.18223 (2023)

[20] Fedus, W., Zoph, B., Shazeer, N.: Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. The Journal of Machine Learning Research **23**(1), 5232–5270 (2022)

[21] Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I., *et al.*: Language models are unsupervised multitask learners. OpenAI Blog **1**(8), 9 (2019)

[22] Kaplan, J., McCandlish, S., Henighan, T., Brown, T.B., Chess, B., Child, R., Gray, S., Radford, A., Wu, J., Amodei, D.: Scaling laws for neural language models. arXiv:2001.08361 (2020)

[23] Budzianowski, P., Vulić, I.: Hello, it's GPT-2–how can I help you? Towards the use of pretrained language models for task-oriented dialogue systems.

arXiv:1907.05774 (2019)

[24] Wei, J., Tay, Y., Bommasani, R., Raffel, C., Zoph, B., Borgeaud, S., Yogatama, D., Bosma, M., Zhou, D., Metzler, D., et al.: Emergent abilities of large language models. arXiv:2206.07682 (2022)

[25] Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J.D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., *et al.*: Language models are few-shot learners. Advances in Neural Information Processing Systems **33**, 1877–1901 (2020)

[26] Shanahan, M.: Talking about large language models. arXiv:2212.03551 (2022)

[27] Hoffmann, J., Borgeaud, S., Mensch, A., Buchatskaya, E., Cai, T., Rutherford, E., Casas, D.d.L., Hendricks, L.A., Welbl, J., Clark, A., et al.: Training compute-optimal large language models. arXiv:2203.15556 (2022)

[28] Epstein, Z., Levine, S., Rand, D.G., Rahwan, I.: Who gets credit for AI-generated art? Iscience **23**(9) (2020)

[29] Gozalo-Brizuela, R., Garrido-Merchan, E.C.: ChatGPT is not all you need. A state of the art review of large generative AI models. arXiv:2301.04655 (2023)

[30] OpenAI: GPT-4 Technical Report. arXiv:2303.08774v3 (2023)

[31] Altman, S.: Planning for AGI and beyond. OpenAI Blog, February (2023)

[32] Khowaja, S.A., Khuwaja, P., Dev, K.: ChatGPT needs spade (sustainability, privacy, digital divide, and ethics) evaluation: A review. techrxiv.22619932 (2023)

[33] Challenge, B.: Sample-efficient pretraining on a developmentally plausible corpus. Github (2023). https://babylm.github.io/ Accessed July 5, 2023

[34] Huang, S., Dong, L., Wang, W., Hao, Y., Singhal, S., Ma, S., Lv, T., Cui, L., Mohammed, O.K., Liu, Q., et al.: Language is not all you need: Aligning perception with language models. arXiv:2302.14045 (2023)

[35] Driess, D., Xia, F., Sajjadi, M.S., Lynch, C., Chowdhery, A., Ichter, B., Wahid, A., Tompson, J., Vuong, Q., Yu, T., et al.: Palm-e: An embodied multimodal language model. arXiv:2303.03378 (2023)

[36] Borji, A.: A categorical archive of ChatGPT failures. arXiv:2302.03494 (2023)

[37] Garfinkle, A.: ChatGPT on track to surpass 100 million users faster than Tik-Tok or Instagram: UBS. Yahoo Finance (Feb 2023). https://finance.yahoo.com/ Accessed July 5, 2023

44

[38] Huang, J., Chang, K.C.-C.: Towards reasoning in large language models: A survey. arXiv:2212.10403 (2022)

[39] Cao, Y., Li, S., Liu, Y., Yan, Z., Dai, Y., Yu, P.S., Sun, L.: A comprehensive survey of AI-generated content (AIGC): A history of generative AI from GAN to ChatGPT. arXiv:2303.04226 (2023)

[40] Zhou, J., Ke, P., Qiu, X., Huang, M., Zhang, J.: ChatGPT: Potential, prospects, and limitations. Frontiers of Information Technology & Electronic Engineering, 1–6 (2023)

[41] Zhang, C., Zhang, C., Li, C., Qiao, Y., Zheng, S., Dam, S.K., Zhang, M., Kim, J.U., Kim, S.T., Choi, J., et al.: One small step for generative AI, one giant leap for AGI: A complete survey on ChatGPT in AIGC era. arXiv:2304.06488 (2023)

[42] Zhang, C., Zhang, C., Zheng, S., Qiao, Y., Li, C., Zhang, M., Dam, S.K., Thwal, C.M., Tun, Y.L., Huy, L.L., et al.: A complete survey on generative AI (AIGC): Is ChatGPT from GPT-4 to GPT-5 all you need? arXiv:2303.11717 (2023)

[43] Zhou, Z.-H.: Machine Learning, (2021)

[44] LeCun, Y., Bengio, Y., Hinton, G.: Deep learning. Nature **521**(7553), 436–444 (2015)

[45] Chowdhary, K., Chowdhary, K.: Natural language processing. Fundamentals of Artificial Intelligence, 603–649 (2020)

[46] Federico, M., Cettolo, M., Brugnara, F., Antoniol, G.: Language modelling for efficient beam-search. Computer Speech and Language **9**(4), 353–380 (1995)

[47] Meister, C., Cotterell, R.: Language model evaluation beyond perplexity. arXiv:2106.00085 (2021)

[48] Malladi, S., Gao, T., Nichani, E., Damian, A., Lee, J.D., Chen, D., Arora, S.: Fine-tuning language models with just forward passes. arXiv:2305.17333 (2023)

