

Large Language Models: A Comprehensive Survey of its Applications, Challenges, Limitations, and Future Prospects

Muhammad Usman Hadi¹, Qasem Al Tashi¹, Rizwan Qureshi², Abbas Shah¹, Amgad Muneer¹, Muhammad Irfan¹, Anas Zafar¹, Muhammad Bilal Shaikh¹, Naveed Akhtar¹, Syed Zohaib Hassan¹, Maged Shoman¹, Jia Wu¹, Seyedali Mirjalili¹, and Mubarak Shah¹

¹Affiliation not available

²MD Anderson Cancer Center

July 10, 2023

Abstract

Within the vast expanse of computerized language processing, a revolutionary entity known as Large Language Models (LLMs) has emerged, wielding immense power in its capacity to comprehend intricate linguistic patterns and conjure coherent and contextually fitting responses. LLMs are a type of artificial intelligence (AI) that have emerged as powerful tools for a wide range of tasks, including natural language processing (NLP), machine translation, vision applications, and question-answering. This survey provides a comprehensive overview of LLMs, including their history, architecture, training methods, applications, and challenges. We begin by discussing the fundamental concepts of generative AI and the architecture of generative pre-trained transformers (GPT). We then provide an overview of the history of LLMs, their evolution over time, and the different training methods that have been used to train them. We then discuss the wide range of tasks where they are used and also discuss applications of LLMs in different domains, including medicine, education, finance, engineering, media, entertainment, politics, and law. We also discuss how LLMs are shaping the future of AI and their increasing role in scientific discovery, and how they can be used to solve real-world problems. Next, we explore the challenges associated with deploying LLMs in real-world scenarios, including ethical considerations, model biases, interpretability, and computational resource requirements. This survey also highlights techniques for enhancing the robustness and controllability of LLMs and addressing bias, fairness, and quality issues in Generative AI. Finally, we conclude by highlighting the future of LLM research and the challenges that need to be addressed in order to make this technology more reliable and useful. This survey is intended to provide researchers, practitioners, and enthusiasts with a comprehensive understanding of LLMs, their evolution, applications, and challenges. By consolidating the state-of-the-art knowledge in the field, this article is anticipated to serve as a valuable resource for learning the current state-of-the-art as well as further advancements in the development and utilization of LLMs for a wide range of real-world applications. The GitHub repo for this project is available at <https://github.com/anas-zafar/LLM-Survey>

A Survey on Large Language Models: Applications, Challenges, Limitations, and Practical Usage

Muhammad Usman Hadi^{1,*}, Qasem Al-Tashi^{2,12*}, Rizwan Qureshi^{2,*}, Abbas Shah^{3,*}, Amgad Muneer², Muhammad Irfan⁴, Anas Zafar⁵, Muhammad Bilal Shaikh⁶, Naveed Akhtar⁷, Mohammed Ali Al-Garadi⁸, Syed Zohaib Hassan⁹, Maged Shoman⁹, Jia Wu², and Seyedali Mirjalili^{10,11}, Mubarak Shah¹²

¹School of Engineering, Ulster University, Belfast, BT15 1AP, United Kingdom (m.hadi@ulster.ac.uk)

²Department of Imaging Physics, The University of Texas MD Anderson Cancer Center, Houston, TX 77030, USA (qaal@mdanderson.org; fnu.rizwan@ucf.edu; amabdulraheem@mdanderson.org; JWu11@mdanderson.org)

³Department of Electronics Engineering, Mehran University of Engineering and Technology, Jamshoro, 76062 Pakistan (zaigham.shah@faculty.muet.edu.pk)

⁴Faculty of Electrical Engineering, Ghulam Ishaq Khan Institute (GIKI) of Engineering Sciences and Technology, Swabi, 23460 Pakistan (mirfan@giki.edu.pk)

⁵Department of Computer Science, National University of Computer and Emerging Sciences, Karachi, Pakistan (anaszafar98@gmail.com)

⁶Center for Artificial Intelligence and Machine Learning (CAIML), Edith Cowan University, 270 Joondalup Drive, Joondalup, WA 6027, Perth, Australia (mbshaikh@our.ecu.edu.au)

⁷Computing and Information Systems, The University of Melbourne, 700 Swanston Street, Carlton 3010, VIC Australia (naveed.akhtar1@unimelb.edu.au)

⁸Department of Biomedical Informatics, Vanderbilt University Medical Center, Nashville, TN, USA

⁹Department of Civil, Environmental, and Construction Engineering, The University of Central Florida, Orlando, Florida, USA (zohaibneduet@gmail.com; magedshoman@gmail.com)

¹⁰Centre for Artificial Intelligence Research and Optimization, Torrens University Australia, Fortitude Valley, Brisbane, QLD 4006, Australia (ali.mirjalili@torrens.edu.au)

¹¹University Research and Innovation Center, Obuda University, 1034 Budapest, Hungary * Corresponding Author

¹²Center for Research in Computer Vision, The University of Central Florida, Orlando, Florida, USA (shah@crcv.ucf.edu)

Abstract

Within the vast expanse of computerized language processing, a revolutionary entity known as Large Language Models (LLMs) has emerged, wielding immense power in its capacity to comprehend intricate linguistic patterns and conjure coherent and contextually fitting responses. LLMs are a type of artificial intelligence (AI) that have emerged as powerful tools for a wide range of tasks, including natural language processing (NLP), machine translation, vision applications, and question-answering. This survey provides a comprehensive overview of LLMs, including their history, architecture, datasets, training methods, applications, challenges, and future prospects. We begin by discussing the fundamental concepts of generative AI and the architecture of generative pre-trained transformers (GPT). We then provide an overview of the history of LLMs, their evolution over time, and the different training methods. We also present benchmark dataset for training and fine-tuning and evaluating LLMs. We then discuss the wide range of tasks where they are used and also discuss applications of LLMs in different domains, including medicine, education, finance, engineering, agriculture, media, entertainment, politics, and law. We also discuss how LLMs are shaping the future of AI and their increasing role in scientific discovery, and how they can be used to solve real-world problems. Next, we explore the challenges associated with deploying LLMs in real-world scenarios, including ethical considerations, model biases, interpretability, privacy concerns, and computational resource requirements. This survey also highlights techniques for enhancing the robustness and controllability of LLMs and addressing bias, fairness, and quality issues in Generative AI. Finally, we conclude by highlighting the future of LLM research and the challenges that need to be addressed in order to make this technology more reliable and useful. This survey is intended to provide researchers, practitioners, and enthusiasts with a comprehensive understanding of LLMs, their evolution, applications, and challenges. By consolidating the state-of-the-art knowledge in the field, this article is anticipated to serve as a valuable resource for learning the current state-of-the-art as well as further advancements in the development and utilization of LLMs for a wide range of real-world applications. The GitHub repo for this project is available at Github-Repo.

Index Terms

Large Language Models, Large Vision Models, Generative AI, Conversational AI, LangChain, Natural language processing, Computer Vision, GPT, ChatGPT, Bard, AI chatbots

A Survey on Large Language Models: Applications, Challenges, Limitations, and Practical Usage

I. INTRODUCTION

Language modeling (LM) is a fundamental task in natural language processing (NLP) that aims to predict the next word or a character in a given sequence of text [1], [2]. It involves developing algorithms and models that can understand and generate coherent human language. The primary objective of LM is to capture the probability distribution of words in a language, which allows the model to generate new text [3], complete sentences [4], and predict the likelihood of different word sequences [5], [6]. They are broadly categorized into statistical language models, machine learning models, deep learning models, and transformer based models as shown in Fig. 1. Early language models, such as n-gram models [7], were based on simple statistical techniques that estimated the probabilities of word sequences using frequency counts [8], [9]. However, with the rise of deep learning in NLP [10], the availability of enormous amounts of public datasets [11], and powerful computing devices [12] to process these big data with complex algorithms, has led to the development of large language models.

Large Language Models (LLMs) [13], sometimes referred to as "transformative [14]" or "next-generation [15]" language models, represent a significant breakthrough in NLP [16]. These models leverage deep learning techniques, particularly transformer architectures [17], to learn and understand the complex patterns and structures present in language data [18]. A key characteristic of LLMs is their ability to process vast amounts of data, including unstructured text, and capture semantic relationships between words and phrases [19]. These models can also process visual [20], audio [21], audiovisual [22], as well as multi-modal data [23] and learn the semantic relationships between them. These models have significantly enhanced the capabilities of machines to understand and generate human-like language [24].

The history of LLMs can be traced back to the early development of language models and neural networks [25]. The journey begins with the era of statistical language models [26]. In this stage, researchers primarily relied on probabilistic approaches [27] to predict word sequences. Classic examples include n-grams, Hidden Markov Models (HMMs) [28] and Maximum Entropy Models [29]. N-grams, for instance, are sequences of adjacent words or tokens that are used to predict the likelihood of the next word based on the preceding ones [30]. While rudimentary by today's standards, these models marked a crucial starting point in the field of natural language understanding. They allowed for basic text generation and word prediction but were limited in their ability to capture complex contextual relationships [31] [32]. [33]. Then a shift towards more data-driven methodologies has

been witnessed [34]. Researchers began to explore machine learning algorithms to improve language understanding [35]. These models learned patterns and relationships within large text corpora. Support Vector Machines (SVMs) is a notable example from this [36]. Machine learning models brought a more sophisticated approach to NLP tasks, allowing for the development of applications like spam detection [37] and sentiment analysis [38]. Moreover the availability of large-scale Twitter¹ datasets has brought a revolution in real time sentiment analysis [39].

The emergence of deep learning marked a pivotal moment in the development of LLMs [40]. Neural networks, particularly Recurrent Neural Networks (RNNs) and Long Short-Term Memory (LSTM) networks gained prominence [41]. These deep learning architectures delved deeper into the data, allowing them to capture more intricate features and long-range dependencies within text. This stage significantly improved the models' ability to understand context, making them suitable for tasks like machine translation and speech recognition [16], [42]. However, deep learning also faced challenges with vanishing gradients [43] and long-term dependencies [44], limiting their effectiveness.

The breakthrough in LLMs came with the introduction of the Transformer architecture in the seminal work "Attention is All You Need" by Vaswani et al. in 2017 [45]. The Transformer model, based on the self-attention mechanism [46], enabled parallelization and efficient handling of long-range dependencies. It laid the foundation for models like OpenAI's GPT (Generative Pre-trained Transformer) series [47] and BERT (Bidirectional Encoder Representations from Transformers) [48] by Google, which achieved groundbreaking results in a wide range of language tasks. These mechanisms enabled models to consider the entire context of a sentence or document, allowing for true contextual understanding [49]. Transformer-based models, often pre-trained on massive text corpora, can generate coherent and contextually relevant text, revolutionizing applications like chatbots [50], text summarization [51], and language translation [52].

ChatGPT, Llama, and Falcon are all remarkable variants of the GPT (Generative Pre-trained Transformer) model [53], which is developed and pioneered by OpenAI. These models represent OpenAI's ongoing efforts to push the boundaries of natural language processing and understanding. They share a common foundation in their training methodology, which involves pre-training on vast corpora of text data followed by fine-tuning for specific tasks. During pre-training, the

¹X, as a microblogging platform, allows users to express their thoughts and opinions in short, concise messages called tweets. These tweets often contain rich, real-time information about various topics, making Twitter an excellent source for sentiment analysis

models are exposed to diverse internet text to learn grammar, facts, reasoning abilities, and some degree of common-sense knowledge [54], [55]. This process equips them with a broad understanding of language. Subsequently, fine-tuning is carried out on narrower datasets to specialize the models for particular applications. ChatGPT, for instance, is fine-tuned for conversational contexts, making it suitable for chatbots and virtual assistants [56], [57], [58]. Llama and Falcon, though not as widely known, represent potential advancements or specialized versions, possibly tailored for specific use cases or research objectives. These models collectively exemplify the cutting-edge advancements in natural language processing, enabling more human-like interactions and understanding through the power of AI-driven language models [59], [60], [61].

The training process for models like ChatGPT, Llama, and Falcon consists of several key stages [62], [53]. It begins with a phase known as pre-training, where these models are exposed to a vast and varied dataset of internet text, enabling them to learn grammar, vocabulary, world knowledge, and context [63]. The underlying architecture, based on the Transformer model, is crucial for understanding the relationships between words in sentences. Following pre-training, the models undergo fine-tuning on specific datasets tailored to particular tasks, such as text generation or conversation in the case of ChatGPT. Fine-tuning refines their capabilities for these specialized tasks, with hyperparameter tuning to optimize performance. Ethical considerations are an integral part of the process, aiming to minimize harmful or biased outputs. It's important to note that this training is an iterative and resource-intensive endeavor, continually improved and monitored to enhance both performance and safety [64], [65].

LLMs have undergone several developmental stages, with models increasing in size and complexity. The GPT series, starting with GPT-1 and continuing with GPT-2 and GPT-3 [66], has successively grown in the number of parameters, starting from hundreds of millions (GPT-1) to 1.7 Trillion (GPT-4) [67], allowing for more sophisticated language understanding and generation capabilities [68]. Similarly, BERT-inspired models have seen advancements in pre-training strategies, such as ALBERT [69] (A Lite BERT) and RoBERTa [70], which further improved performance and efficiency.

Furthermore, advancements in LLMs have extended to more specific domains, with models designed for specialized tasks like medical language processing [71], scientific research [72], website development [73] and code generation [74]. Moreover, efforts have been made to address ethical concerns [75], interpretability [76], and reducing biases in LLMs to ensure responsible and equitable use [77]. The development stages of large models have witnessed a constant quest for larger models, improved pre-training strategies, and specialized domain adaptations [78], [79]. As research continues, the potential applications and impact of LLMs on various fields, including education, healthcare, and human-computer interaction, continue to expand, inspiring further innovations and advancements.

In summary and as can be seen from Fig 1; LM research has received widespread attention and has undergone four significant development stages including: statistical language

models, machine learning models, deep learning models and transformer-based models². In this research, we mainly focus on LLMs and foundation AI models for language and vision tasks.

Modern language model called ChatGPT [80] was developed by OpenAI [81] and launched in 2022. It is based on the GPT-3.5 architecture [82] and was trained using a sizable amount of internet-sourced text data, including books, articles, wikis and websites (Table I) [83]. ChatGPT is exceptional at producing human-like responses and having conversations with users.

In computer vision (CV), researchers are also actively engaged in the development of vision-language models inspired by the capabilities of ChatGPT. These models are specifically designed to enhance multimodal dialogues, where both visual and textual information are important [84]. Moreover, the advancements in the field have led to the introduction of GPT-4 [82], which has further expanded the capabilities of language models by seamlessly integrating visual information as part of the input. This integration of visual data empowers the model to effectively understand and generate responses that incorporate both textual and visual cues, enabling more contextually rich and nuanced conversations in multimodal settings.

A. Survey Motivation

The revolutionary ChatGPT has captivated the attention of the community, sparking a wealth of fascinating reviews and discussions on the advancements of LLMs and artificial intelligence [85], [86], [87], [88], [89], [90], [91]. For example, the role of ChatGPT in education is evaluated in [92], healthcare and medicine in [93], [71], protein sequence modeling in [94] and protein generation in [95]. A survey on generative AI is presented in [96], scientific text modeling in [97] and text generation in [98]. The use of LLM in finance is evaluated in [99], impact on labor market in [100] and supply chain in [101], telecom in [102], on code writing capabilities in [103], deep fakes in [104], legal aspects in [105], AI for drug discovery in [106], clinical prediction with LLM in [107], ML for cancer biomarkers in [108], and integration of biotechnology and AI applications to address global challenges in [109]. The advancements in pre-training, fine-tuning, utilization and capability evaluation of LLMs is presented in [85] and a survey on autonomous agents in [110]. The recent progress in visio-language pre-trained models is discussed in [86] and the knowledge graphs construction and reasoning are explained in [111], selection inference is in [76], and quantum-inspired machine learning in [112]. The survey vision language pre-trained models [86] presents an overview of various techniques for encoding raw images and texts into single-modal embeddings as a fundamental aspect, and also discusses prevalent architectures of vision-language pre-trained models (VL-PTMs), focusing on their ability to effectively model the interaction between text and image representations.