[49] Liu, P., Yuan, W., Fu, J., Jiang, Z., Hayashi, H., Neubig, G.: Pre-train, Prompt, and Predict: A Systematic Survey of Prompting Methods in Natural Language Processing (2021)

[50] Daull, X., Bellot, P., Bruno, E., Martin, V., Murisasco, E.: Complex QA and language models hybrid architectures, survey. arXiv:2302.09051 (2023)

[51] Tene, O., Polonetsky, J.: Big data for all: Privacy and user control in the age of analytics. Nw. J. Tech. & Intell. Prop. **11**, 239 (2012)

[52] Carlini, N., Tramer, F., Wallace, E., Jagielski, M., Herbert-Voss, A., Lee, K., Roberts, A., Brown, T., Song, D., Erlingsson, U., *et al.*: Extracting training data

from large language models. In: 30th USENIX Security Symposium (USENIX Security 21), pp. 2633–2650 (2021)

[53] Chowdhery, A., Narang, S., Devlin, J., Bosma, M., Mishra, G., Roberts, A., Barham, P., Chung, H.W., Sutton, C., Gehrmann, S., et al.: Palm: Scaling language modeling with pathways. arXiv:2204.02311 (2022)

[54] Mehlin, V., Schacht, S., Lanquillon, C.: Towards energy-efficient deep learning: An overview of energy-efficient approaches along the deep learning lifecycle. arXiv:2303.01980 (2023)

[55] McGee, R.W.: Is chat GPT biased against conservatives? an empirical study. An Empirical Study (February 15, 2023) (2023)

[56] Liebrenz, M., Schleifer, R., Buadze, A., Bhugra, D., Smith, A.: Generating scholarly content with ChatGPT: ethical challenges for medical publishing. The Lancet Digital Health **5**(3), 105–106 (2023)

[57] Ray, P.P.: ChatGPT: A comprehensive review on background, applications, key challenges, bias, ethics, limitations and future scope. Internet of Things and Cyber-Physical Systems (2023)

[58] OpenAI: Introducing ChatGPT. OpenAI Blog (2022). https://openai.com/ Accessed July 5, 2023

[59] Thoppilan, R., De Freitas, D., Hall, J., Shazeer, N., Kulshreshtha, A., Cheng, H.-T., Jin, A., Bos, T., Baker, L., Du, Y., et al.: Lamda: Language models for dialog applications. arXiv:2201.08239 (2022)

[60] Schick, T., Dwivedi-Yu, J., Jiang, Z., Petroni, F., Lewis, P., Izacard, G., You, Q., Nalmpantis, C., Grave, E., Riedel, S.: PEER: A collaborative language model. arXiv:2208.11663 (2022)

[61] Aminabadi, R.Y., Rajbhandari, S., Awan, A.A., Li, C., Li, D., Zheng, E., Ruwase, O., Smith, S., Zhang, M., Rasley, J., et al.: Deepspeed-inference: enabling efficient inference of transformer models at unprecedented scale. In: SC22: International Conference for High Performance Computing, Networking, Storage and Analysis, pp. 1–15 (2022). IEEE

[62] Daras, G., Dimakis, A.G.: Discovering the hidden vocabulary of DALLE-2. arXiv:2206.00169 (2022)

[63] Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: International Conference on Machine Learning, pp. 8748–8763 (2021). PMLR

[64] Kingma, D., Salimans, T., Poole, B., Ho, J.: Variational diffusion models. Advances in Neural Information Processing Systems **34**, 21696–21707 (2021)

[65] Saharia, C., Chan, W., Saxena, S., Li, L., Whang, J., Denton, E.L., Ghasemipour, K., Gontijo Lopes, R., Karagol Ayan, B., Salimans, T., *et al.*: Photorealistic text-to-image diffusion models with deep language understanding. Advances in Neural Information Processing Systems **35**, 36479–36494 (2022)

[66] Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 10684–10695 (2022)

[67] Chang, H., Zhang, H., Barber, J., Maschinot, A., Lezama, J., Jiang, L., Yang, M.-H., Murphy, K., Freeman, W.T., Rubinstein, M., et al.: Muse: Text-to-image generation via masked generative transformers. arXiv:2301.00704 (2023)

[68] Yu, J., Li, X., Koh, J.Y., Zhang, H., Pang, R., Qin, J., Ku, A., Xu, Y., Baldridge, J., Wu, Y.: Vector-quantized image modeling with improved VQGAN. arXiv:2110.04627 (2021)

[69] Borsos, Z., Marinier, R., Vincent, D., Kharitonov, E., Pietquin, O., Sharifi, M., Teboul, O., Grangier, D., Tagliasacchi, M., Zeghidour, N.: Audiolm: a language modeling approach to audio generation. arXiv:2209.03143 (2022)

[70] Mujtaba, G., Lee, S., Kim, J., Ryu, E.-S.: Client-driven animated GIF generation framework using an acoustic feature. Multimedia Tools and Applications, 1–18 (2021)

[71] Dhariwal, P., Jun, H., Payne, C., Kim, J.W., Radford, A., Sutskever, I.: Jukebox: A generative model for music. arXiv:2005.00341 (2020)

[72] Ding, S., Gutierrez-Osuna, R.: Group latent embedding for vector quantized variational autoencoder in non-parallel voice conversion. In: INTERSPEECH, pp. 724–728 (2019)

[73] Radford, A., Kim, J.W., Xu, T., Brockman, G., McLeavey, C., Sutskever, I.: Robust speech recognition via large-scale weak supervision. arXiv:2212.04356 (2022)