²Due to the dominance of Transformer based models, we consider it a different stage, not as a subset of deep learning

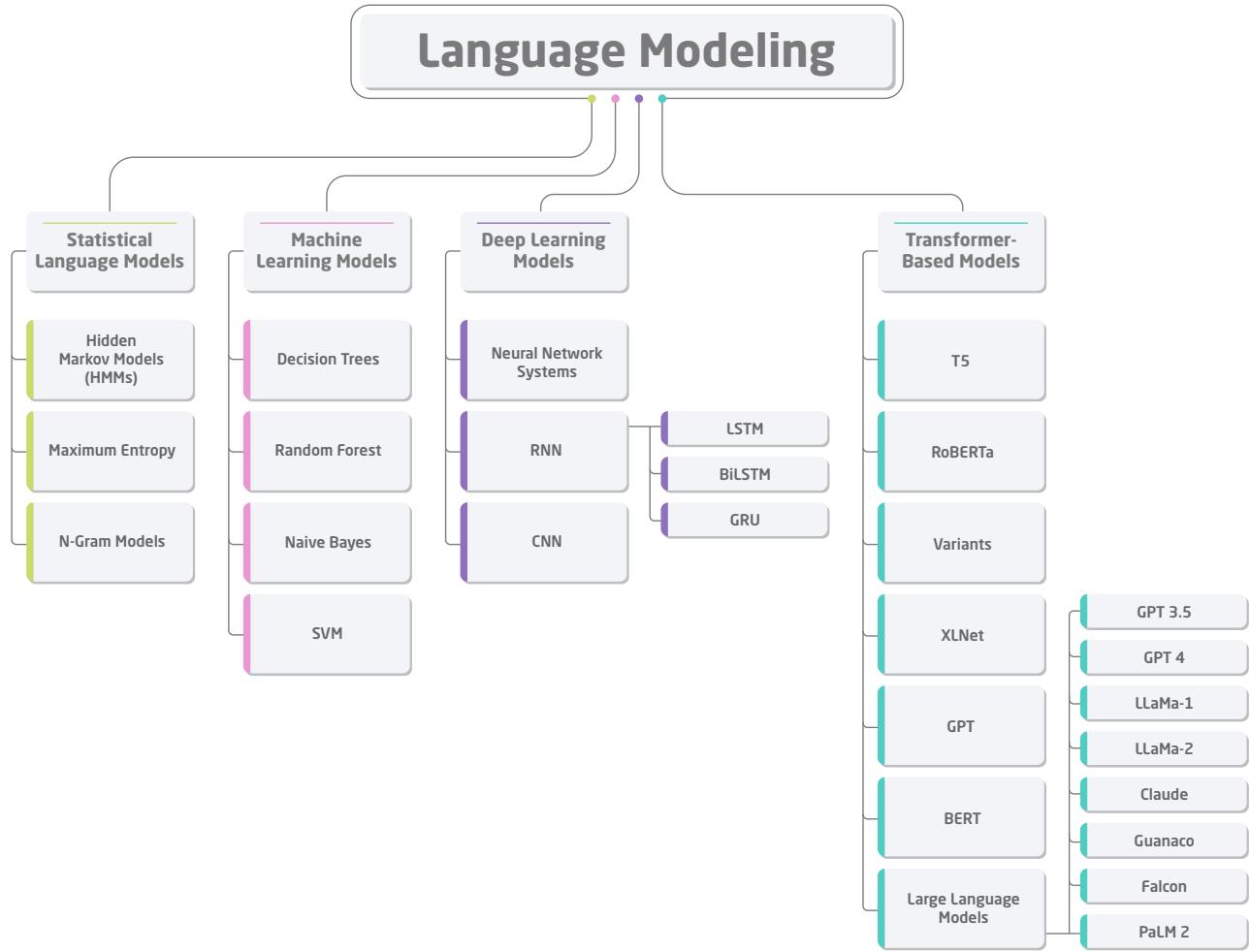


Fig. 1: Types of language modeling. The division of LLMs is categorized into four major blocks: Statistical language models, Machine learning models, Deep learning models and Transformer-based models.

Despite the growing number of studies on LLMs, there remains a scarcity of research focusing on their technical intricacies and effective utilization. Also, the field is progressing at a very fast pace, so a review article with practical applications will contribute a lot to the field. Therefore we also write this paper in the form of an application oriented review. Our primary objective is to explore, learn, and evaluate language models across various domains. We delve into the working principles of language models, analyze different architectures of the GPT family and others, and discuss strategies for their optimal utilization. Furthermore, we provide detailed insights about generative AI, writing prompts, and visual prompting techniques, leveraging GPT-plug-ins, and harnessing other AI/LLM tools. These aspects are generally not covered by the existing related articles. Our comprehensive examination also encompasses a discussion on the limitations associated with the LLMs, including considerations related to security, ethics, economy, and the environment. In addition, we present a set of guidelines to steer future research and development in the effective use of LLMs. We hope that this paper will contribute to a better understanding and utilization of LLMs.

B. Contributions

The main contributions of this article are as follows:

- 1) Providing a comprehensive overview of GenAI and LLMs, including their technical details, advancements, challenges, capabilities and limitations. We present state of the art analysis and comparisons of different LLMs.
- 2) Addressing ethical concerns about LLMs, including their computational requirements and potential for perpetuating biases. We also discuss the limitations of LLMs; including, limited understanding of the physical world, tokenization problems, infomration hallucination, finetuning and risk of foundation models.
- 3) Offering insights into the future potential of LLMs and their impact on society and demonstrating the applications of LLM through four practical use cases in the fields of medicine, education, finance, law, politics, media, entertainment, engineering, and others.
- 4) This article is uniquely presented in a manner to promote practical usage of LLMs, showcasing the actual LLM outputs to corroborate the discussions.

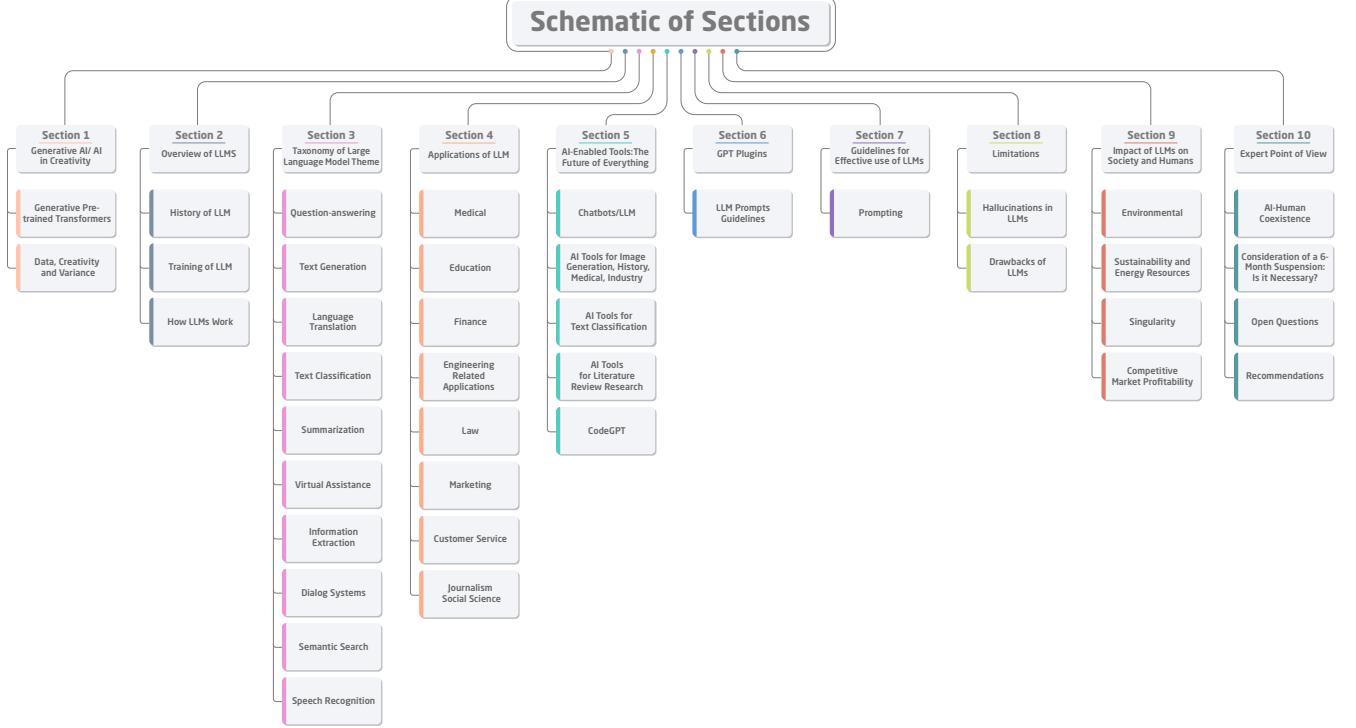


Fig. 2: Schematic Representation of Article Sections - A visual overview of the key sections comprising the structure of the article, providing readers with a roadmap for navigating the content effectively.

TABLE I: **Pre-training data.** Mixtures of data used for pre-training LLaMA [15].

Dataset	Sampling prop.	Epochs	Disk size
CommonCrawl	67.0%	1.10	3.3TB
C4	15.0%	1.06	783GB
Github	4.5%	0.64	328GB
Wikipedia	4.5%	2.45	83GB
Books	4.5%	2.23	85GB
ArXiv	2.5%	1.06	92GB
Stock Exchange	2.0%	1.03	78GB

The paper is organized as the following sections. Section II provides an introduction to the role of generative AI. Section III presents an overview of LLMs, summarizing a brief history of LLMs and discussing their training and functionality. Benchmark datasets for LLM training and instruction fine-tuning are presented in Section IV. Taxonomy of LLMs is presented in Section V and major applications of LLM through different use cases is discussed in Section VI. Section VIII explores AI-enabled tools. Section IX discusses the practical use cases of GPT plugins. Section X presents guidelines and working examples using prompting techniques. Section XI proposes the limitations and drawbacks of the current state-of-the-art LLM. Section XII presents open questions on the subject matter and the authors' perspective on open unanswered avenues. Section XIII concludes the survey paper. The overall structure of the article is presented in Fig. 2 for a quick reference at a glance.

II. GENERATIVE AI

Generative AI (GenAI) [113] is perhaps the most disruptive [114] and generalized technology of this decade [115], already influenced many industries, including, Media [116], Marketing [117], Game development and Metaverse [118], Education [119], Software development [120], and Medical [121], construction technology [122], and pharmaceuticals [123]. Unlike general AI systems that perform specific tasks such as data classification [124], clustering [125], object detection [126] and segmentation [127] or predictions [128]; GenAI can generate meaningful new content of multiple data modalities [129]; including, text [3], speech [130], images [131], and videos [132]. Some common examples of GenAI systems are image generators (Midjourney or stable diffusion), Chatbots (ChatGPT, Bard, Palm), code generators (CodeX, Co-Pilot [133]) audio generators(VALL-E)Vall-e [134], and video generators (Gen-2) [135].

During the past few years, GenAI models size has been scaled from a few million parameters(BERT [48], 110M) to hundreds of billions of parameters (GPT [136], 175B). Generally speaking, as the size of the model (number of parameters) increases, the performance of the model also increases [137], and it can be generalized for a variety of tasks [138], for example, Foundation models [139]. However, smaller models can also be fine-tuned for a more focused task [140].

LLMs, such as ChatGPT by OpenAI, Bard by Google, and Llama by Meta, are a type of GenAI models, specifically designed to generate human-like language in response to a

given prompt [141]. These models are trained on massive amounts of data (see Table I), using techniques to learn the statistical patterns of language. However, many people accord the capabilities provided by GPT models to “more data and computing power” instead of “better ML research” [142].

GenAI works by leveraging complex algorithms and statistical models to generate new content that mimics the patterns and characteristics of the training data [143]. These algorithms may include probabilistic techniques; such as Autoregressive model [144] and Variations Auto-encoders [145], or more recently, Generative Adversarial Networks [146] and Diffusion models [147] or Reinforcement Learning Human Feedback (RLHF) [148].

GenAI has captured significant interest in recent years due to its remarkable performance across an extensive array of applications in text, image, video and generation [149]. Constructed upon the foundation of the transformer architecture [45], these models exhibit an extraordinary capacity to process and generate human-like content by leveraging massive volumes of training data for various topics [150].

A. Data, Generation, Variance, and Performance measures

To comprehend the intricacies of GenAI systems, it is important to delve into the concepts of data, generation, and variance, and the interplay between them, as they form the foundation of generative systems [151].

1) Data: The core of generative AI systems is data. Training models that can successfully capture the underlying patterns and structures of the target domain require high-quality and diverse training data. The generating performance is influenced by the amount, quality, and representation of the training data [152], [153], [154], [155]. Furthermore, the availability of large-scale, labeled datasets allows for the development of more accurate and coherent samples [156], whereas restricted or biased training data may yield sub-optimal results [157].

2) Generation Process: GenAI uses the gained knowledge from the training data to generate samples with similar statistical patterns [158]. The generative models are designed to capture the underlying distributions of the training data and generate reliable and realistic samples with properties consistent with the original dataset [159]. The generating process involves approaches; such as adversarial training [160], latent space interpolation [161], and autoregressive modeling [162].

3) Variance: Variance is another important factor in defining the diversity and quality of generated samples [163], which shows the variability in the generated samples. A low variance generative AI system may produce similar or repetitive samples, resulting in poor generation, whereas, a high variance, may yield diversified but unrealistic or incoherent samples [164]. Striking a balance between variation and fidelity is difficult in generative AI [165], since it requires managing the trade-off between exploring and exploiting the learnt data distribution [166], [167], [168].

Understanding and regulating the relationship between data, variation, and generation process, is essential for the development of efficient GenAI systems [169]. It entails dealing

with issues; including dataset biases [170], [171], mode collapse [172], and balancing exploration and exploitation [?]. GenAI systems may generate high-quality, diversified, and realistic samples that correspond with the desired aims and applications by refining the training data, optimizing the generation procedures, and regulating variation [173].³

4) Performance Metrics: Evaluating the quality and diversity of generated samples is critical for assessing the models’ performance [174], [175]. Several strategies for assessing the quality, diversity, and authenticity of generated samples have been established. Here are some common evaluation techniques:

- Visual inspection [176], which is a subjective a evaluation method where human experts or users examine the generated samples and provide qualitative feedback.
- Inception Score (IS) [177] a widely used quantitative evaluation metric for, which measures the quality of generated samples based on their visual appeal and the diversity of the generated classes. Higher IS scores indicate better quality and diversity of the generated samples.
- Frechet Inception Distance (FID) [178] compares the distributions of real and generated samples by calculating the Fréchet distance between their feature representations extracted from a pre-trained Inception model.

Other performance metrics for GenAI systems, include PR-curves [179], coverage metrics [180], user studies and others [181].

B. Generative AI Design Cycle

A typical design cycle of Generative AI is shown in Figure. 3 as adopted from [129], [182]. GenAI development cycle may be broken into four key steps; (i) define the problem, (ii) model selection or developing from scratch (iii) adapt and align the model, or fine-tuning if necessary, and (iv) finally the deployment and optimization stage [183].

The first stage, scope entails deciding the the nature of the target GenAI model. For e.g., is the target to make it perform well at multiple tasks or only a single task? The nexts stage is the selection of a model. In this stage, a GenAI model developer needs to decide whether to use an existing model for the application or to pre-train one from scratch [184]. In this step, the developer can go for general techniques, for e.g. RNNs, transformers, or pretraining their own model and/or delve deep into creation methods involving more nuanced modifications of the model being used [185].

Following the second stage is an iterative phase of aligning and adapting the model for the scope chosen [186]. This includes steps of prompt engineering [79], which may consist of zero-shot learning, one-shot learning, or a few show learning techniques [187], or even fine-tuning based on the scope [188], and evaluating the model performance.

The last stage is the optimization and deployment for the target application [189]. Since prompting is a fundamental aspect of using and developing GenAI models, we provide a detailed discussion on prompting techniques in this section.

³genai meets copyright science paper

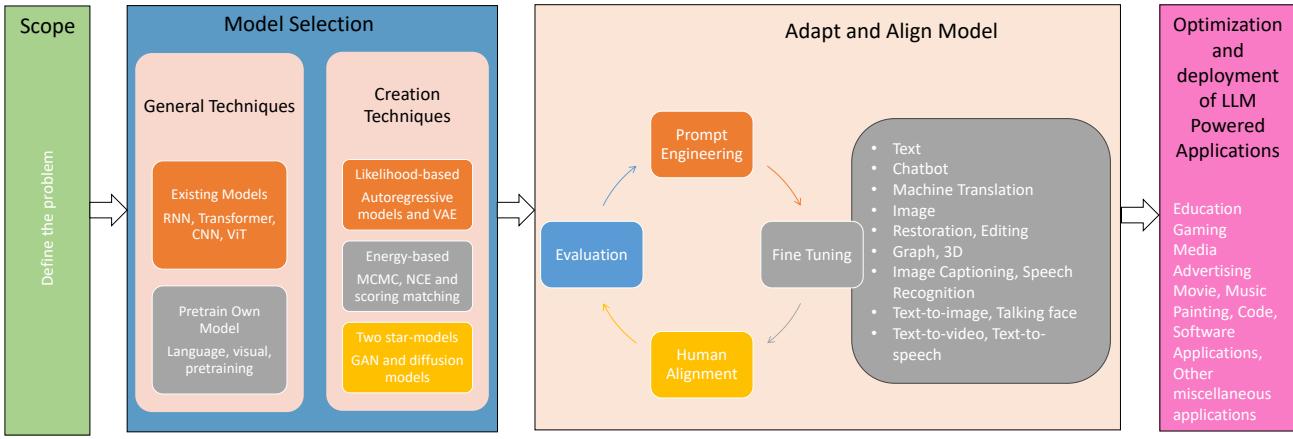


Fig. 3: Generative AI Design Process Schematic. Generative AI models can optimized for a variety of tasks, including, education, healthcare, entertainment and others.

1) Prompting: LLMs have given rise to what's called "Prompt Engineering". Prompts are the instructions provided to an LLM to make it follow specified rules, automation of processes and to ensure that the output generated is of a specific quality or quantity [79], [190]. While there is a lack of a formal definition, prompt engineering refers to the designing and wording of prompts given to LLMs so as to get a desired response from them. Writing a prompt appropriately is therefore very important if one needs to use LLMs to assist with tasks in the best manner possible [191].

While some formal techniques such as Explicit instruction (providing a clear direction to the LLM to do something) [192], System Specific Instruction (asking a question from the LLM to answer), Formatting with an example (providing a sample question and its answer and asking the LLM to provide an answer in the same manner), Control tokens (use special keywords in the prompt to help the LLM provide an answer while considering special provided criteria) [193] and Interaction and iteration/chaining (interact with model iteratively to reach to a good answer by fine-tuning on each reply) have been presented [79].

Several different frameworks have been suggested in lieu of prompt patterns for LLMs, these are generic prompt patterns targeting a specific category such as prompt improvement, input semantics etc [79], [194], [195], or prompting for software engineering tasks [196], [197], however, in this work, we aim to present some sets of commands to help users get the most out of the LLMs capabilities from a generic perspective.

- Defining the role/context:* This should be the first prompt for the LLM. An example of this prompt could be: "Act as a secretary to the Chair of the department", "Act as a Lawyer" or "Act as my programming tutor for Python". By defining a role for the LLM, one can direct it to provide replies or do tasks as a human would do when provided information to work on [198]. A similar first prompt could be providing the context. This can be performed to give the LLM a background of the

conditions in which the LLM is supposed to work. For e.g., "We are a company performing mobile application development for Fortune 500 organizations". This can then be followed up with aspects like actions, tasks to perform, steps to follow, etc as mentioned before [199].

- Prompt creation:* Another interesting prompt command is to ask the model to generate prompts for a certain task [200]. This way, the LLM can be used to generate optimized prompts for tasks that need to be done. An example of this could be: "You are a large language model and are an expert in generating prompts for ChatGPT. Please generate the best prompts on extracting important information from my time series data".
- Chain of thoughts:* Chain of thoughts prompting [201] in the context of Language Models (LMs) refers to the practice of providing a series of related prompts or partial sentences to guide the generation of coherent and connected text. Instead of providing a single prompt, a chain of thoughts prompt involves providing multiple prompts in succession to encourage the LM to continue generating text that follows a specific line of thinking or narrative [202].
- Other interesting directions in which Prompts can be given are explanation prompts [203] (e.g., "Explain the concept of infinity", Instructional Guides (e.g., "How do I tie my shoe laces"), Extract information (e.g.: one can paste a passage and ask the model to provide answers to questions that one might have), Solve Math problems (e.g., "Find the roots for the quadratic equation, $2x^2 + 3x + 10$ ") and Code help (e.g., "Find the syntax error in the following code") [204].

One concept within prompt engineering is in-context learning [205] in terms of the user "teaching" the LLM to act in a certain manner. The typical prompting scheme in which the LLM is asked to perform a task is an example of zero-shot inference [206], that is, within the context of the current task being worked on, the LLM is asked to perform the task

without providing any sample solution for it. An example of this type prompt could happen in the task of classifying tweets. To perform zero-shot inference, a user will have to just provide the text of the tweet to the LLM and ask it to classify it as positive or negative in sentiment. Another type of prompting could be one-shot inference [207]. In such a case, the user would give an example of a task solution o the LLM and then ask it to perform the task. In the tweet sentiment analysis example previously, this would be the user providing a sample of a tweet and information the LLM that the sentiment is positive and then providing it a second tweet to determine the sentiment of. The third type of prompt is few-shot inference [208], herein, the user provides a few examples of task solutions to teach the LLM about the kind of operation the user wants it to do. For the tweet sentiment analysis example above, it would be providing a tweet/tweets with a positive sentiment and indicating its sentiment and doing so with a negative tweet/tweets as well. Finally, the user can then use the LLM for tweet classification. Using in-context learning allows a user to “fine-tune” the LLM for the specific tasks being performed in the application [209].

2) *Negative Prompting*: Negative prompting [210], [211], [212] provides directions to the LLM about aspects of the prompt that it should avoid generating or deliberately excluding during the generation process [210]. Through the use of negative prompts, one can fine-tune the results generated by the LLM in response to a prompt while being able to keep the prompt generation generic [213]. Another advantage of the use of negative prompting is that it allows for moderation of the output content generated by the model thereby preventing harmful or inappropriate from being generated. “Don’t write anything that is offensive or harmful, or factually incorrect.” This prompt tells the model to avoid generating text that could be offensive or harmful to others and inaccurate. Notably, the authors in [214] conducted experiments for text based image tranlation and found that negative prompting to be very useful when working with textureless images. Moreover, this type of prompting is very useful when working on text to image generation scenarios and has been incorporated in text to image generation methods such as Muse [215].

3) *Visual Prompting*: Visual prompting [216] refers to the use of visual prompts (such as images or non-visual ones such as music) when providing directions to a model in addition to plain text prompts. The aim is to provide the AI model with a starting point or an example/reference that it can use for the given generative task. For images, this may be given to modify the image provided or generate something that is similar in style, color, or texture etc [217]. This can help in generating content that is closer to a user’s expectation from the generative AI being used.

An image-based example of visual prompting could be providing a picture of an office and asking the AI to generate a different theme for it, maybe more nature-centric or in a different color or organizational style [218]. Visual prompting provides greater control of the generated output and therefore results in a more accurate result. It should be noted that visual prompting is not related to images only, this is currently being explored for a host of different applications, including, text

generation (generating something based on a sample text so as to copy its style of writing for e.g.) [98], the composition of music (wherein the supplied music piece can be used as a reference for the type of music to compose) [219], game development [220] (where a defined game environment may be provided to the model as a starting point and the model is asked to generate new and unique content) and virtual and augmented reality (wherein a set of augmented/virtual reality environments can be provided to further populate/create current/new environments) [221].

III. OVERVIEW OF LLMS

Large Language Models have transformed the way we interact with and process language, opening up new possibilities for natural language understanding, generation, and communication [222]. Early language models primarily relied on rule-based methods, including Noam Chomsky’s theory of grammar [223]. They continue to evolve, pushing the boundaries of what is possible in the realm of language processing and artificial intelligence [224]. In this Section, we briefly discuss the history, evolution, and training of LLMs.

A. History of LLM

The history of LLMs can be traced back to the early days of NLP research [225], [226]. The first language models were developed in the 1950s and 1960s [227]. These models were rule-based [228] and relied on hand-crafted linguistic rules and features to process language [229]. They were limited in their capabilities and were not able to handle the complexity of NLP [230].

In the 1980s and 1990s, statistical language models were developed [31]. These models used probabilistic methods to estimate the likelihood of a sequence of words in a given context [231]. They were able to handle larger amounts of data and were more accurate than rule-based models [232]. However, they still had limitations in their ability to understand the semantics and context of language [233].

The next major breakthrough in language modeling came in the mid-2010s with the development of neural language models [234]. These models used deep learning techniques to learn the patterns and structures of language from large amounts of text [10]. The first neural language model was the recurrent neural network language model (RNNLM) [42], which was developed in 2010. RNNLM was able to model the context of words and produce more natural-sounding text than previous models [235]. In 2015, Google introduced the first large-scale neural language model called the Google Neural Machine Translation (GNMT) system [236].

The development of LLMs continued with the introduction of the Transformer model in 2017 [45]. The Transformer was able to learn the longer-term dependencies in language and allowed for parallel training [237] on multiple Graphical Processing Units (GPUs), making it possible to train much larger models [238]. The release of OpenAI’s GPT-1 [239] in 2018, marked a significant advance in NLP with its transformer-based architecture. With 117 million parameters, GPT-1 could generate contextually relevant sentences, demonstrating the

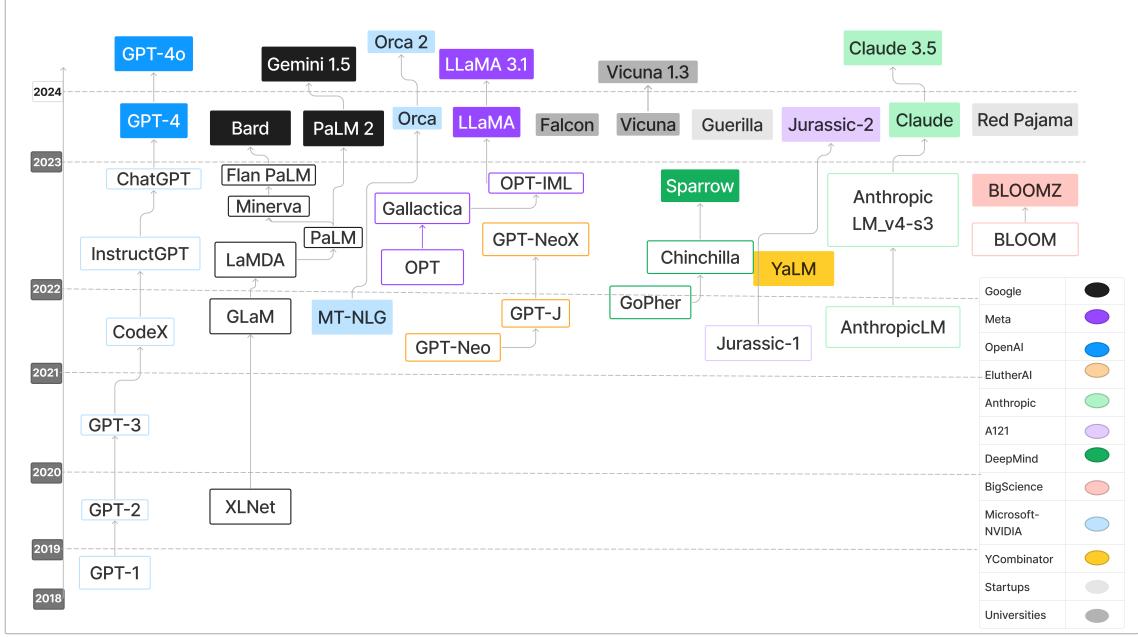


Fig. 4: Illustration of the evolution of Large Language Models (LLMs) over time, highlighting their development across a range of research and commercial organizations. Starting from the initial advancements made in this field, the figure maps out the journey of LLMs, outlining the key milestones, breakthroughs, and model iterations along the way. **need to be updated**

potential of transformers in revolutionizing NLP tasks [240]. While GPT-1 had its limitations, it set the stage for subsequent, more powerful models, propelling a new era of AI research and highly-competitive research in LLMs (see Fig. 6).

In 2020, OpenAI released GPT-3 [241], which was able to generate highly coherent and natural-sounding text [242]. GPT-3 demonstrated the potential of LLMs for a wide range of NLP tasks [243]. Inspired by the success of GPT-3, OpenAI released the next iteration of their language model, GPT-4 [244] with the ability to generate even more coherent and natural-sounding text. Following GPT-4's success, Meta also introduced Llama [15], a family of open-source foundation models. Google introduced Bard [245], Amazon introduced AI features in the Alexa [246] models, and Huawei introduced Pangu models [144], joining the AI race.

B. State-of-the-art LLMs

LLMs over the past few years reflects significant advancements in natural language processing (NLP) and machine learning techniques. This progression is characterized by increases in model size, improvements in multilingual capabilities, enhanced few-shot learning, and innovations in training methodologies.

1) Early Stages (2019-2020): The initial stages of LLM development saw the introduction of models like T5 (2019) and GPT-3 (2020). T5, developed by Google AI, emphasized the effectiveness of shared encoder-decoder parameters and fine-tuning adapter layers. GPT-3, with its massive 175 billion parameters, demonstrated the power of few-shot learning, establishing LLMs as capable meta-learners. Another notable

model from this period is mT5 (2020), which highlighted the performance parity between large multilingual and single-language models.

2) Expansion and Multilingual Capabilities (2021): The year 2021 marked significant advancements with the introduction of models like PanGu- α , CPM-2, ERNIE 3.0, Jurassic-1, HyperCLOVA, and Gopher. These models varied widely in size, from 10 billion to 280 billion parameters, and emphasized few-shot learning, modular architectures, and enhanced representation capabilities. Notably, mT5's performance in multilingual tasks set a precedent for subsequent models like ERNIE 3.0 and HyperCLOVA, which further advanced multilingual and task-specific tuning.

3) Scaling and Efficiency (2022): In 2022, models such as ERNIE 3.0 Titan, GPT-NeoX-20B, and GLaM showcased increased parameter sizes and innovations in efficiency. ERNIE 3.0 Titan incorporated adversarial loss for better performance, while GPT-NeoX-20B improved training speed with parallel layers. GLaM's use of mixture-of-experts (MoE) in transformer layers demonstrated how to maintain model capacity with less computation. This period also saw a focus on balancing model size and training data, as evidenced by DeepMind's Chinchilla, which stressed proportional scaling of model size and training tokens.

4) Specialization and Cost Efficiency (2023): The LLM landscape in 2023 featured a diverse array of models, including UL2, GPT-4, Claude Instant, and GPT-4 Turbo. These models emphasized mode switching, multimodal capabilities, and cost efficiency. GPT-4, with its 1.76 trillion parameters, showcased enhanced reliability and creativity across various modalities.

Models like Claude Instant and GPT-4 Turbo prioritized optimization for speed and cost, reflecting a trend towards more practical and scalable AI solutions. Furthermore, the introduction of models tailored for specific applications, such as StarCoder for coding tasks and Mistral Tiny for mobile devices, highlighted the growing specialization within the LLM domain.

5) *Innovations in Instruction Tuning and Safety (2024):* Looking ahead to 2024, models such as Gemini 1.5, WebGPT, and LLaMA-2-Chat illustrate ongoing innovations in instruction tuning and safety. WebGPT, for instance, focuses on generating factually accurate responses with references, while LLaMA-2-Chat incorporates reinforcement learning from human feedback (RLHF) to enhance response safety and reduce susceptibility to jailbreak attacks. The trend towards incorporating multilingual training for improved generalization, as seen with mT0 and BLOOMZ, continued to shape the development of future LLMs.

The evolution of LLMs from 2019 to 2024 is marked by exponential growth in model size, enhanced multilingual and few-shot capabilities, innovations in training efficiency, and a focus on practical, specialized applications. These advancements underscored the transformative impact of LLMs on NLP and their potential to drive future innovations in AI.

TABLE II: Summary of LLM models, their parameters, insights, and developers.

Release Year	LLM Model	No. of Parameters	Key Insights	Developer
2019	T5	Up to 11 billion	Shared encoder-decoder parameters perform equivalently to unshared ones, and fine-tuning adapter layers surpasses training only classification layers	Google AI
2020	GPT-3	175 billion	The superior few-shot performance of LLMs compared to zero-shot indicates their capability as meta-learners	OpenAI
2020	mT5	Up to 13 billion	Large multilingual models perform comparably to single-language models on downstream tasks, whereas smaller multilingual models underperform	Google AI
2021	PanGua- α	200 billion	Few shot capabilities of LLMs are good	Huawei
2021	CPM-2	11 billion	Prompt fine-tuning, though slower to converge, updates few parameters, matches full model performance, and enhances understanding of sentence relationships by providing context and aggregating information	Tsinghua Univ.
2021	ERNIE 3.0	10 billion	A modular LLM architecture with universal and task-specific representation modules, optimized during fine-tuning, efficiently leverages the power of pre-trained models	Baidu
2021	Jurassic-1	178 billion	LLM performance is tied to network size; enhancing runtime with parallel operations and using a larger, unrestricted vocabulary for tokenization improves text representation and few-shot learning	AI21 Labs
2021	HyperCLOVA	204 billion	Employing prompt-based tuning can enhance model performance, often exceeding that of state-of-the-art models when backward gradients of inputs are accessible	Naver
2021	Yuan 1.0	245 billion	The model architecture that performs well in pre-training and fine-tuning scenarios may exhibit differing behavior in zero-shot and few-shot learning contexts	Inspur
2021	Gopher	280 billion	Relative encodings enable the model to process sequences longer than those encountered during training	DeepMind
2022	ERNIE 3.0 Titan	260 billion	Incorporating an additional self-supervised adversarial loss to distinguish between real and generated text enhances model performance compared to ERNIE 3.0	Baidu
2022	GPT-NeoX-20B	20 billion	Parallel attention and feed-forward layers accelerate training by 15% while maintaining performance, and training on File outperforms GPT-3 in five-shot evaluations	EleutherAI
2022	OPT	175 billion	If loss diverges, restart training from an earlier checkpoint with a lower learning rate, as the model is prone to generating repetitive text and getting stuck in loops	Meta
				Continued on next page

TABLE II: Summary of LLM models, their parameters, insights, and developers (continued).