[74] Villegas, R., Babaeizadeh, M., Kindermans, P.-J., Moraldo, H., Zhang, H., Saffar, M.T., Castro, S., Kunze, J., Erhan, D.: Phenaki: Variable length video generation from open domain textual description. arXiv:2210.02399 (2022)

[75] Chen, G., Zheng, Y.-D., Wang, J., Xu, J., Huang, Y., Pan, J., Wang, Y., Wang, Y., Qiao, Y., Lu, T., et al.: Videollm: Modeling video sequence with large language models. arXiv:2305.13292 (2023)

[76] Poole, B., Jain, A., Barron, J.T., Mildenhall, B.: Dreamfusion: Text-to-3D using 2d diffusion. arXiv:2209.14988 (2022)

[77] Lin, C.-H., Gao, J., Tang, L., Takikawa, T., Zeng, X., Huang, X., Kreis, K., Fidler, S., Liu, M.-Y., Lin, T.-Y.: Magic3D: High-resolution text-to-3D content creation. arXiv:2211.10440 (2022)

[78] Chen, M., Tworek, J., Jun, H., Yuan, Q., Pinto, H.P.d.O., Kaplan, J., Edwards, H., Burda, Y., Joseph, N., Brockman, G., et al.: Evaluating large language models trained on code. arXiv:2107.03374 (2021)

[79] Li, Y., Choi, D., Chung, J., Kushman, N., Schrittwieser, J., Leblond, R., Eccles, T., Keeling, J., Gimeno, F., Dal Lago, A., et al.: Competition-level code generation with Alphacode. Science **378**(6624), 1092–1097 (2022)

[80] Alayrac, J.-B., Donahue, J., Luc, P., Miech, A., Barr, I., Hasson, Y., Lenc, K., Mensch, A., Millican, K., Reynolds, M., et al.: Flamingo: A visual language model for few-shot learning. Advances in Neural Information Processing Systems **35**, 23716–23736 (2022)

[81] Chen, J., Guo, H., Yi, K., Li, B., Elhoseiny, M.: VisualGPT: Data-efficient adaptation of pretrained language models for image captioning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 18030–18040 (2022)

[82] Fawzi, A., Balog, M., Huang, A., Hubert, T., Romera-Paredes, B., Barekatain, M., Novikov, A., R Ruiz, F.J., Schrittwieser, J., Swirszcz, G., et al.: Discovering faster matrix multiplication algorithms with reinforcement learning. Nature **610**(7930), 47–53 (2022)

[83] Reed, S., Zolna, K., Parisotto, E., Colmenarejo, S.G., Novikov, A., Barth-Maron, G., Gimenez, M., Sulsky, Y., Kay, J., Springenberg, J.T., et al.: A generalist agent. arXiv:2205.06175 (2022)

[84] Schramowski, P., Turan, C., Andersen, N., Rothkopf, C.A., Kersting, K.: Large pre-trained language models contain human-like biases of what is right and wrong to do. Nature Machine Intelligence **4**(3), 258–268 (2022)

[85] Défossez, A., Caucheteux, C., Rapin, J., Kabeli, O., King, J.-R.: Decoding speech from non-invasive brain recordings. arXiv:2208.12266 (2022)

[86] Lin, D.C.-E., Germanidis, A., Valenzuela, C., Shi, Y., Martelaro, N.: Soundify: Matching sound effects to video. arXiv:2112.09726 (2021)

[87] Tevet, G., Raab, S., Gordon, B., Shafir, Y., Cohen-Or, D., Bermano, A.H.: Human motion diffusion model. arXiv:2209.14916 (2022)

[88] Obermeyer, Z., Powers, B., Vogeli, C., Mullainathan, S.: Dissecting racial bias in an algorithm used to manage the health of populations. Science **366**(6464), 447–453 (2019)

[89] Caliskan, A., Bryson, J.J., Narayanan, A.: Semantics derived automatically from language corpora contain human-like biases. Science **356**(6334), 183–186 (2017)

[90] Simon, J.: Large Language Models: A New Moore's Law? Hugging Face (2021). https://huggingface.co/ Accessed July 5, 2023

[91] Narayanan, D., Shoeybi, M., Casper, J., LeGresley, P., Patwary, M., Korthikanti, V., Vainbrand, D., Kashinkunti, P., Bernauer, J., Catanzaro, B., *et al.*: Efficient large-scale language model training on GPU clusters using megatron-LM. In: Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis, pp. 1–15 (2021)

[92] Yin, J., Dash, S., Gounley, J., Wang, F., Tourassi, G.: Evaluation of pre-training large language models on leadership-class supercomputers. The Journal of Supercomputing, 1–22 (2023)

[93] Lakim, I., Almazrouei, E., Abualhaol, I., Debbah, M., Launay, J.: A holistic assessment of the carbon footprint of noor, a very large arabic language model. In: Proceedings of BigScience Episode 5–Workshop on Challenges & Perspectives in Creating Large Language Models, pp. 84–94 (2022)

[94] Kasirzadeh, A., Gabriel, I.: In conversation with artificial intelligence: aligning language models with human values. Philosophy & Technology **36**(2), 1–24 (2023)

[95] Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., *et al.*: Training language models to follow instructions with human feedback. Advances in Neural Information Processing Systems **35**, 27730–27744 (2022)

[96] Yuan, A., Coenen, A., Reif, E., Ippolito, D.: Wordcraft: story writing with large language models. In: 27th International Conference on Intelligent User Interfaces, pp. 841–852 (2022)

[97] Wahlster, W.: Understanding computational dialogue understanding. Philosophical Transactions of the Royal Society A **381**(2251), 20220049 (2023)

[98] Nori, H., King, N., McKinney, S.M., Carignan, D., Horvitz, E.: Capabilities of GPT-4 on medical challenge problems. arXiv:2303.13375 (2023)

[99] Shue, E., Liu, L., Li, B., Feng, Z., Li, X., Hu, G.: Empowering beginners in bioinformatics with ChatGPT. bioRxiv, 2023–03 (2023)

[100] Lopez-Lira, A., Tang, Y.: Can ChatGPT forecast stock price movements? return predictability and large language models. arXiv:2304.07619 (2023)

[101] Espejel, J.L., Ettifouri, E.H., Alassan, M.S.Y., Chouham, E.M., Dahhane, W.: GPT-3.5 vs GPT-4: Evaluating ChatGPT's reasoning performance in zero-shot learning. arXiv:2305.12477 (2023)

[102] Motiian, S., Piccirilli, M., Adjeroh, D.A., Doretto, G.: Unified deep supervised domain adaptation and generalization. In: Proceedings of the IEEE International Conference on Computer Vision, pp. 5715–5725 (2017)

[103] Shen, Y., Heacock, L., Elias, J., Hentel, K.D., Reig, B., Shih, G., Moy, L.: ChatGPT and other large language models are double-edged swords. Radiology, 230163 (2023)

[104] OpenAI: New GPT-3 capabilities: Edit & insert. OpenAI Blog (2022). https://openai.com/ Accessed July 5, 2023

[105] Radford, A., Narasimhan, K., Salimans, T., Sutskever, I.: Improving language understanding with unsupervised learning (2018)

[106] Banerjee, P., Srivastava, A., Adjeroh, D., Reddy, Y.R., Karimian, N.: Understanding ChatGPT: Impact analysis and path forward for teaching computer science and engineering. TechRxiv:22639705.v1 (2023)

[107] King, M.R., ChatGPT: A conversation on artificial intelligence, chatbots, and plagiarism in higher education. Cellular and Molecular Bioengineering **16**(1), 1–2 (2023)

[108] Baidoo-Anu, D., Owusu Ansah, L.: Education in the era of generative artificial intelligence (ai): Understanding the potential benefits of ChatGPT in promoting teaching and learning. Available at SSRN 4337484 (2023)

[109] AlAfnan, M.A., Dishari, S., Jovic, M., Lomidze, K.: ChatGPT as an educational tool: Opportunities, challenges, and recommendations for communication, business writing, and composition courses. Journal of Artificial Intelligence and Technology (2023)

[110] Mhlanga, D.: Open AI in education, the responsible and ethical use of ChatGPT towards lifelong learning. Education, the Responsible and Ethical Use of ChatGPT Towards Lifelong Learning (2023)

[111] Kung, T.H., Cheatham, M., Medenilla, A., Sillos, C., De Leon, L., Elepaño, C., Madriaga, M., Aggabao, R., Diaz-Candido, G., Maningo, J., *et al.*: Performance of ChatGPT on usmle: Potential for AI-assisted medical education using large language models. PLoS Digital Health **2**(2), 0000198 (2023)

[112] Khan, R.A., Jawaid, M., Khan, A.R., Sajjad, M.: ChatGPT-reshaping medical education and clinical management. Pakistan Journal of Medical Sciences **39**(2) (2023)

[113] Hosseini, M., Gao, C.A., Liebovitz, D.M., Carvalho, A.M., Ahmad, F.S., Luo, Y., MacDonald, N., Holmes, K.L., Kho, A.: An exploratory survey about using ChatGPT in education, healthcare, and research. medRxiv, 2023–03 (2023)

[114] Zhong, Q., Ding, L., Liu, J., Du, B., Tao, D.: Can ChatGPT understand too? A comparative study on ChatGPT and fine-tuned BERT. arXiv:2302.10198 (2023)

[115] Omar, R., Mangukiya, O., Kalnis, P., Mansour, E.: ChatGPT versus traditional question answering for knowledge graphs: Current status and future directions towards knowledge graph chatbots. arXiv:2302.06466 (2023)

[116] Diefenbach, D., Lopez, V., Singh, K., Maret, P.: Core techniques of question answering systems over knowledge bases: a survey. Knowledge and Information Systems **55**, 529–569 (2018)

[117] Shabbir, J., Arshad, M.U., Shahzad, W.: Nubot: Embedded knowledge graph with rasa framework for generating semantic intents responses in Roman Urdu. arXiv:2102.10410 (2021)

[118] Rospigliosi, P..: Artificial intelligence in teaching and learning: what questions should we ask of ChatGPT? Interactive Learning Environments **31**(1), 1–3 (2023)

[119] Rahaman, M., Ahsan, M., Anjum, N., Rahman, M., Rahman, M.N., et al.: The AI race is on! Google's Bard and OpenAI's ChatGPT head to head: An opinion article. Mizanur and Rahman, Md Nafizur, TheAI Race is on (2023)

[120] Zou, L., Zhang, S., Cai, H., Ma, D., Cheng, S., Wang, S., Shi, D., Cheng, Z., Yin, D.: Pre-trained language model based ranking in Baidu search. In: Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining, pp. 4014–4022 (2021)

[121] Menick, J., Trebacz, M., Mikulik, V., Aslanides, J., Song, F., Chadwick, M., Glaese, M., Young, S., Campbell-Gillingham, L., Irving, G., et al.: Teaching language models to support answers with verified quotes. arXiv:2203.11147 (2022)