Release Year	LLM Model	No. of Parameters	Key Insights	Developer
2022	Galactica	Not available	Galactica's performance has consistently improved across validation sets, in-domain, and out-of-domain benchmarks, even with repeated corpus exposures, surpassing existing research on LLMs	Meta
2022	GPT-3.5 Turbo	Not available 1.2 trillion	More cost-effective and performance-tuned version of GPT-3.5 Using mixture-of-experts (MoE) in transformer layers maintains capacity with less computation, and filtered data training significantly boosts NLG and NLU performance, especially for NLG	OpenAI Google Research
2022	LaMDA	137 billion	The model can be fine-tuned to effectively utilize various external information resources and tools To enhance effectiveness and efficiency, a transformer model can utilize a shallower encoder and deeper decoder, while massively scaling upsampling and filtering and clustering samples into a compact set improves performance	Google AI DeepMind
2022	AlphaCode	Not available	The model size and the number of training tokens should be scaled proportionately; specifically, doubling the model size should be accompanied by a doubling of the number of training tokens	DeepMind
2022	Chinchilla	70 billion	English-centric models excel in English translations, generalized models match specialized ones, larger models memorize more, and performance improves with scale beyond 540B parameters	Amazon AI
2022	PaLM	540 billion	Incorporating the Causal Language Modeling (CLM) task can enhance the model's efficiency in in-context learning, and placing layer normalization at the beginning of each transformer layer improves training stability	Google AI
2022	AlexaTM	Not available	Training with a mixture of denoisers not only surpasses PaLM when extended for additional FLOPs but also enhances the model's infilling ability and diversity in open-ended text generation	Google AI
2022	U-PaLM	Not available	Mode switching training enhances performance on downstream tasks, while Chain-of-Thought (CoT) prompting surpasses standard prompting techniques for the UL2 model	Google AI
2023	UL2	Not available	High-performance model for complex tasks	Cohere
2023	Command GLM-130B	Not available 130 billion	Incorporating a minor proportion of multi-task instruction data in pre-training datasets enhances the overall performance of the model	Tsinghua Univ.
2023	GPT-4	1.76 trillion	Multimodal capabilities with enhanced reliability and creativity	OpenAI

Continued on next page

TABLE II: Summary of LLM models, their parameters, insights, and developers (continued).

Release Year	LLM Model	No. of Parameters	Key Insights	Developer
2023	Claude Instant	Not available	Optimized for speed and cost-efficiency	Anthropic
2023	CodeGen	16 billion	Employing multi-step prompting for code synthesis enhances the comprehension of user intent and improves the quality of code generation	Salesforce
2023	Command Light	Not available	Focused on lightweight applications with efficient performance	Cohere
2023	LLaMA	7B, 13B, 70B	Good performance of a smaller model can be achieved by provided increased training and computing time	Meta
2023	PangU-Σ	1.6 trillion	At a low cost, a sparse model can provide the benefits of a large model	Huawei
2023	BloombergGPT	50 billion	Harming model capabilities can be avoided by pre-training the model with general-purpose and task-specific data yielding improved performance	Bloomberg
2023	Claude 2.1	Not available	Improved reasoning and safety over previous versions	Anthropic
2023	Mistral Tiny	Not available	Optimized for mobile and embedded devices	Mistral
2023	Llama 2 (Meta)	13 billion	Open-source model with extensive fine-tuning capabilities	Meta
2023	Llama 2-13B	13 billion	Balanced model for a variety of tasks	Meta
2023	XuanYuan 2.0	200 billion	Catastrophic forgetting can be avoided by combining pre-training and fine-tuning	Baidu
2023	PPLX-70B	70 billion	Tailored for enterprise applications	People.ai
2023	CodeT5+	220 million	In encoder-decoder architectures, Causal LM is crucial for a model's generation capability	Salesforce
2023	Mistral Small	Not available	Balances performance and resource efficiency	Mistral
2023	Mistral 7B	7 billion	High efficiency and performance for general use	Mistral
2023	StarCoder	Not available	Model can follow instruction even without fine-tuning by helpful, honest, and harmless (HHH) prompting	BigCode Project
2023	Command-R Plus	Not available	Advanced version with enhanced retrieval capabilities	Cohere
2023	Mixtral 8x22B	8 models of 22B each	Ensemble model for robust performance	Mistral
2023	LLaMA-2	7B, 13B, 70B	Un-filtered training of model can yet a toxic but after fine-tuning it can perform better on downstream tasks	Meta
2023	GPT-4 Turbo	Not available	Improved cost and performance over GPT-4 with a 128K context window	OpenAI
2023	Mistral Medium	Not available	Enhanced performance for mid-tier applications	Mistral
2023	Inflection-2.5	2.5 billion	Lightweight model for efficient inference	Inflection AI
2023	Jamba	Not available	Designed for integrating with productivity tools	Notion AI
2023	PaLM-2	Not available	Larger models are outperformed by smaller models trained for larger iterations	Google AI
2023	Gemini Pro	Not available	Advanced model with enhanced reasoning capabilities	Google DeepMind
2023	Claude 3 Opus	Not available	Advanced conversational abilities with high safety standards	Anthropic

TABLE II: Summary of LLM models, their parameters, insights, and developers (continued).

Release Year	LLM Model	No. of Parameters	Key Insights	Developer
2023	Claude 3 Sonnet	Not available	Optimized for creative writing tasks	Anthropic
2023	Claude 3 Haiku	Not available	Specialized in generating concise and poetic text	Anthropic
2023	Stable LM 2	Not available	Focus on stability and reliability for production environments	Stability AI
2023	T0	11 billion	Baselines are outperformed by multi-task prompting enabling zero-shot generalization	BigScience
2024	Gemini 1.5	Not available	Mid-tier model with enhanced contextual understanding	Google DeepMind
2024	WebGPT	Not available	Labelers can easily deduce the factual accuracy of answers by generating responses with references	OpenAI
2024	Tk-INSTRUCT	1.5 billion	For unseen tasks instruction tuning leads to a stronger generalization	AI2
2024	mT0 and BLOOMZ	176 billion	Multi-lingual training (non-English and English) improves zero-shot generalization for both languages	BigScience
2024	Llama 3 (Meta)	Not available	Successor to Llama 2 with improved capabilities	Meta
2024	OPT-IML	175 billion	Using pretraining data in small amount during fine-tuning is effective. for generalization, training datasets should also be proportional	Meta
2024	Sparrow	Not available	Preference win rates and resilience against adversarial probing is improved by reinforcement learning (RL) with reranking yielding the optimal performance	DeepMind
2024	GPT-4o	Over one trillion	GPT-4o excels in real-time, multimodal processing, integrating text, audio, and images, with superior multilingual capabilities. It offers enhanced performance across various tasks, being 50% cheaper and twice as fast as GPT-4 Turbo	OpenAI
2024	Flan	Not available	CoT tuning was used to improve zero-shot reasoning. Multi-tasking improves zero-shot generalization abilities	Google AI
2024	WizardCoder	Not available	Improves performance with re-written instruction-tuning data into a complex for fine-tuning	Notion AI
2024	LlAMA-2-Chat	Not available	The model learns to generate safe responses through fine-tuning on safe examples, and an additional RLHF step further enhances its safety, making it less susceptible to jailbreak attacks.	Meta
2024	LIMA	Not available	Fine-tuning of LIMA can be performed with less high quality data	Meta

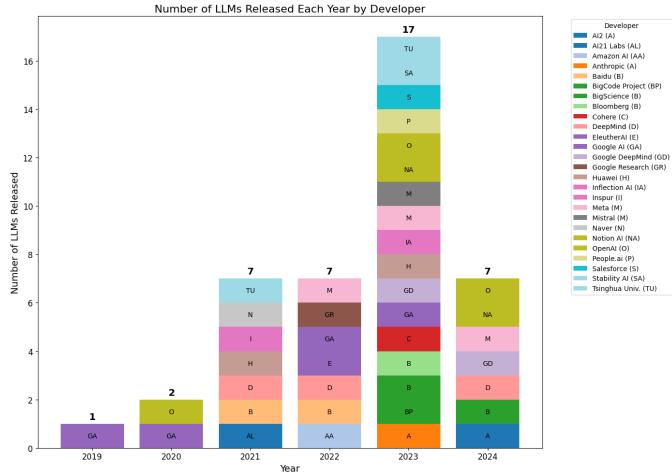


Fig. 5: Graph representing the summary of LLMs released during 2019 and 2024 from different developers, with yearly total. Please note that the LLMs shown here are summarized until July 2024.

C. Training of LLMs

Training large language models involves several key steps [248]. The process typically begins with the collection and preprocessing of a large amount of text data from diverse sources [249], such as books, articles, websites, and other textual corpora. The curated dataset [250] serves as the foundation for training the LLMs. After the removal of duplicates [251], noisy and poisonous data [252] and ensuring privacy reduction [253], the training process involves unsupervised learning, where the model learns to predict the next word in a sequence given the preceding context assuming the language generation as a random process [254].

Currently, LLMs utilize Transformers which enable them to model long-range dependencies [255], understand text data [256], enable them to generate new content in the style and characteristics of a genre or author [222]. The training objective is to optimize the model’s parameters to maximize the likelihood of generating the correct next word in a given context [85]. This optimization is typically achieved through stochastic gradient descent (SGD) [257] or its variants, combined with backpropagation [258], which computes gradients to update the model’s parameters iteratively. Some of the popular transformer-based LLMs are discussed below.

- Bidirectional Encoder Representations from Transformer (BERT): BERT [48] is a prominent language model with significantly advanced NLP tasks. Its training process comprises pretraining and fine-tuning stages [259]. During pretraining, BERT learns a general language representation from large-scale unlabeled text data. It employs masked language modeling (MLM) [260] and next-sentence prediction (NSP) tasks [261].

MLM involves masking a portion of input tokens and training the model to predict the original masked tokens, fostering bidirectional context understanding [262]. NSP trains BERT to predict whether a second sentence follows the first, enhancing coherence comprehension. After pre-

training, BERT undergoes fine-tuning on specific tasks with labeled data. Fine-tuning tailors BERT’s learned representations to target tasks, such as sentiment analysis or named entity recognition [263]. Training BERT demands significant computational resources [240], utilizing high-performance hardware like GPUs or Tensor Processing Units (TPUs) or field programmable gate arrays (FPGAs) [264], [265], [266]. Techniques such as layer normalization [267], residual connections [268], and attention mechanisms inherent in the transformer architecture further enhance BERT’s capacity to capture intricate dependencies and long-range contextual relationships.

- eXtreme Language understanding Network (XLNet): XLNet [269] is a generalized autoregressive [270] pre-training method that surpasses the limitations of traditional left-to-right or right-to-left language modeling. XLNet is trained using a permutation-based approach that differs from traditional autoregressive models [271]. The training of XLNet involves two key steps: unsupervised pretraining and supervised fine-tuning. During unsupervised pretraining, XLNet learns to predict words conditioned on the entire input context by maximizing the expected log-likelihood over all possible permutations. This is achieved using a variant of the transformer architecture, similar to models like BERT. The permutation-based objective function used in XLNet training presents unique challenges [272]. Unlike traditional autoregressive models that can rely on the causal order of words for prediction [273], XLNet needs to consider all possible permutations, resulting in an exponentially large number of training instances.
 - Text-to-Text Transfer Transformer (T5): T5, developed by Google, is a versatile language model that is trained in a "text-to-text" framework [274]. The key innovation of T5 is the formulation of all tasks as text generation problems. This means that every task, including text classification, summarization, translation, and question answering, is cast into a text-to-text format [275]. For example, instead of training T5 to answer questions directly, it is trained to generate the complete answer given the question and relevant context. In the pretraining phase, T5 is trained using a variant of the transformer architecture. The pretraining objective is typically based on maximum likelihood estimation, where T5 is trained to predict the target text given the source text. Once pretraining is complete, T5 undergoes fine-tuning on specific downstream tasks [274].
 - Conditional Transformer Language Model (CTRL): CTRL is a language model designed to generate text based on specific control codes or prompts [276]. One of the unique aspects of CTRL is its conditioning of control codes or prompts. These control codes guide the model's text generation process, allowing users to specify the desired style, topic, or other characteristics of the generated text. The fine-tuning phase of CTRL is crucial for adapting the model to specific tasks or domains [277].

TABLE III: Comparison of Traditional ML, Deep Learning, and Large Language Models

Aspect	Traditional ML	Deep Learning	Large Language Models
Training Data Size	Medium	Large	Very Large
Feature Engineering	Required	Partial	Not Required
Model Complexity	Low	High	Very High
Interpretability	High	Low	Low
Performance	Moderate	High	Very High
Hardware Requirements	Standard	High	Very High
Scalability	Limited	Moderate	High
Fine-tuning	Not Applicable	Possible	Possible
Transfer Learning	Limited	Effective	Effective
Contextual Understanding	Limited	Moderate	High

TABLE IV: State-of-the-art for LLM training pipeline [247]. Notations: RM: Reward Modeling, RL: Reinforcement Learning, SFT: Supervised Fine-tuned.

Stage	Pretraining	Supervised-Finetuning	Reward Modeling	Reinforcement Learning
Dataset	Raw Internet I	Demonstration	Comparisons	Prompts
Algorithm	Language Modeling	Language Modeling	Binary Classification	Reinforcement Learning
Model	Base Model	SFT Model	RM Model	RL Model
Resources	100s of GPUs months of training deployable	1-100 of GPUs days of training deployable	1-100 of GPUs days of training not deployable	1-100 of GPUs days of training deployable

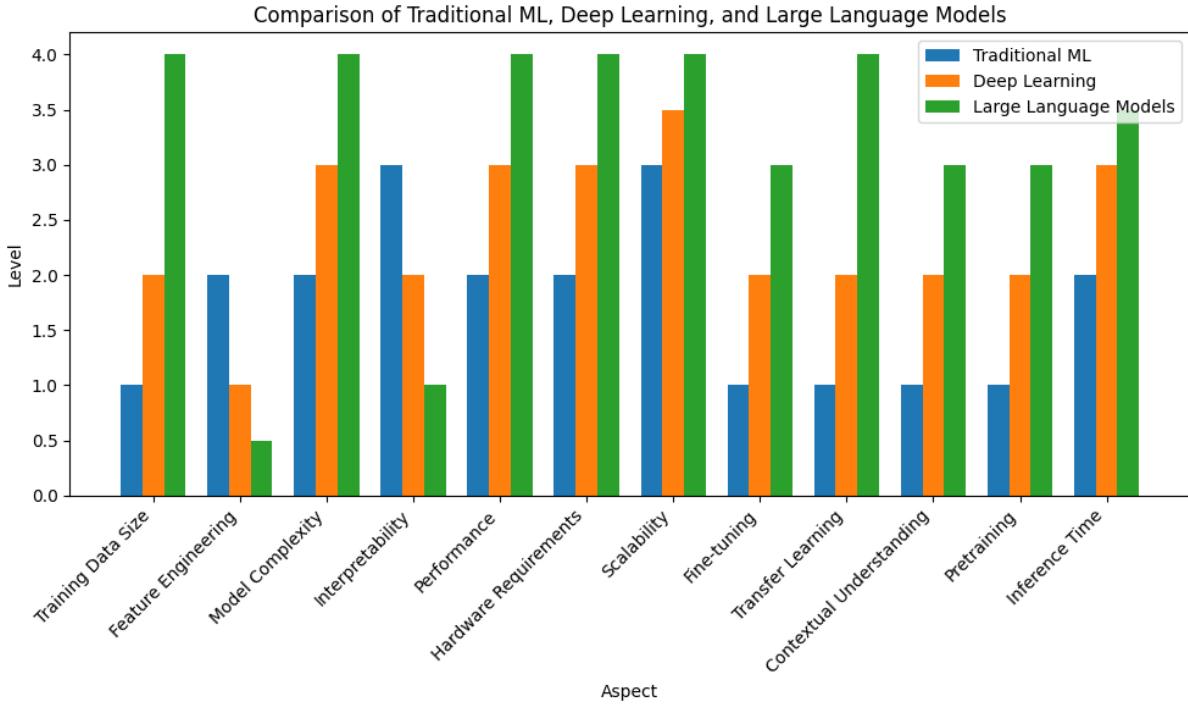


Fig. 6: Comparision of Traditional ML, Deep Learning, and Large Language Models

IV. LLM’s DATASETS

The landscape of large language model (LLM) development is significantly shaped by the diversity and quality of datasets used throughout various stages of model training, fine-tuning, preference assessment, and evaluation. This section provides a detailed examination of the datasets employed in the development of LLMs as presented in Figure 7, categorized into five key perspectives: Pre-training Corpora, Instruction Fine-tuning Datasets, Preference Datasets, Evaluation Datasets, and Traditional Natural Language Processing (NLP) Datasets. Each category plays a pivotal role in shaping the capabilities,

robustness, and generalization of LLMs across different tasks and domains.