[122] Anderson, N., Belavy, D.L., Perle, S.M., Hendricks, S., Hespanhol, L., Verhagen, E., Memon, A.R.: AI did not write this manuscript, or did it? Can we trick the AI text detector into generated texts? The potential future of ChatGPT and AI in sports & exercise medicine manuscript generation. BMJ Open Sport & Exercise Medicine **9**(1), 001568 (2023)

[123] GPTZero: GPTZero: The Global Standard for AI Detection Humans Deserve the Truth. GPTZero (2023). https://gptzero.me/ Accessed July 5, 2023

[124] GPTRadar: Detect AI-generated text in a click. GPTRadar (2023). https://gptradar.com/ Accessed July 5, 2023

[125] Turnitin: Turnitin's AI writing detection available now. Turnitin (2023). https://www.turnitin.com/ Accessed July 5, 2023

[126] Originally.AI: Most Accurate AI & Plagiarism Detector for Serious Content Publishers. Originally.AI (2023). https://originality.ai/ Accessed July 5, 2023

[127] Wiggers, K.: Most sites claiming to catch AI-written text fail spectacularly. TechCrunch (2023). https://techcrunch.com/ Accessed July 5, 2023

[128] Sadasivan, V.S., Kumar, A., Balasubramanian, S., Wang, W., Feizi, S.: Can AI-generated text be reliably detected? arXiv:2303.11156 (2023)

[129] Kirchenbauer, J., Geiping, J., Wen, Y., Katz, J., Miers, I., Goldstein, T.: A watermark for large language models. arXiv:2301.10226 (2023)

[130] Mitrović, S., Andreoletti, D., Ayoub, O.: ChatGPT or human? Detect and explain. Explaining decisions of machine learning model for detecting short ChatGPT-generated text. arXiv:2301.13852 (2023)

[131] Zhang, W., Deng, Y., Liu, B., Pan, S.J., Bing, L.: Sentiment analysis in the era of large language models: A reality check. arXiv:2305.15005 (2023)

[132] Abdullah, M., Madain, A., Jararweh, Y.: Chatgpt: Fundamentals, applications and social impacts. In: 2022 Ninth International Conference on Social Networks Analysis, Management and Security (SNAMS), pp. 1–8 (2022). IEEE

[133] Subagja, A.D., Ausat, A.M.A., Sari, A.R., Wanof, M.I., Suherlan, S.: Improving customer service quality in MSMEs through the use of ChatGPT. Jurnal Minfo Polgan **12**(2), 380–386 (2023)

[134] M Alshater, M.: Exploring the role of artificial intelligence in enhancing academic performance: A case study of ChatGPT. Available at SSRN (2022)

[135] Deng, J., Lin, Y.: The benefits and challenges of ChatGPT: An overview. Frontiers in Computing and Intelligent Systems **2**(2), 81–83 (2022)

[136] Hassani, H., Silva, E.S.: The role of ChatGPT in data science: how AI-assisted conversational interfaces are revolutionizing the field. Big Data and Cognitive Computing **7**(2), 62 (2023)

[137] Zhai, X.: ChatGPT for next generation science learning. XRDS: Crossroads, The ACM Magazine for Students **29**(3), 42–46 (2023)

[138] Kocoń, J., Cichecki, I., Kaszyca, O., Kochanek, M., Szydło, D., Baran, J., Bielaniewicz, J., Gruza, M., Janz, A., Kanclerz, K., et al.: ChatGPT: Jack of all trades, master of none. arXiv:2302.10724 (2023)

[139] Hariri, W.: Unlocking the potential of ChatGPT: A comprehensive exploration of its applications. Technology **15**(2), 16 (2023)

[140] Haleem, A., Javaid, M., Singh, R.P.: An era of ChatGPT as a significant futuristic support tool: A study on features, abilities, and challenges. BenchCouncil Transactions on Benchmarks, Standards and Evaluations **2**(4), 100089 (2022)

[141] Azaria, A.: ChatGPT Usage and Limitations (2022). https://hal.science/hal-03913837

[142] Yang, X., Li, Y., Zhang, X., Chen, H., Cheng, W.: Exploring the limits of ChatGPT for query or aspect-based text summarization. arXiv:2302.08081 (2023)

[143] Shahriar, S., Hayawi, K.: Let's have a chat! a conversation with ChatGPT: Technology, applications, and limitations. arXiv:2302.13817 (2023)

[144] Aljanabi, M., Ghazi, M., Ali, A.H., Abed, S.A., *et al.*: ChatGPT: Open possibilities. Iraqi Journal For Computer Science and Mathematics **4**(1), 62–64 (2023)

[145] Huang, F., Kwak, H., An, J.: Is ChatGPT better than human annotators? Potential and limitations of ChatGPT in explaining implicit hate speech. arXiv:2302.07736 (2023)

[146] Thorp, H.H.: ChatGPT is fun, but not an author. Science **379**(6630), 313–313 (2023)

[147] Guo, B., Zhang, X., Wang, Z., Jiang, M., Nie, J., Ding, Y., Yue, J., Wu, Y.: How close is ChatGPT to human experts? Comparison corpus, evaluation, and detection. arXiv:2301.07597 (2023)

[148] Susnjak, T.: ChatGPT: The end of online exam integrity? arXiv:2212.09292 (2022)

[149] Kohnke, L., Moorhouse, B.L., Zou, D.: ChatGPT for language teaching and learning. RELC Journal, 00336882231162868 (2023)