A. Pre-training Corpora

Pre-training corpora form the foundational layer of LLM training, providing vast amounts of textual data that enable models to learn language structures, patterns, and general knowledge. These corpora are constructed using various methods, each with distinct advantages and limitations:

- 1) *General Pre-training*: General pre-training datasets consist of large-scale collections of text that cover a wide range of

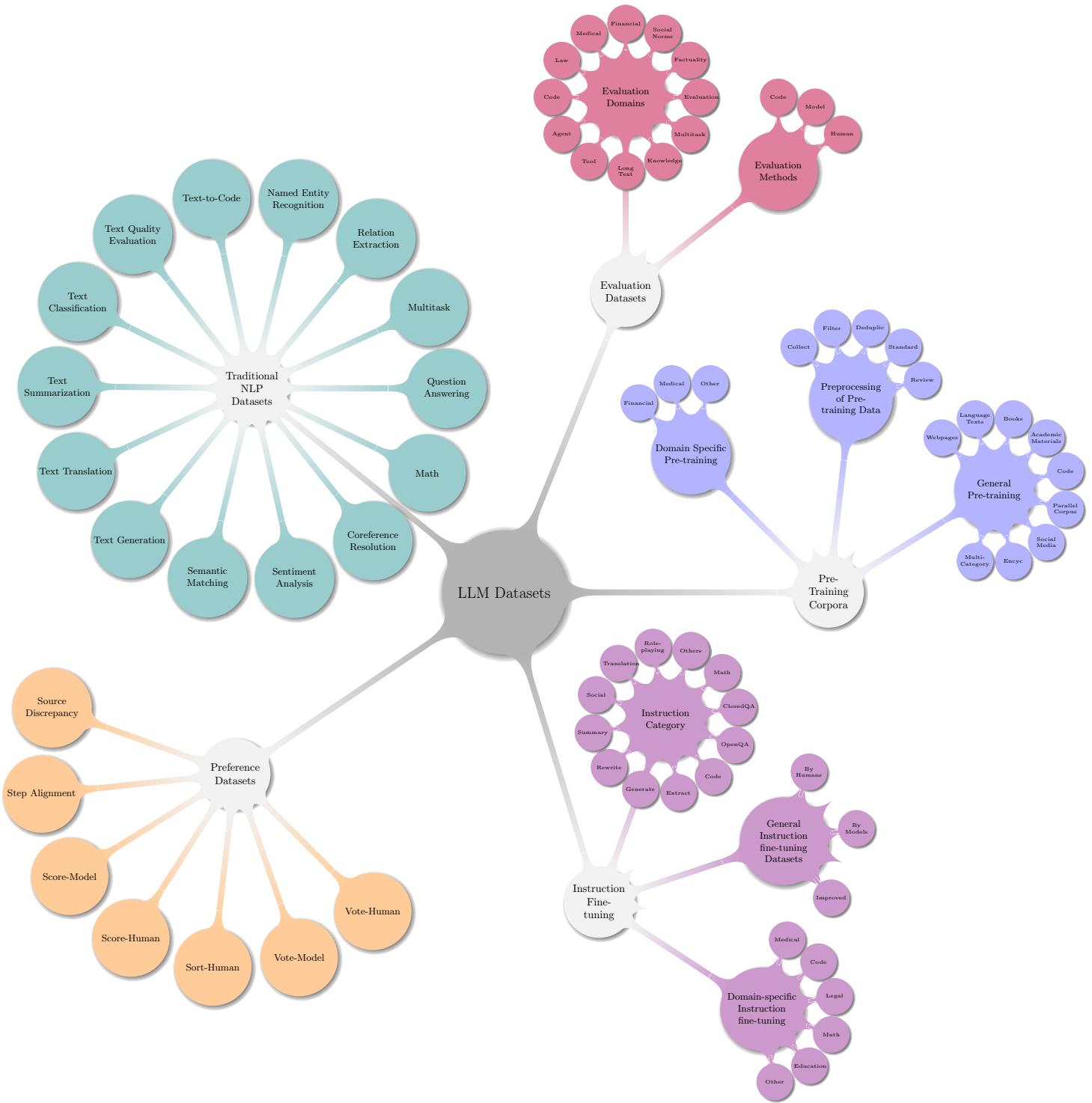


Fig. 7: Comprehensive Overview of LLM Datasets, spanning from pre-training corpora, instruction fine-tuning, traditional NLP datasets, to evaluation and preference datasets.

topics and languages. These datasets are essential for training LLMs to develop a broad understanding of language, enabling them to perform well on diverse tasks:

- Common Crawl [278]: A massive web archive dataset, widely used for language model training due to its extensive coverage of internet text, offering a diverse range of language structures and topics.
- CC100 [274]: A multilingual dataset derived from Common Crawl, containing text in 100 languages, facilitating the training of models on diverse linguistic data.
- Project Gutenberg [279]: A collection of over 60,000 public domain books, providing high-quality, classical literature for language models to learn structured, coherent narratives.
- C4 [274]: The Colossal Clean Crawled Corpus, derived from the web, specifically filtered to remove noise, making it a high-quality resource for training general-purpose language models.
- The Pile [280]: A large-scale, diverse dataset designed for training language models, composed of 22 smaller datasets, including academic papers, code, and internet text, offering comprehensive coverage across domains.
- Europarl [281]: A parallel corpus containing transcripts from the European Parliament, widely used in training machine translation models due to its multilingual content.
- Reddit and Twitter datasets [282]: Social media datasets providing real-world conversational data, which are crucial for training models in understanding informal language and conversational nuances.
- Wikipedia dumps [283]: A collection of articles from Wikipedia, offering encyclopedic knowledge across various subjects, serving as a foundational knowledge base for language models.
- OSCAR [284]: A large-scale multilingual corpus, sourced from Common Crawl, featuring clean and deduplicated text in multiple languages, supporting multilingual and cross-lingual model training.

2) *Domain-Specific Pre-training*: Domain-specific pre-training datasets are tailored to specific fields, such as finance or medicine. These datasets enhance a model's ability to understand and generate domain-specific language, making them crucial for specialized applications:

- BBT-FinCorpus [285]: A dataset specifically curated for financial texts, designed to enhance models' understanding of finance-related language and terminology.
- Medical-pt [286]: A specialized corpus focused on biomedical literature, aiding in the development of models proficient in medical and clinical terminology.
- PubMed [287]: A comprehensive resource of biomedical articles, extensively used for training models in the medical domain, particularly in tasks like named entity recognition and information retrieval.

3) *Preprocessing of Pre-training Data*: The preprocessing of pre-training data involves various techniques to enhance the quality and utility of the datasets. This includes collecting

large volumes of data, filtering out noise, removing duplicates, and ensuring consistency across datasets:

- High-volume data collection [280]: Involves the aggregation of vast amounts of text data from diverse sources to ensure comprehensive language coverage.
- Quality filtering processes [274]: Filtering techniques applied to remove noise and irrelevant content, thereby enhancing the overall quality of the dataset.
- Deduplication [280]: The process of identifying and removing duplicate text to avoid redundancy in the training data, ensuring more effective model training.
- Consistency across datasets [274]: Standardization efforts that ensure uniformity in data formatting and annotation across different datasets.
- Review [288]: A final step where human annotators or models review the dataset to identify and correct potential issues, ensuring the highest quality standards.

B. Instruction Fine-tuning Datasets

Instruction fine-tuning datasets are crucial for aligning LLMs with specific user instructions and enhancing their ability to follow complex commands. These datasets are constructed through various methods:

1) *General Instruction Fine-tuning Datasets*: General instruction fine-tuning datasets provide a broad range of instructional prompts, helping models improve their ability to understand and execute specific tasks. These datasets often cover diverse topics and are used to fine-tune models for better alignment with human instructions:

- Alpaca data [289]: A dataset created by human annotators providing a range of instructional tasks, aimed at fine-tuning models to better understand and execute specific instructions.
- Self-Instruct [290]: A dataset generated by models, leveraging self-instruction techniques to produce diverse instructional prompts, thereby enhancing scalability.
- BELLE Generated Chat [280]: A refined dataset derived from existing instructions, focusing on improving the quality and variety of instruction-following tasks.

2) *Instruction Category*: Instruction categories encompass datasets that focus on specific types of tasks, such as reasoning, text generation, or code comprehension. These datasets allow for targeted fine-tuning, enhancing models' performance in specific areas:

- Reasoning tasks datasets [291]: Datasets focused on evaluating and improving the reasoning capabilities of models, covering logical, numerical, and causal reasoning.
- Mathematical reasoning datasets [292]: Datasets that challenge models with mathematical problems, ranging from basic arithmetic to complex word problems.
- Cloze tasks like CLOTH [293]: Designed to test a model's ability to understand and complete text passages by filling in missing words, crucial for language comprehension.
- Open-ended question-answering datasets [294]: These datasets present open-ended questions that require models to generate appropriate responses based on given contexts.

- Code generation and understanding datasets [103]: Focused on enabling models to write and comprehend code, these datasets are vital for tasks in software development and automation.
- Information extraction tasks [295]: Datasets designed to train models in extracting specific information from text, such as names, dates, and entities.
- Text generation tasks [296]: These datasets provide structured input, requiring models to generate coherent text based on given keywords or topics.
- Paraphrasing and rewriting tasks [279]: Focused on teaching models to rephrase or rewrite sentences while preserving the original meaning, enhancing versatility in language use.
- Text summarization tasks [297]: Datasets that challenge models to condense longer texts into concise summaries, retaining key information and context.
- Social task datasets [298]: Designed to help models understand and engage in social interactions, including etiquette and empathy-driven responses.
- Machine translation datasets [299]: Essential for training models to accurately translate text between different languages while preserving meaning and nuance.
- Role-play and dialogue datasets [300]: Datasets where models engage in role-playing scenarios, useful for developing conversational agents and interactive systems.
- Miscellaneous instruction tasks [288]: A collection of various instructional prompts that do not fit into specific categories but are important for broadening a model's instruction-following capabilities.

3) *Domain-specific Instruction Fine-tuning:* Domain-specific instruction fine-tuning datasets are tailored to enhance a model's performance in specific fields such as medicine, law, or education. These datasets provide specialized instructions that align with the language and tasks relevant to particular domains:

- Medical instruction datasets [301]: Datasets focused on medical and clinical instructions, enhancing the model's ability to process and generate medical information.
- Code-specific instruction datasets [103]: Designed for fine-tuning models to understand and follow coding-related instructions, improving their performance in programming tasks.
- Legal domain instruction datasets [302]: These datasets provide legal scenarios and instructions, useful for developing models in the legal domain.
- Specialized math datasets [292]: Focused on enhancing a model's mathematical problem-solving skills, particularly in educational and technical contexts.
- Educational tasks datasets [303]: Instructional datasets aimed at training models to assist with various educational tasks, from answering questions to generating teaching material.
- Miscellaneous domain-specific tasks [288]: Cover a broad range of domains, ensuring that models can be fine-tuned for specialized industry or academic needs.

C. Preference Datasets

Preference datasets play a critical role in refining LLM outputs according to human or model preferences, enhancing the overall quality and alignment of generated content. These datasets help ensure that the outputs of LLMs are aligned with desired outcomes, based on specific preferences or criteria:

- Human-based voting datasets [303]: Datasets where human annotators vote on the best response, helping align model outputs with human preferences.
- Model-based voting datasets [148]: In these datasets, models vote on their own or other models' outputs, a scalable method for preference fine-tuning.
- Human-based sorting datasets [304]: Datasets where human annotators rank responses based on criteria such as relevance, clarity, or style.
- Human-based scoring datasets [148]: Involves human annotators assigning scores to model outputs, providing fine-grained feedback for further tuning.
- Model-based scoring datasets [305]: Similar to human scoring, but performed by models to evaluate the quality of generated content.
- Step-wise task alignment datasets [148]: Datasets that evaluate models on the alignment of outputs with step-by-step task instructions, ensuring process adherence.
- Discrepancy evaluation datasets [305]: Used to identify and correct discrepancies between model-generated outputs and expected results, enhancing reliability.

D. Evaluation Datasets

Evaluation datasets are integral to assessing the performance of LLMs across various tasks, including code generation, question answering, and reasoning. These datasets provide the benchmarks needed to evaluate how well models perform on specific tasks, ensuring they meet the desired standards.

1) *Evaluation Domains:* Evaluation domains categorize datasets based on the specific tasks they assess, such as general language understanding, reasoning, or domain-specific knowledge. These datasets provide comprehensive coverage of the different areas where LLMs are expected to perform well:

- General evaluation datasets [291]: Datasets that assess a model's overall capabilities across a wide range of tasks, providing a comprehensive evaluation.
- Academic exam datasets [306]: These datasets simulate academic exams, challenging models with questions typically encountered in educational settings.
- Subject-specific datasets [291]: Focused on particular subjects like history, science, or literature, assessing a model's depth of knowledge in specific areas.
- Natural language understanding datasets [307]: Essential for evaluating a model's ability to comprehend and process natural language across various contexts.
- Logical and contextual reasoning datasets [308]: Designed to test a model's reasoning abilities, particularly in understanding complex contexts and making logical inferences.
- Knowledge-based evaluation datasets [309]: These datasets assess a model's general knowledge and its

- ability to apply this knowledge to answer questions accurately.
- Long-form text analysis datasets [303]: Focused on evaluating a model's ability to generate, comprehend, and analyze long-form text, crucial for tasks like essays and reports.
 - Tool-use evaluation datasets [310]: Designed to evaluate models that interact with tools or APIs, testing their ability to generate and execute correct commands.
 - Agent-based interaction datasets [242]: Assess models' capabilities in simulated environments where they must interact with agents or follow complex workflows.
 - Code synthesis and execution datasets [103]: These datasets challenge models to generate executable code from natural language descriptions, crucial for AI-driven software development.
 - Legal knowledge evaluation datasets [311]: Focused on assessing a model's understanding of legal texts and its ability to apply legal reasoning.
 - Medical knowledge evaluation datasets [301]: Used to evaluate models on tasks requiring medical expertise, including diagnosis and treatment recommendations.
 - Financial domain evaluation datasets [285]: Assess models on their understanding of financial texts, important for applications in finance and economics.
 - Social norms evaluation datasets [312]: Datasets that assess how well models understand and adhere to social norms in their outputs, ensuring ethical and responsible AI behavior.
 - Factuality check datasets [313]: Used to evaluate the accuracy and factual correctness of model-generated content, essential for maintaining trust in AI outputs.
 - Multitask evaluation datasets [314]: Comprehensive datasets that assess a model's performance across multiple tasks, providing insights into its generalization capabilities.

2) *Evaluation Methods:* Evaluation methods refer to the techniques used to assess model performance, ranging from human evaluations to automated scoring by models. These methods ensure that models are evaluated accurately and consistently:

- Code evaluation methods [103]: Methods that focus on evaluating the correctness and efficiency of code generated by models.
- Model-based evaluation methods [148]: Techniques where models themselves are used to evaluate other models' outputs, providing scalability in the evaluation process.
- Human-based evaluation methods [305]: Involves human annotators assessing the quality of model outputs, ensuring that models meet human standards of quality and relevance.