[150] Farrokhnia, M., Banihashem, S.K., Noroozi, O., Wals, A.: A swot analysis of ChatGPT: Implications for educational practice and research. Innovations in Education and Teaching International, 1–15 (2023)

[151] Bishop, L.: A computer wrote this paper: What ChatGPT means for education, research, and writing. Research, and Writing (January 26, 2023) (2023)

[152] Atlas, S.: ChatGPT for higher education and professional development: A guide to conversational AI (2023)

[153] Lin, Z.: Why and how to embrace AI such as ChatGPT in your academic life (2023)

[154] Fergus, S., Botha, M., Ostovar, M.: Evaluating academic answers generated using ChatGPT. Journal of Chemical Education (2023)

[155] Costello, E.: ChatGPT and the educational AI chatter: Full of bullshit or trying to tell us something? Postdigital Science and Education, 1–6 (2023)

[156] Cox, C., Tzoc, E.: ChatGPT: Implications for academic libraries. College & Research Libraries News **84**(3), 99 (2023)

[157] Firat, M.: How chat gpt can transform autodidactic experiences and open education. Department of Distance Education, Open Education Faculty, Anadolu Unive (2023)

[158] Lee, H.: The rise of ChatGPT: Exploring its potential in medical education. Anatomical Sciences Education (2023)

[159] Sallam, M., Salim, N., Barakat, M., Al-Tammemi, A.: ChatGPT applications in medical, dental, pharmacy, and public health education: A descriptive study highlighting the advantages and limitations. Narra J **3**(1), 103–103 (2023)

[160] Hisan, U.K., Amri, M.M.: ChatGPT and medical education: A double-edged sword. Journal of Pedagogy and Education Science **2**(01), 71–89 (2023)

[161] Gilson, A., Safranek, C.W., Huang, T., Socrates, V., Chi, L., Taylor, R.A., Chartash, D., *et al.*: How does ChatGPT perform on the united states medical licensing examination? the implications of large language models for medical education and knowledge assessment. JMIR Medical Education **9**(1), 45312 (2023)

[162] Bitzenbauer, P.: ChatGPT in physics education: A pilot study on easy-to-implement activities. Contemporary Educational Technology **15**(3) (2023)

[163] Qadir, J.: Engineering education in the era of ChatGPT: Promise and pitfalls of generative AI for education (2022)

[164] Halaweh, M.: ChatGPT in education: Strategies for responsible implementation (2023)

[165] Tlili, A., Shehata, B., Adarkwah, M.A., Bozkurt, A., Hickey, D.T., Huang, R., Agyemang, B.: What if the devil is my guardian angel: ChatGPT as a case study of using chatbots in education. Smart Learning Environments **10**(1), 15 (2023)

[166] Kasneci, E., Seßler, K., Küchemann, S., Bannert, M., Dementieva, D., Fischer, F., Gasser, U., Groh, G., Günnemann, S., Hüllermeier, E., *et al.*: ChatGPT for good? on opportunities and challenges of large language models for education. Learning and Individual Differences **103**, 102274 (2023)

[167] Sallam, M.: ChatGPT utility in healthcare education, research, and practice: Systematic review on the promising perspectives and valid concerns. In: Healthcare, vol. 11, p. 887 (2023). MDPI

[168] Rudolph, J., Tan, S., Tan, S.: ChatGPT: Bullshit spewer or the end of traditional assessments in higher education? Journal of Applied Learning and Teaching **6**(1) (2023)

[169] Ivanov, S., Soliman, M.: Game of algorithms: ChatGPT implications for the future of tourism education and research. Journal of Tourism Futures (2023)

[170] King, M.: Can GPT-4 formulate and test a novel hypothesis? Yes and no (2023)

[171] Quintans-Júnior, L.J., Gurgel, R.Q., Araújo, A.A.d.S., Correia, D., Martins-Filho, P.R.: ChatGPT: The new panacea of the academic world. Revista da Sociedade Brasileira de Medicina Tropical **56**, 0060–2023 (2023)

[172] Macdonald, C., Adeloye, D., Sheikh, A., Rudan, I.: Can ChatGPT draft a research article? An example of population-level vaccine effectiveness analysis. Journal of Global Health **13**, 01003 (2023)

[173] Alkaissi, H., McFarlane, S.I.: Artificial hallucinations in ChatGPT: Implications in scientific writing. Cureus **15**(2) (2023)

[174] Teubner, T., Flath, C.M., Weinhardt, C., Aalst, W., Hinz, O.: Welcome to the era of ChatGPT et al. the prospects of large language models. Business & Information Systems Engineering, 1–7 (2023)

[175] Aydın, Ö., Karaarslan, E.: OpenAI ChatGPT generated literature review: Digital twin in healthcare. Available at SSRN 4308687 (2022)

[176] Reed, D.A., Zhao, Y., Bagheri Varzaneh, M., Soo Shin, J., Rozynek, J., Miloro, M., Han, M.: Ng2/cspg4 regulates cartilage degeneration during tmj osteoarthritis. Frontiers in Dental Medicine, 69 (2022)

[177] BagheriVarzaneh, M., Zhao, Y., Rozynek, J., Han, M., Reed, D.: Disrupting mechanical homeostasis promotes matrix metalloproteinase-13 mediated processing of neuron glial antigen 2 in mandibular condylar cartilage

[178] Eggmann, F., Weiger, R., Zitzmann, N.U., Blatz, M.B.: Implications of large language models such as ChatGPT for dental medicine. Journal of Esthetic and Restorative Dentistry (2023)

[179] Lyu, Q., Tan, J., Zapadka, M.E., Ponnatapuram, J., Niu, C., Wang, G., Whitlow, C.T.: Translating radiology reports into plain language using ChatGPT and GPT-4 with prompt learning: Promising results, limitations, and potential. arXiv:2303.09038 (2023)

[180] Lecler, A., Duron, L., Soyer, P.: Revolutionizing radiology with GPT-based models: Current applications, future possibilities and limitations of ChatGPT. Diagnostic and Interventional Imaging (2023)

[181] Varzaneh, M.B., Rahmani, H., Jahanian, R., Mahdavi, A.H., Perreau, C., Perrot, G., Brézillon, S., Maquart, F.-X.: The influence of oral copper-methionine on matrix metalloproteinase-2 gene expression and activation in right-sided heart failure induced by cold temperature: A broiler chicken perspective. Journal of Trace Elements in Medicine and Biology **39**, 71–75 (2017)

[182] Bagheri Varzaneh, M., Rahmani, H., Jahanian, R., Mahdavi, A.H., Perreau, C., Perrot, G., Brézillon, S., Maquart, F.-X.: Effects of dietary copper-methionine on matrix metalloproteinase-2 in the lungs of cold-stressed broilers as an animal model for pulmonary hypertension. Biological Trace Element Research **172**, 504–510 (2016)

[183] Sabeti, P.C., Reich, D.E., Higgins, J.M., Levine, H.Z., Richter, D.J., Schaffner, S.F., Gabriel, S.B., Platko, J.V., Patterson, N.J., McDonald, G.J., *et al.*: Detecting recent positive selection in the human genome from haplotype structure. Nature **419**(6909), 832–837 (2002)

[184] Biswas, S.S.: Role of ChatGPT in public health. Annals of Biomedical Engineering, 1–2 (2023)

[185] Luo, R., Sun, L., Xia, Y., Qin, T., Zhang, S., Poon, H., Liu, T.-Y.: BioGPT: Generative pre-trained transformer for biomedical text generation and mining. Briefings in Bioinformatics **23**(6) (2022)

[186] Li, J., Dada, A., Kleesiek, J., Egger, J.: ChatGPT in healthcare: A taxonomy and systematic review. medRxiv, 2023–03 (2023)

[187] Jin, Q., Yang, Y., Chen, Q., Lu, Z.: Genegpt: Augmenting large language models with domain tools for improved access to biomedical information. ArXiv (2023)

[188] Vaishya, R., Misra, A., Vaish, A.: ChatGPT: Is this version good for healthcare and research? Diabetes & Metabolic Syndrome: Clinical Research & Reviews **17**(4), 102744 (2023)

[189] Cascella, M., Montomoli, J., Bellini, V., Bignami, E.: Evaluating the feasibility of ChatGPT in healthcare: an analysis of multiple clinical and research scenarios. Journal of Medical Systems **47**(1), 1–5 (2023)

[190] Ufuk, F.: The role and limitations of large language models such as ChatGPT in clinical settings and medical journalism. Radiology, 230276 (2023)

[191] Asch, D.A.: An interview with ChatGPT about health care. NEJM Catalyst Innovations in Care Delivery **4**(2) (2023)

[192] Patel, S.B., Lam, K.: ChatGPT: the future of discharge summaries? The Lancet Digital Health **5**(3), 107–108 (2023)

[193] King, M.R.: The future of AI in medicine: A perspective from a chatbot. Annals of Biomedical Engineering **51**(2), 291–295 (2023)

[194] Biswas, S.: ChatGPT and the future of medical writing. Radiology, 223312 (2023)

[195] Ali, S.R., Dobbs, T.D., Hutchings, H.A., Whitaker, I.S.: Using ChatGPT to write patient clinic letters. The Lancet Digital Health **5**(4), 179–181 (2023)

[196] Sharma, M., Sharma, S.: Transforming maritime health with ChatGPT-powered healthcare services for mariners. Annals of Biomedical Engineering, 1–3 (2023)

[197] Iftikhar, L., *et al.*: Docgpt: Impact of ChatGPT-3 on health services as a virtual doctor. EC Paediatrics **12**(1), 45–55 (2023)

[198] Mann, D.L.: Artificial intelligence discusses the role of artificial intelligence in translational medicine: A jacc: Basic to translational science interview with ChatGPT. Basic to Translational Science (2023)

[199] Hopkins, A.M., Logan, J.M., Kichenadasse, G., Sorich, M.J.: Artificial intelligence chatbots will revolutionize how cancer patients access information: ChatGPT represents a paradigm-shift. JNCI Cancer Spectrum **7**(2), 010 (2023)

[200] Mijwil, M., Aljanabi, M., Ali, A.H.: ChatGPT: Exploring the role of cybersecurity in the protection of medical information. Mesopotamian Journal of Cybersecurity **2023**, 18–21 (2023)

[201] George, A.S., George, A.H.: A review of ChatGPT AI's impact on several business sectors. Partners Universal International Innovation Journal **1**(1), 9–23 (2023)

[202] Dowling, M., Lucey, B.: ChatGPT for (Finance) research: The bananarama conjecture. Finance Research Letters **53**, 103662 (2023)

[203] Zaremba, A., Demir, E.: ChatGPT: Unlocking the future of NLP in finance. Available at SSRN 4323643 (2023)

[204] Wenzlaff, K., Spaeth, S.: Smarter than humans? Validating how OpenAIs ChatGPT model explains crowdfunding, alternative finance and community finance.