E. Traditional Natural Language Processing (NLP) Datasets

Traditional NLP datasets provide a rich resource for evaluating specific language tasks, including sentiment analysis, text classification, and machine translation. These datasets are

foundational to LLM development, offering structured and annotated data that supports model training and testing across core NLP tasks. Notable datasets include:

- SQuAD [295] and TriviaQA [315]: Benchmark datasets for question answering tasks, requiring models to generate answers from provided text passages.
- GSM8K [292]: A dataset focused on mathematical problem-solving, challenging models to tackle complex word problems with logical reasoning.
- CoNLL2012 [316]: A dataset used for coreference resolution, where models identify which words in a text refer to the same entity.
- IMDB [317] and SST-2 [318]: Sentiment analysis datasets used to determine the emotional tone of a given text, whether positive, negative, or neutral.
- MRPC [319] and STS-B [320]: Datasets for semantic matching, where models assess the similarity or relatedness of sentence pairs.
- CommonGen [296]: A text generation dataset that provides structured inputs, challenging models to generate coherent sentences based on given keywords.
- WMT [299]: A dataset for machine translation, widely used for training and evaluating models in translating text between languages.
- XSum [321] and CNN/DM [297]: Text summarization datasets that require models to generate concise summaries from longer text documents.
- AGNews [322] and TREC [323]: Text classification datasets that test a model's ability to categorize text into predefined categories based on content.
- CoLA [324]: A dataset focused on evaluating the grammatical correctness of sentences, essential for text quality evaluation.
- MBPP [325]: A text-to-code dataset that provides natural language programming problems and assesses a model's ability to generate correct code solutions.
- CoNLL2003 [326]: A named entity recognition dataset where models identify and classify proper nouns, such as names of people, organizations, and locations, within a text.
- TACRED [327]: A relation extraction dataset where models identify relationships between entities in text, such as "located in" or "works for."
- GLUE [328] and SuperGLUE [307]: Multitask benchmark datasets that evaluate models across a range of NLP tasks, providing a comprehensive assessment of their generalization abilities.

V. TAXONOMY OF LARGE LANGUAGE MODEL TASKS

LLMs have a wide array of uses for the tasks of processing natural language including but not limited to writing, summarization, translation, retrieving information as shown in Fig. 8. In this section, the various tasks of LLMs towards working with developing such systems have been discussed

A. Question-answering

Question-answering (QA) systems [329] allow users to obtain direct answers to questions posed in natural language.

LLMs have become a key component in building robust QA systems [330]. LLMs can be effectively pretrained on large text corpora and then fine-tuned on QA labeled datasets [331]. This adapts them to extract or generate answers from passages of text. The broad linguistic knowledge learned during pre-training allows LLMs to understand the semantics of questions and use that to reason about potential answers. Fine-tuning on QA data teaches the models to identify relevant context passages and output the correct response [329]. Key benefits of using LLMs include handling complex questions, synthesizing answers from multiple context documents, and generating clarifying responses when a query is ambiguous. LLM-based QA systems have achieved high accuracy on benchmark datasets [332], surpassing previous state-of-the-art methods. They can be deployed via voice assistants [333], search engines [334], and other interfaces to provide users with quick access to information through natural dialog. Ongoing research is focused on improving reasoning abilities, explainability, and efficiency of LLM question answering.

B. Text Generation

Text generation is a useful application of large language models, which can automate the process of generating content for various purposes [335], such as articles [336], blogs [337], research papers, social media posts, product descriptions, source codes, emails, and more. With their ability to comprehend and generate natural language, these models can produce high-quality content that is both accurate and coherent [338].

C. Language Translation

LLMs possess the capability to translate text from one language to another with exceptional accuracy and fluency [339]. This feature is beneficial for a range of users, including language service providers, global companies, and individuals, who can utilize these models for real-time translation, localization, and overcoming language barriers in communication [340]. The impressive accuracy and fluency of these models make them a valuable tool for facilitating effective communication across different languages and cultures [341]. This feature has the potential to enhance global collaboration and increase access to information, making it an important area of research and development in the field of NLP [342].

D. Text Classification

In addition to their text generation and translation abilities, LLMs are also equipped with exceptional organizational capabilities, such as text classification, analysis, and categorization based on predefined labels or topics [343], [344]. This feature enables the models to effectively manage large volumes of textual data, making them highly valuable for a range of tasks such as sentiment analysis [345], spam detection [346], content moderation [347], and customer feedback analysis [348]. By automating these processes, language models can streamline data management [349], reduce manual labor, and improve the accuracy and efficiency of analysis [350]. These capabilities are particularly useful for businesses and organizations that deal with large amounts of textual data and require effective methods for organizing and analyzing it.

E. Summarization

LLMs possess the ability to generate concise and coherent summaries of lengthy texts or documents [351]. This feature is highly advantageous for a variety of uses, such as summarizing news articles, research papers, legal documents, and other types of content where extracting essential information is crucial. Summarization by language models can save time and effort while ensuring that the most important points are captured accurately [352]. This feature has the potential to enhance the efficiency and effectiveness of content absorption and subsequent creation, making it a valuable tool for individuals and organizations.

F. Virtual Assistance

In the realm of virtual assistants and chatbots [353], LLMs play a critical role. These models possess the ability to comprehend user queries, provide relevant information, and engage in natural language conversations. By leveraging large language models, virtual assistants and chatbots can provide highly effective and responsive support to users while also reducing the workload for human operators [354]. This area of research and development is of significant importance, as it has the potential to transform the way users interact with technology and improve the effectiveness and efficiency of customer support and service delivery.

G. Information Extraction (IE)

The use of LLMs in IE is significant for populating knowledge bases. By leveraging fine-tuned LLMs, entities such as people, organizations, and locations, as well as the relationships between them, can be accurately extracted from unstructured text. This process can facilitate the creation of structured knowledge graphs that can be utilized for various purposes [355]. In addition, LLMs assist in event extraction, enabling the identification of key occurrences described in text documents [356]. This feature has the potential to enhance the efficiency and accuracy of information extraction, making it a valuable tool for businesses and organizations that deal with large amounts of textual data. Further research and development in this area can lead to improvements in the quality and effectiveness of IE using LLMs.

H. Dialog systems

In the context of dialog systems, large language models play a crucial role in facilitating language understanding. The development of large pretrained models like Google's Meena and Microsoft's Blender has led to significant improvements in the naturalness and coherence of open-domain chatbots [357]. These models possess the ability to generate informative, interesting, and harmless responses, making conversational agents much more usable. The application of LLMs in dialog systems has the potential to transform the way users interact with technology, creating more engaging and effective conversational experiences [358]. Further research in this area can lead to improvements in the quality and effectiveness of dialog systems, making them even more valuable for a range of applications and industries.

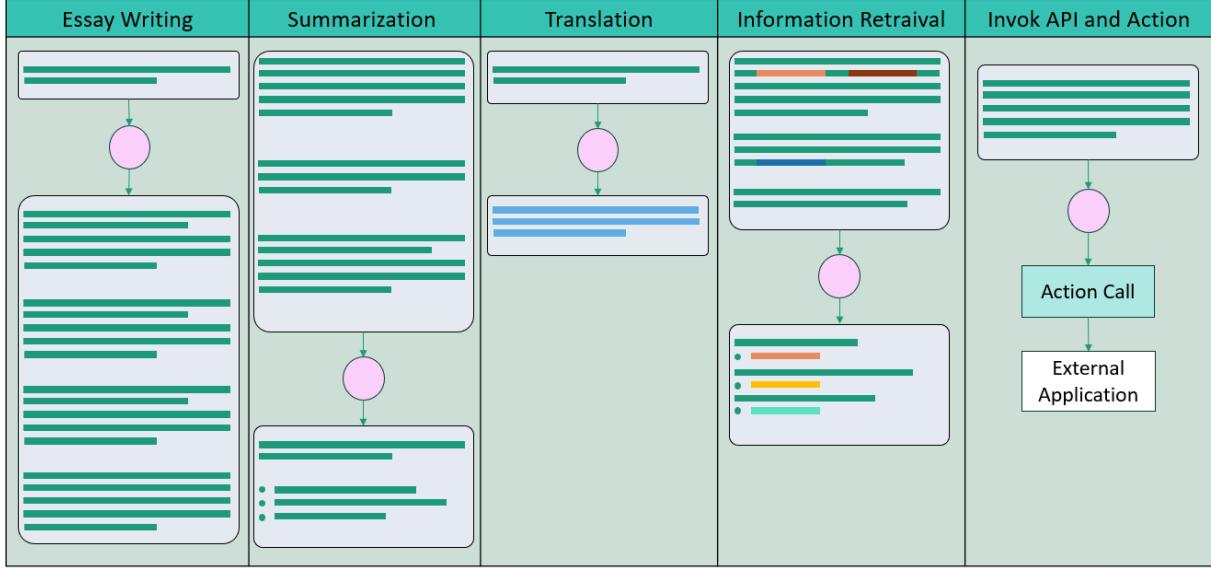


Fig. 8: Examples of LLM. Prompts can be to write an essay on X topic, summarize a paragraph, translate into X language etc.

I. Semantic Search

In the field of Semantic Search, query understanding is of utmost importance, and LLMs are unparalleled in their ability to discern the underlying intent and meaning of user search queries [359]. This ability enables next-generation search capabilities that go beyond simple keyword matching. For instance, LLMs can recognize that the phrases "best budget laptop"; and "affordable student computer" convey the same information need. This feature has the potential to enhance the accuracy and relevance of search results, making it easier for users to find the information they need. Further research and development in the area of Semantic Search can lead to the creation of more effective and efficient search systems, making LLMs a valuable tool for a range of applications and industries [57].

J. Speech recognition

Automated speech recognition is a crucial aspect of voice interfaces and transcription. While traditional systems relied on hidden Markov models or Gaussian mixture models, the emergence of deep learning has seen large neural network models like LLMs take center stage in state-of-the-art results. Fine-tuning these models on transcribed speech data using connectionist temporal classification loss enables them to learn acoustic-to-text mappings. This leads to significant improvements in the accuracy of automated speech transcription, even in the presence of accented speech or domain specific vocabulary [360]. The contextual knowledge and continual learning abilities of LLMs make them ideally suited for handling the variability and ambiguity inherent in speech signals. As LLMs continue to increase in scale, they are becoming the standard for building high performance and robust automated speech recognition systems [361].

To sum up the discussion on different LLMs, Table V provides information on the performance of various LLMs on different reasoning tasks.

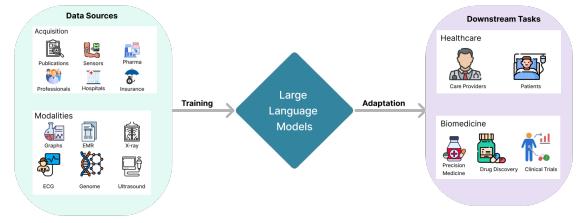


Fig. 9: Interactive framework: LLMs leveraging multimodal healthcare data for diverse tasks in healthcare and biomedicine.

VI. APPLICATIONS OF LARGE LANGUAGE MODELS

Given LLMs wide range of applications, we provide a discussion of their usecases in the fields based on their significance, relevance, and potential impact within their respective domains [266], [370]. An interactive framework for the integration of LLMs in the healthcare ecosystem is shown in Figure 9.

A. Medical

LLMs like ChatGPT have exhibited remarkable potential in diverse healthcare applications [371], [372], [373]. In medical education, ChatGPT has emerged as an interactive tool that aids learning and problem-solving [374], [375]. Notably, ChatGPT's performance in the United States Medical Licensing Exam (USMLE) was comparable to or exceeded the passing threshold, without requiring specialized training or reinforcement [374]. Moreover, ChatGPT's explanations displayed a high level of concordance and insightful understanding [372].

In [332], introduced MultiMedQA, a new benchmark dataset for evaluating LLMs on clinical tasks. MultiMedQA combines six existing medical question-answering datasets spanning

TABLE V: Comparison of LLMs’ Reasoning Performance. Notations: MMLU [362]: high school and college knowledge, GSM8K: elementary school math, MATH: very hard math and natural science. All current models struggle, BBH [363]: a collection of 27 hard reasoning problems, HumanEval [364]: a classical dataset for evaluating coding capability, C-Eval [365]: a collection of 52 disciplines of knowledge test in Chinese, TheoremQA [366]: a question-answering dataset driven by STEM theorems. [367], [362], [15], [368], [369], [365]

Model	Param.	Type	GSM8K	MATH	MMLU	BBH	HumanEval	C-Eval	TheoremQA
GPT-4	-	RLHF	92.0	42.5	86.4	-	67.0	68.7*	43.4
claude-v1.3	-	RLHF	81.8*	-	74.8*	67.3*	-	54.2*	24.9
PaLM-2	-	Base	80.7	34.3	78.3	78.1	-	-	31.8
GPT-3.5-turbo	-	RLHF	74.9*	-	67.3*	70.1*	48.1	54.4*	30.2
claude-instant	-	RLHF	70.8*	-	-	66.9*	-	45.9*	23.6
text-davinci-003	-	RLHF	-	-	64.6	70.7	-	-	22.8
code-davinci-002	-	Base	66.6	19.1	64.5	73.7	47.0	-	-
text-davinci-002	-	SIFT	55.4	-	60.0	67.2	-	-	16.6
Minerva	540B	SIFT	58.8	33.6	-	-	-	-	-
Flan-PaLM	540B	SIFT	-	-	70.9	66.3	-	-	-
Flan-U-PaLM	540B	SIFT	-	-	69.8	64.9	-	-	-
PaLM	540B	Base	56.9	8.8	62.9	62.0	26.2	-	-
LLaMA	65B	Base	50.9	10.6	63.4	-	23.7	38.8*	-
PaLM	64B	Base	52.4	4.4	49.0	42.3	-	-	-
LLaMA	33B	Base	35.6	7.1	57.8	-	21.7	-	-
InstructCodeT5+	16B	SIFT	-	-	-	-	35.0	-	11.6
StarCoder	15B	Base	8.4	15.1	33.9	-	33.6	-	12.2
Vicuna	13B	SIFT	-	-	-	-	-	-	12.9
LLaMA	13B	Base	17.8	3.9	46.9	-	15.8	-	-
Flan-T5	11B	SIFT	16.1*	-	48.6	41.4	-	-	-
Alpaca	7B	SIFT	-	-	-	-	-	-	13.5
LLaMA	7B	Base	11.0	2.9	35.1	-	10.5	-	-
Flan-T5	3B	SIFT	13.5*	-	45.5	35.2	-	-	-

professional medicine, research, and consumer queries. They introduced the concept of instruction prompt tuning [376], which can be used to improve the performance of LLMs on a variety of clinical tasks.

In [377], the potential of ChatGPT in radiologic decision-making is emphasized, showcasing its feasibility and potential benefits in enhancing clinical workflow and promoting responsible utilization of radiology services. Similarly, Kung et al. [372] concluded in their research that LLMs, including ChatGPT, have the capacity to enhance the delivery of individualized, compassionate, and scalable healthcare. These models can assist in medical education and potentially aid in clinical decision-making.

In the domain of clinical genetics, Duong and Solomon [378] found that ChatGPT’s performance did not significantly differ from humans when answering genetics-related questions. Furthermore, Fijacko [379] evaluated ChatGPT’s accuracy in answering questions related to life support and resuscitation. The findings revealed that ChatGPT demonstrated the ability to provide accurate answers to a majority of the questions on the American Heart Association’s Basic Life Support and Advanced Cardiovascular Life Support exams. In neurosurgical research and patient care, LLMs have been investigated for their potential role in various aspects, including gathering patient data, administering surveys or questionnaires, and providing information about care and treatment [380]. Furthermore, AI-powered chatbots hold the potential to enhance patient outcomes by facilitating communication between patients and healthcare professionals. Leveraging NLP, these chatbots can provide patients with information about their care and treatment in a more accessible manner [381]. There are

several tools already in use that allow the system to interact with patients such as Ada Health [382], Babylon Health [383], and Buoy Health [384]. There are tools developed to assist medical practitioners. One such tool is XrayGPT [385], it can be used for automated analysis of X-ray images and have the user/patient ask questions about the analysis. Through the chats, the user can get insight into their condition through an interactive chat dialogue. Another big development is the segment anything (SAM) model by meta, which may be fine-tuned for a variety of medical image tasks [386]. In the drug discovery domain, DrugGPT [387] is developed, which can be used to design potential ligands, targeting specific proteins, using text prompts. The authors in [388] propose generalist medical AI (GMAI). GMAI models are trained on large, diverse datasets of medical data, and they are able to perform a wide range of tasks, such as diagnosis, prognosis, and treatment planning.