Validating How OpenAIs ChatGPT Model Explains Crowdfunding, Alternative Finance and Community Finance (2022)

[205] Yue, T., Au, D., Au, C.C., Iu, K.Y.: Democratizing financial knowledge with ChatGPT by openai: Unleashing the power of technology. Available at SSRN 4346152 (2023)

[206] Alshurafat, H.: The usefulness and challenges of chatbots for accounting professionals: application on ChatGPT. Available at SSRN 4345921 (2023)

[207] Ali, H., Aysan, A.F.: What will ChatGPT revolutionize in financial industry? Available at SSRN 4403372 (2023)

[208] Chuma, E., Bang, M., Alfredson, J.: Business AI decision-making tools: Case ChatGPT evaluation (2023)

[209] Rathore, B.: Future of textile: Sustainable manufacturing & prediction via Chat-GPT. Eduzone: International Peer Reviewed/Refereed Multidisciplinary Journal **12**(1), 52–62 (2023)

[210] Chui, M., Roberts, R., Yee, L.: Generative AI is here: How tools like ChatGPT could change your business. Quantum Black AI by McKinsey (2022)

[211] Singh, D.: ChatGPT: A new approach to revolutionise organisations. International Journal of New Media Studies (IJNMS) **10**(1), 57–63 (2023)

[212] Pettinato Oltz, T.: ChatGPT, professor of law. Professor of Law (2023)

[213] Armstrong, A.B.: Whos afraid of ChatGPT? an examination of chatgpts implications for legal writing. An Examination of ChatGPTs Implications for Legal Writing (2023)

[214] Hargreaves, S.: words are flowing out like endless rain into a paper cup: Chatgpt & law school assessments. The Chinese University of Hong Kong Faculty of Law Research Paper (2023-03) (2023)

[215] Choi, J.H., Hickman, K.E., Monahan, A., Schwarcz, D.: ChatGPT goes to law school. Available at SSRN (2023)

[216] Biswas, S.S.: Potential use of chat gpt in global warming. Annals of Biomedical Engineering, 1–2 (2023)

[217] Jiao, W., Wang, W., Huang, J.-t., Wang, X., Tu, Z.: Is ChatGPT a good translator? a preliminary study. arXiv:2301.08745 (2023)

[218] Hoes, E., Altay, S., Bermeo, J.: Using ChatGPT to fight misinformation: ChatGPT nails 72% of 12,000 verified claims (2023)

[219] McGee, R.W.: Political philosophy and ChatGPT. Technical report, DOI: 10.13140 (2023)

[220] Susnjak, T.: Applying bert and ChatGPT for sentiment analysis of lyme disease in scientific literature. arXiv:2302.06474 (2023)

[221] McGee, R.W.: Capitalism, socialism and ChatGPT. Available at SSRN 4369953 (2023)

[222] Chomsky, N., Roberts, I., Watumull, J.: Noam chomsky: The false promise of ChatGPT. The New York Times **8** (2023)

[223] Samuel, A.: A Guide to Collaborating With ChatGPT for Work. The Wall Street Journal (2023). https://www.wsj.com/ Accessed July 5, 2023

[224] Papachristou, L., Deutsch, J.: ChatGPT Advances Are Moving So Fast Regulators Cant Keep Up. Bloomberg (2023). https://www.bloomberg.com/ Accessed July 5, 2023

[225] NEWS, C.: 'Godfather of artificial intelligence' weighs in on the past and potential of AI. CBS News (2023). https://www.cbsnews.com/ Accessed July 5, 2023

[226] Moran, C.: ChatGPT is making up fake Guardian articles. Heres how were responding. The Guardian (2023). https://www.theguardian.com/ Accessed July 5, 2023

[227] Economist: ChatGPT could replace telemarketers, teachers and traders. The Economist (2023). https://www.economist.com/ Accessed July 5, 2023

[228] Brown, J.: ChatGPT delivers sermon to packed German church, tells congregants not to fear death. Fox News (2023). https://www.foxnews.com/ Accessed July 5, 2023

[229] Anderson, J., Rainie, L.: Closing thoughts on ChatGPT and other steps in the evolution of humans, digital tools and systems by 2035. Pew Research Center (2023). https://www.pewresearch.org/ Accessed July 5, 2023

[230] Roose, K.: Dont Ban ChatGPT in Schools. Teach With It. The New York Times (2023). https://www.nytimes.com/ Accessed July 5, 2023

[231] O'Brien, M.: EXPLAINER: What is ChatGPT and why are schools blocking it? AP News (2023). https://apnews.com/ Accessed July 5, 2023

[232] McCallum, S.: ChatGPT banned in Italy over privacy concerns. BBC (2023). https://www.bbc.com/ Accessed July 5, 2023

[233] Johnson, K.: Spooked by ChatGPT, US Lawmakers Want to Create an AI

Regulator. Wired (2023). https://www.wired.com/ Accessed July 5, 2023

[234] Reed, B.: Authors file a lawsuit against OpenAI for unlawfully ingesting their books. The Guardian (2023). https://www.theguardian.com/ Accessed July 5, 2023

[235] Meaker, M., Johnson, K.: The Global Battle to Regulate AI Is Just Beginning. Wired (2023). https://www.wired.com/ Accessed July 5, 2023

[236] Ambartsoumean, V.M., Yampolskiy, R.V.: AI risk skepticism, a comprehensive survey. arXiv:2303.03885 (2023)

[237] Mujtaba, G., Malik, A., Ryu, E.-S.: LTC-SUM: Lightweight client-driven personalized video summarization framework using 2D CNN. IEEE Access **10**, 103041–103055 (2022)