B. Education

AI’s impact on education, a widely discussed topic, shows promise in areas like providing student feedback, assisting teacher development, and creating personalized learning experiences [389], [390], [391]. However, effectively implementing tech solutions for inclusive quality education poses significant challenges. General-purpose Large models applicable across tasks and subjects, like foundation models, offer a potential solution [392]. Already, foundation models like MathBERT improve performance in specific educational tasks, such as “knowledge tracing” and the “feedback challenge”

Since the advent of ChatGPT developed by OpenAI, the way students interact with educational materials, assignments and coursework has become different [393] [394] [395]. The

accuracy rate for the exams discussed in [393] was below 70 percent indicating its inability to pass the AHA exams. However, this conclusion was drawn due to a design limitation in their study, where they only generated a single response using ChatGPT, introducing bias and severely underestimating ChatGPT's capabilities in this domain. However, the latest study revealed that ChatGPT's accuracy rate increased to 96 and 92.1 percent for the Basic Life Support (BLS) and Advanced Cardiovascular Life Support (ACLS) exams, respectively, allowing ChatGPT to pass both exams with outstanding results [396].

One of the main advantages of using LLM such as ChatGPT in education is that they can help students complete their assignments more efficiently and provide personalized learning experiences [397]. Additionally, AI bots can help automate the grading process, which can reduce the workload for teachers and enable them to provide more detailed feedback to students [398]. Khan Academy, a nonprofit educational organization, has shown interest in utilizing ChatGPT for its business. They have developed an AI chatbot called Khanmigo, which serves as a virtual tutor and classroom assistant [399]. The goal of incorporating ChatGPT into their platform is to enhance tutoring and coaching experiences by providing one-on-one interactions with students. Khan Academy⁴ has been an early adapter of GPT-4 based LLMs working as online tutors, becoming the largest case study for the evaluations of LLMs in an educational context in the process. The incorporation of AI in tutoring and teaching proves that it can be a valuable tool in reducing negativity, particularly the perception that its main purpose is for cheating. Undoubtedly, AI technology is still in its nascent phase, yet it shows great potential in supporting students and catering to their individual requirements. [400].

However, there are also some potential drawbacks to using ChatGPT and AI bots in education. One concern is that these technologies may lead to a loss of creativity and critical thinking skills among students [401]. If students rely too heavily on AI bots to complete their assignments and exams, they may not be developing the skills necessary to think critically and solve problems on their own [397].

1) Learning in the age of AI: AI chatbots can serve as a valuable tool to aid in various aspects of syllabus preparation [402]. Course objectives can be generated, relevant topics identified, curricula structured, learning resources gathered and reviewed, assessment methods defined, engaging learning activities established, and a well-balanced course schedule created.

2) Major issues for AI in Education: One of the major concerns is the utilization of these tools without proper training. It is crucial to address the issue of inadequate training and contextual fine-tuning for LLMs, as their potential utilization without such preparations raises significant concerns [403]. While it is true that LLMs possess the ability to provide answers to a wide range of questions and assist users in generating responses effortlessly, it is essential for students and scientists [404] to receive adequate training specific to

⁴<https://blog.khanacademy.org/harnessing-ai-so-that-all-students-benefit-a-nonprofit-approach-for-equal-access/>

their needs in order to fully harness the capabilities of LLMs. Neglecting the necessity for context-specific training and fine-tuning can render these tools less effective and limit their true potential.

Another concern is the use of AI bots that could lead to a decrease in the number of teaching jobs available, which could further widen the gap between those who have access to education and those who do not [405], [406].

C. Agriculture

Combining LLMs with proximal (ground-based) sensing technologies can harness advanced data processing and analysis capabilities to enhance agricultural decision-making [407] for crops, soil, and environmental conditions [408]. This integration facilitates precise monitoring and management of agricultural practices, optimizing resource use and improving crop yields[409].

LLM in agriculture is transforming various crop-related industries (e.g., precision agriculture, farming management, crop monitoring, pest control etc.). In an attempt to improve understanding of plant science-related topics, LLMs such as PLLaMa [410] presented an advanced open-source language model derived from LLaMa-2 with an extensive repository of over 1.5 million scholarly articles specifically focused on plant science. This substantial enhancement endows PLLaMa with a deep reservoir of knowledge and expertise in plant and agricultural sciences, significantly broadening its capabilities in agriculture science. For farm management practices, LLMs aided in designing robotic system for optimizing tomato picking [411].

In precision agriculture, Tan et. al [412] used LLM-based model to enhance crop yield prediction accuracy. LLM through prompt engineering is used to integrate and analyze diverse agricultural data sources. The approach offered significant improvements over traditional methods with 90% contextual accuracy. However, authors indicated complexity in designing effective prompts for LLMs and dependence of the method on the quality and comprehensiveness of the input data. Future directions included refining the prompt designs and incorporating more diverse data sources. Accurately predicting crop yields, which is essential for agricultural planning and food security

In agriculture safety, ChatGPT is used in demonstrating five specific use cases for LLMs in agricultural safety and health: responding to technical questions, interpreting research articles and technical information, summarizing incident reports and identifying trends, guiding the use of engineering and regulatory standards, and brainstorming content for presentations and articles. Each use case included practical examples of how LLMs can enhance professional workflows in this domain [413].

For disease monitoring, an LLM-based intelligent Agribot [414] is trained to address issues such as crop disease detection, weather forecasting, and frequently asked agricultural queries. The model was trained using the Kisan Call Center (KCC) dataset for conversational tasks and the Plant Village dataset for disease detection. In another attempt, use of GPT-based models for diagnosing plant diseases through image

analysis is also demonstrated [415]. The integration of GPT's natural language processing capabilities allows for the synthesis of diagnostic information and actionable recommendations.

Farmers must be versatile, mastering plant growth, protection, legislation, and economics to ensure successful and sustainable agricultural practices. AIChatbots [416] for agriculture are designed to assist in understanding decision support models for managing plant diseases. To enhance crop disease recognition and providing farmer support a multi-faceted approach [417] combining LLMs with Generative Adversarial Networks (GANs), Convolutional Neural Networks (CNNs) is presented. An advanced chatbot was developed in the study using the Langchain Llama Model, which leveraged the capabilities of LLMs to offer an interactive and user-friendly interface for farmers. This chatbot provided real-time responses, insights, and access to government support, making advanced AI technology accessible and beneficial for farmers. In a particular case study for disease recognition for cucumber, end-to-end multi-modal language model (ITMLP) is developed [418]. The model combined image-text multi-modal contrastive learning, image self-supervised contrastive learning, and label information to improve recognition accuracy in a small-sample setting.

To study the potential of LLMs for managing unstructured metadata, converting it between formats, and identifying data collection errors. researcher [419] highlighted how AGI can revolutionize agriculture through technologies like image processing, NLP, and robotics. These advancements aimed to enhance crop yields, reduce waste, and promote sustainable farming. The use of LLMs in metadata management underscores AGI's ability to improve agricultural data handling and decision-making.

To expand institutional capacity and reach LLM-based ChatGPT is used to simplify scientific information and deliver personalized, location-specific, and data-driven agricultural advice to Nigerian cassava farmers [420]. The model (ChatGPT) showed the potential of transforming agricultural extension services by making scientific knowledge more accessible and easier to understand.

In addition to ground-based sensors, various remote sensing (RS) platforms, such as satellites, aircraft's, and unmanned aerial vehicles (UAVs), are employed to gather data with different spatial, temporal, and spectral resolutions for monitoring and management across industries, including oil and gas [421], municipal solid waste landfills [422], [423], and precision agriculture [424]. UAVs enhance the capabilities of satellites by improving the spatial and temporal resolution of data [425].

The optimal resolution of RS data for precision agriculture is influenced by several factors, including management objectives, crop types and growth stages, field size, and the ability of farm machinery to adjust inputs such as fertilizers, pesticides, and irrigation [426]. Conventional RS platforms (e.g., airborne and satellites) face limitations in spatial, spectral, and temporal resolutions. While airborne systems provide high-quality data, they are expensive and impractical for frequent monitoring. UAVs, on the other hand, offer a cost-effective solution, enhancing geospatial data collection and advancing RS capabilities [427].

Recent advancements in the UAV industry have demonstrated that UAV-based RS applications are highly effective in various aspects of precision agriculture (PA), including weed detection [428], [429], [430], soil analysis [431], [432], pathogen monitoring [433], [434], [435], and crop monitoring [436], [437], [438]. During the past decade, there has been a strong increase in the use of drones in agriculture highlighting the integration of UAVs with other subsectors like environmental sustainability, IoT, and AI, demonstrating the coordinated advancements in agricultural technology [439].

Given that LLMs are relatively new, there has been limited work to fully exploit their potential in proximal and remote sensing agriculture monitoring and management. Integrating LLMs with proximal and remote sensing can revolutionize precision farming by enhancing data analysis, decision-making, and real-time monitoring of crop health, pest infestations, and soil conditions. This integration can improve efficiency, productivity, and sustainability in farming through targeted applications of fertilizers and pesticides.

However, LLM-based responses have limitations, including potential biases, inaccuracies, and a lack of domain-specific knowledge. These shortcomings can affect the reliability of LLM-driven insights and decisions in agricultural practices

D. Finance

LLMs are making significant advancements in the finance industry [440] with applications ranging from financial NLP tasks [441], risk assessment, algorithmic trading [442], market prediction [443] and financial reporting [444]. LLM's such as BloombergGPT[99], a 50 billion parameter large language model trained on large diversified financial corpus, has revolutionized financial NLP tasks including but not limited to news classification, entity recognition and question answering. By utilizing the huge amount of financial data available, it is able to enhance customer services drastically by efficiently handling customer queries and providing them with excellent financial advisory.

In addition, LLMs are being used for risk assessment and management, by analyzing past market trends and data, it is able to identify potential risks and provide mitigation steps through different financial algorithms. Financial institutions can use it for better decision making such as credit risk assessment [445], loan approvals and investments [446]. Algorithmic Trading [447] is another application that can leverage LLMs to identify potential opportunities in the trading market by using its predictive and analyzing capabilities.

However, due to the sensitivity of the financial information and privacy concerns, techniques like data encryption, redaction, and data protection policies should be implemented so that these LLMs can be used efficiently in accordance with data protection policies. In this regard, a recent proposition suggested is FinGPT [448] which is an open-source LLM tailored for finance. It is expected that more work will be carried out in this domain.

E. Engineering and similar

LLMs have gained substantial attention across various fields, and their potential applications in engineering domains

are increasingly being explored [449]. In software engineering, CodeGPT can be employed to generate code snippets based on natural language descriptions of desired functionality. Additionally, CodeGPT can assist in debugging code by leveraging its language understanding capabilities to identify errors and suggest potential fixes, thereby streamlining the debugging process and reducing development time[450].

The possibility of LLMs utilization to various calculations in mechanical engineering was attempted in [451]. However [451] encountered instances where incorrect procedures, formulas, or results were provided. None of the tasks yielded an exact solution, leading them to discontinue further research.

In manufacturing, Wang et al. [452] conducted an evaluation of LLMs capabilities in supporting design, manufacturing, and engineering tasks. The results indicate that ChatGPT is impressive in providing information, generating coherent and structured content, and proposing initial solutions. Similarly, Badini et al. [453], performed a study in additive manufacturing troubleshooting and evaluated ChatGPT's expertise in technical matters, focusing on the evaluation of printing parameters and bed detachment, warping, and stringing issues for Fused Filament Fabrication (FFF) methods using thermoplastic polyurethane polymer as feedstock material. It was found that ChatGPT provided remarkable accuracy, correctness, and organization in its responses and its approach to problem-solving offered valuable insights in addressing hurdles.

F. Media and Entertainment Industry

The media and entertainment sector is currently undergoing a transformative phase that revolves around data and prioritizes consumer-centric experiences [454]. Companies of all sizes are now striving to introduce groundbreaking innovations that enable personalized, one-to-one interactions on a large scale [455], [456], [457]. LLMs not only enable the creation of original content but also demonstrate a profound grasp of intricate information and the ability to simulate human-like interactions. This includes MediaGPT, a large language model for the Chinese media domain which was presented recently. It can generate high-quality and relevant outputs for various tasks in the Chinese media domain [458]. Similarly, Robertuito [459] was proposed for Spanish social media.

Large AI models can also be utilized for generating attractive advertisements and marketing [100], political speeches, slogans and social media posts [460], and promotional videos [117]. Similarly, leading entertainment networks and applications are using LLM based algorithms that can analyze user data to offer personalized recommendations for movies, TV shows, and music. This helps entertainment companies to retain customers and improve their engagement with their content [461]. Moreover, LLMs automate content curation fostering user satisfaction, retention, and monetization. Recently, many companies have developed and offered their services for media and entertainment purposes. One of the prime examples of such services is Dolly, an LLM-trained model developed by databricks Incorporation [462].

The creation of AI-based newscasters [463] is a recent concept that consists of virtual news presenters or anchors that

TABLE VI: Unveiling the AI Revolution in Entertainment: Real-world Illustrations

Tools	Function
Scriptbook	A cutting-edge AI-powered script analysis tool, is harnessed by film studios to forecast the commercial triumph of a screenplay.
Aiva	AIVA (Artificial Intelligence Virtual Artist) represents an AI-driven music composition tool that generates original music tracks tailored to user preferences.
LyricFind	LyricFind takes center stage as an AI-powered lyrics search engine, empowering users to find song lyrics using natural language queries.
Ziva Dynamics	Ziva Dynamics showcases an AI-powered software tool tailored for creating authentic 3D character models in films and video games.
DeepMotion	DeepMotion introduces an AI-powered animation tool capable of producing lifelike 3D animations for video games and films.
Speechify	Speechify is one of the most popular and efficient first AI Voice Over generators for using famous singer's voices for singing different songs.

are generated using AI technologies, particularly LLMs [464]. In April 2023, a Kuwaiti media outlet unveiled a virtual news presenter "Fedha" with plans for it to read online bulletins [465]. At the University of Kent's Centre for Journalism, lecturers are grappling with how to prepare the next generation of reporters for the potentially AI-powered newsrooms of the future [466]. AI algorithms have the capability to analyze user data, providing tailored suggestions for movies, TV shows, and music. This enhances customer retention and boosts engagement with entertainment content. Table VI presents the recent tools and applications that are transforming the entertainment industry.

G. Role of LLMs in the Future of Legal Practice

With advancements in AI and the development of tools such as GPT-4, Bard, Gemini, and Bing, it is aimed that these advancements will empower lawyers to enhance legal research, drafting tasks, and decision-making [467]. This has sparked interest among entrepreneurs developing AI tools [468], law firms integrating AI into their workflow, and law professors exploring AI-based techniques for legal aid [469]. A recent example is the Chatlaw [470] model, which is open-source legal language model. A legal informatics approach was introduced in [471] to align AI with human goals and societal values. By incorporating legal knowledge and reasoning into AI systems, the paper contributes to the research agenda of enhancing the integration of AI and law. In [472], the authors propose legal prompt engineering (LPE) as a means to improve LLM performance in legal judgment prediction tasks. The effectiveness of this method has been demonstrated on three multilingual datasets, showcasing the model's capability to handle the intricacies of legal language and reasoning from various sources of information. LLMs' transformative potential in the legal field is evident from their

impressive performance in legal exams. GPT-4 scored in the 90th percentile on the Uniform Bar Examination [61], and ChatGPT autonomously passed four law school final exams at a top law school [473]. These achievements showcase the significant impact of AI language models on legal practice. The authors present Chain-of-Thought (CoT) prompts, which aid LLMs in generating coherent and contextually relevant sentences following a logical structure, simulating a lawyer's analytical approach [474]. The study shows that CoT prompts outperform baseline prompts in the COLIEE entailment task using Japanese Civil Code articles. LLMs have also been utilized to explore fiduciary obligations [55].

In a recent working paper by Choi et al., the authors conducted experiments using ChatGPT to generate answers for four authentic exams administered at the University of Minnesota Law School [475]. In summary, the authors concluded that ChatGPT successfully passed all four exams with an overall average grade of C+. This level of performance would grant it credit towards a JD degree, but it would also place the student on academic probation.

Recently in June 2023, in response to fake case citations generated by ChatGPT and submitted in a court filing, a US judge has imposed a fine of \$5,000 (3,935) on two lawyers, along with their law firm Levidow & Oberman [476] [477], due to the fake citations generated by ChatGPT. The fictitious legal research was utilized in an aviation injury claim, and the lawyer admitted to inventing six non-existent cases citations. In Texas, a judge now requires attorneys to verify that no part of a filing was composed by generative AI or, if it was, that a human has verified its accuracy [478] [477]. However, not all judges share the same stance on chatbots in legal proceedings. For instance, Judge Juan Manuel Padilla, based in Colombia, acknowledged using ChatGPT's assistance in a case concerning an autistic child [478].

In light of these examples and use cases, LLMs indeed offer numerous benefits, but it is crucial to recognize and comprehend their limitations. They can serve as a valuable tool for initial research, explanations, and improving efficiency in legal practice.

H. Marketing

Large language models excel in content generation, creating compelling product descriptions, ad copy, blogs, and social media posts, saving time, analyzing SEO optimization, identify keywords, and enhance social media monitoring, vast data, including feedback and social media, offering insights into trends, sentiment, and resonating with audiences [479], [480], [481]. Personalization is a standout feature, allowing marketers to deliver tailored messages based on customer data, improving satisfaction and loyalty [482].

The adoption of LLMs and ChatGPT in marketing offers numerous benefits, but it also comes with potential risks for marketers, consumers, and other stakeholders [481]. The similarity and lack of uniqueness in ChatGPT's responses to similar prompts from different marketers could undermine the distinct identity of the marketer or brand [483] [479] [484]. AI marketing tools like ChatGPT may draw information

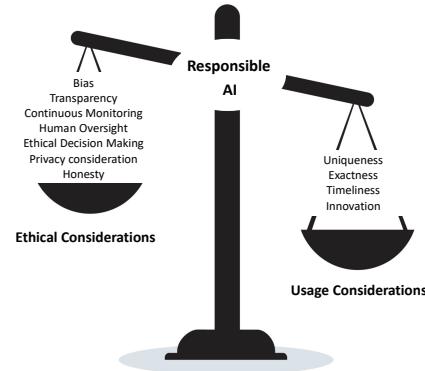


Fig. 10: The Tradeoff for Responsible AI: Example of marketing.

from unreliable sources, leading to the provision of incorrect information that can lead to false outcomes [481] [485]. Ethics is a significant concern as LLMs can generate content that appears human-generated, raising transparency and disclosure issues [486].

Negative consumer perceptions could arise if AI-generated content is overwhelming or perceived as inauthentic, leading to reduced trust in the brand. Compliance with regulations and data protection laws is crucial to avoid potential legal consequences [487]. Moreover, marketers may become dependent on third-party AI providers, leading to vendor lock-in or reliance on external platforms. To mitigate these risks, marketers must use AI responsibly, provide clear disclosure when AI is involved, and maintain human oversight to verify the accuracy and appropriateness of AI-generated content [488]. It is at the interjunction of the risks (market and consumers both) and ethics where we can balance the scale which opens the window of opportunity for responsible AI adoption (see Figure 10).

1) Customer Service : Another application that has witnessed a significant impact is customer service. LLM-powered chatbots and virtual assistants are increasingly being integrated into customer support systems, providing companies with a scalable and efficient means of addressing customer inquiries and concerns [489], [490], [491]. Unlike human representatives, LLM-powered chatbots can process and respond to inquiries instantaneously, enhancing the overall customer experience [492]. The use of LLMs also leads to cost savings for companies [493]. Due to the high cost associated with employing dedicated customer service agents, an increasing number of companies are exploring the use of NLP to assist human agents [494]. NLP enables the auto-generation of responses that can be directly utilized or modified by agents. In this context, LLMs emerge as a natural and suitable solution.

VII. LLM INTEGRATION WITH A DATABASE

As LLMs continue to advance, a new frontier lies in their integration with structured query languages (SQLs) such as PostgreSQL. By coupling LLMs with relational databases, users can navigate complex data sources using natural language prompts rather than conventional query syntax. This approach simplifies data exploration, enhances decision-making,

and lowers the barrier for individuals who lack deep expertise in database query languages.

A. Motivations and Foundations

LLMs excel at understanding human language and generating contextually relevant responses. Rather than writing precise SQL queries, a user might ask, “Which customers from California purchased more than \$5000 worth of products last month?” The LLM can translate this request into SQL behind the scenes. Traditionally, such tasks required knowing the schema, the correct table names, and query syntax. LLMs abstract these details, allowing more intuitive interactions with structured data.

Relational databases store data in well-defined schemas, and SQL provides a declarative language for filtering, aggregating, and joining across multiple tables. LLMs, trained on large corpora, understand linguistic patterns and can adapt these skills to produce SQL queries given a natural language description of the user’s intent. This synergy enables users to access complex, large-scale datasets without specialized training in SQL.

B. Setting up the Environment

Establishing an LLM-to-SQL pipeline involves selecting appropriate tools and frameworks. Credentials, schemas, and access policies must be securely configured. Common relational systems like PostgreSQL or MySQL are often integrated. Role-based access controls, combined with environment variables, ensure that sensitive connection details are protected.

C. Leveraging LangChain and LangGraph

Frameworks such as LangChain simplify connecting LLMs to data sources. LangChain [495] is a framework designed to build applications powered by large language models (LLMs), enabling seamless integration of LLMs with external tools, data sources, and workflows to create dynamic, AI-driven solutions.

Initially designed for diverse retrieval scenarios, LangChain can be adapted for SQL databases by providing schema context and representative examples of user queries alongside their corresponding SQL commands. Its modular design streamlines prompt construction, enabling the LLM to understand how to map natural language requests to SQL queries reliably.

LangGraph, another emerging tool, emphasizes the graph-based representation of data and relationships. While traditional relational databases and SQL queries often follow a tabular, schema-defined structure, certain data domains or enterprise contexts may benefit from representing relationships as graphs. LangGraph can guide the LLM to leverage these graph structures, translating user queries into SPARQL or other query languages that handle graph data. In scenarios where data spans both relational and graph formats, combining LangChain’s SQL-oriented approach with LangGraph’s graph-centric capabilities can yield a powerful hybrid system. This ensures that the LLM can provide intuitive, natural language interfaces over a variety of data models, enhancing flexibility and user-friendliness.

D. Building a Conversational Interface

A conversational UI—such as a chat window—lets users pose questions in everyday language. When the user submits a prompt, the LLM interprets their intent, references database schema details (possibly provided via LangChain’s prompt engineering), and generates the corresponding SQL. The database then executes the query and the results are returned in a readable format. Incorporating vector databases or embeddings into the workflow can further improve the LLM’s ability to handle ambiguous requests and map them to the correct data sources.

E. Enhancing Query Complexity and Accuracy

As use cases evolve, the system must handle more complex queries: multi-table joins, aggregations, and nested subqueries. Prompt engineering, including the careful selection of examples and schema descriptions, guides the LLM toward accurate SQL generation. Over time, reinforcement learning techniques—where users provide feedback on generated queries—can refine the LLM’s accuracy. Additionally, integrating schema-specific hints and foreign key information helps the LLM understand which joins make logical sense, reducing query errors.

F. Addressing Common Challenges

LLM-driven query generation can sometimes produce malformed SQL or inefficient queries. Syntax validation and fallback strategies mitigate such issues. For performance-sensitive environments, indexing frequently accessed columns, caching results, and optimizing schema design maintain responsiveness. Ensuring that the LLM cannot issue destructive queries or bypass permissions requires strict security and authorization checks. Tools like LangChain and LangGraph can assist by providing controlled interfaces that dictate how LLMs access and utilize database metadata, ensuring that only safe, valid queries are produced.

G. Future Directions

Further advancements may include integrating vector search with SQL queries, allowing users to discover data through semantic similarity rather than exact keyword matching. Explainable query generation could help users understand why certain queries were chosen, building trust and transparency. Domain-specific fine-tuning might adapt LLMs for specialized fields, offering more context-sensitive and domain-relevant responses. Similarly, improved support for graph-based queries via LangGraph will enable cross-domain data exploration, as many organizations store related information in both relational and graph formats.

H. Conclusion

By fusing LLMs with SQL and related data technologies, developers can create intuitive interfaces that democratize data access. Users can simply ask questions in natural language and receive structured, meaningful answers without mastering

SQL syntax. With frameworks like LangChain and LangGraph supporting the connection between language models and databases, the ecosystem will continue to evolve - producing more accurate, efficient, and user-friendly data exploration tools. As these methods mature, LLM-integrated data systems are likely to become a standard approach, enabling a wider range of users to harness the full potential of their data.

VIII. AI-ENABLED TOOLS: THE FUTURE OF EVERYTHING

AI tools are becoming increasingly powerful and versatile. They can be used to generate text [496], images [497], and videos [498], translate languages [341], write different kinds of creative content [499], and answer the users questions in an informative way [500]. These powerful tools provide a user-friendly interface for the optimization of daily routine tasks [85]. One such example is the popular website, "There's an AI for THAT"⁵, which contains about $7K$ AI tools for $2K$ different tasks. In this Section, we discuss various AI-enabled tools based on LLMs or text prompts.

A. Chatbots / ChatGPT

Chatbots are frequently used in customer service applications where they can respond to queries, offer assistance, and fix problems [50]. High-tech companies are likely to become even more interested in using chatbots to improve their customer experience and grow their businesses. For example; OpenAI developed ChatGPT [501], Google developed Bard [502], and Meta launched Llama-2 [15]. Here, we critically compare these Chatbots in terms of accuracy, ease of use, cost, integration and others.

1) Comparison between Chatbots: ChatGPT and Google Bard are two of the most popular LLMs available today [503], [504]. The third popular LLM being Bing, which is based on GPT-4. Bard is based on the LaMDA [505] (Language Model for Dialogue Applications) architecture, while ChatGPT is based on the GPT-3 (Generative Pre-trained Transformer 3) architecture. ChatGPT was modified and improved using both supervised and reinforcement learning methods [506], with the assistance of human trainers (RLHF) [80]. The learning includes three steps;(i) supervised fine-tuning [507], reward model [508], and maximum policy optimization [509]. First, initiate with pre-trained GPT-3, fine-tune with labeled data. Generate 4-7 responses for input prompts, rank by labelers, creating scalar values for a reward model. Test on new sequences, evaluate responses with reward model, fine-tune parameters for human-like traits through reinforcement learning [510].

LaMDA is a newer architecture (conversational neural language models) that is specifically designed for dialogue applications. Bard uses LaMDA, which is a hybrid architecture that combines batch processing [511] and streaming processing [136]. This allows BARD to handle both historical and real-time data. It is trained on a massive dataset of text and code, while ChatGPT is trained on a massive dataset of text, which means that Bard has a broader understanding of the

world and can generate more comprehensive and informative responses, while ChatGPT is better at generating creative and interesting responses.

Both models are capable of generating text, translating languages, writing different kinds of creative content, and answering your questions in an informative way. However, there are some key differences between the two models, such as ChatGPT is more creative, while Google Bard is more authentic. Bard is more personalized than ChatGPT, the responses generated by Bard are more tailored to specific needs, and it is also more scalable than ChatGPT. A comparison between ChatGPT, Bard, and BingChat is made in [502] on VNHSGE [512] dataset, which is a Vietnamese High School Graduation Examination Dataset for Large Language Models.

The results indicate that BingChat performed better than Bard, and ChatGPT. All models perform better than the Vietnamese students [513]. In fact, a comparison between the three popular LLMs services, namely, ChatGPT, Bard, and Microsoft Bing has been of interest to researchers and field practitioners. A recent work by Bhardwaz et. al [514] provided a general comparison of these three models considering accuracy, response time, user experience, and engagement. From their experiments, they found that ChatGPT provided the most relevant responses and accuracy, Bard provided the quickest response and Bing provided the best user experience and engagement. Another comparison by Campello et. al [515] experimented with four different chatbots (above three and Quoras Poe [515]) when asked to solve an intelligence test for recruitment in Brazil found that all four chatbots scored above the 95th percentile while ChatGPT and Bing scored 99th percentile. These are in addition to comparisons being made for typical as well as atypical specific use cases such as news fact-checking (GPT-4 performing the best) [516], taxes [517] as well as political leaning [518] and more. To complete the discussion, Table IX presents a comparison between ChatGPT, Google Bard, Llama-2, and Microsoft Bing Chatbots.

B. AI tools for image generation, history, medical, industry

1) Diffusion Models: Diffusion models are a scheme of generative models that have provided excellent performance in a variety of applications, most notably the synthesis of images [519]. Starting from a sample of a target data distribution, a diffusion model works in two steps, a forward diffusion process and a reverse diffusion process. The forward diffusion process gradually adds increasing amounts of Gaussian noise[520] to the sample image successively over time, until it becomes the Gaussian distribution [521]. The model is then tasked to start from this noisy image version and undo the noise addition by going through a reverse process to recreate the original data [522], [20].

The forward process takes the form of a Markov chain [523] where the distribution at a given time instant only depends on the sample from the timestep immediately preceding it. Therefore, the distribution of the corrupted samples at any given point with respect to the original sample is the product of the successive single-step conditionals up till that point.

Moreover, typically the number of passes of noise addition is in the order of a thousand with the increments each time

⁵<https://theresanaiforthat.com/>

being quite small. This is necessary to ensure that the reverse process of “recovering” the original sample is more achievable as it has been shown that with infinitely small step sizes, the reverse form will be able to achieve the same functional form as the forward process [524]. Diffusion models use this observation in the forward process.

A diffusion model can be interpreted as a latent variable generative model similar to a variational autoencoder (VAE). The forward process can be thought of as producing latent from data and the reverse process is as converting latent to data. However, as opposed to VAEs, the forward process for diffusion models is typically fixed so that only a single network needs to be trained that deals with the reverse process. The objective function is the variational lower bound on the log-likelihood of the data. It consists of a log-likelihood term or reconstruction term minus a KL divergence term [525] also called the regularization term [526]. The log-likelihood terms [527] encourage the model to maximize the expected density assigned to the data. The KL divergence term encourages the approximated distribution to the prior distribution on the latent variable. Moreover, diffusion models can also be directed to sample conditionally based on a variable of interest which can be incorporated as an additional input during training. This has been the reason that diffusion models have shown better performance than Generative Adversarial Networks [528] in a variety of image generation tasks including perceptual quality ([529]), text to image generation ([530]), image inpainting [531] and manipulation of images[532].

2) *Image generation:* The images contained in this section were generated by a model incorporating the stable diffusion process in to existing diffusion models as suggested in [533] and uses text to generate photorealistic images. This model was released by stability.ai [534] and was demonstrated to be capable of generating images which were previously difficult to generate, such as images of people with accurate facial features as well as objects with abnormal or impossible shapes. Table VII showcases the output of image generation using various prompts.

In total, nine different prompts were used, these required the AI model to generate humans and natural scenery. The first four prompts tended to depiction of famous personalities (sportsmen and politicians in this case), Muhammad Salah, Lionel Messi, Mike Tyson and Imran Khan. The prompts used were *Mo Salah playing cricket*, *Lionel Messi playing tennis*, *Mike Tyson playing football* and *Imran Khan as a hero*. The second prompt used was regarding the famous painting Monalisa. The prompt was "Generate an image of Monalisa showing her teeth in a wedding ceremony". The third prompt related to natural scenery and was written as *Area of rocks, deep inside the forest, divine domain*. Lastly, the fourth prompt also centered around the generation of humans. In this case, three prompts were given, *A man kissing a girl*, *Generate an image of a guy* and *Generate an image of a woman*.

3) *Video Generation using text prompts:* Text-to-video generation is a challenging task that involves generating video sequences from textual descriptions or prompts [498]. The video generation process involves text understanding [535], video scene generation [536], temporal structure and transi-

tion [537], and video synthesis [538]. One such model is T2V [539], which is a video generation model using text prompts.

C. AI tools for text classification

AI tools are increasingly being used for text classification. Text classification is the process of assigning a category to a piece of text [540]. For example, a text classification tool could be used to classify emails as spam or not spam or to classify news articles as business, sports, or entertainment. Some of the popular libraries include Scikit-learn, NLTK [541], and Spacy [542]. Similarly, Hugging Face’s Transformers library [543] is the state-of-the-art toolkit for developers to implement AI text generation capabilities into their applications; including fine-tune models for sentiment analysis, language translation, and text summarization.. This library offers a collection of pre-trained language models, including GPT-3.5 and various other popular models like Bidirectional Encoder Representations from Transformers BERT [69] and RoBERTa [70]. Openai’s foundation video model, Sora [544] became well known recently for its ability to generate videos up to a minute long while maintaining quality and adherence to the user’s prompt.

D. AI tools for Literature Review Research

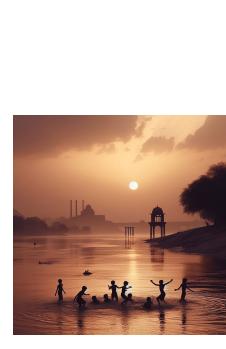
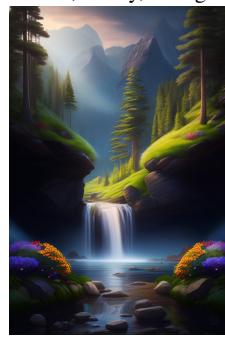
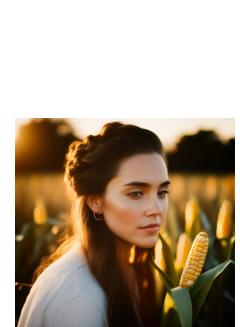
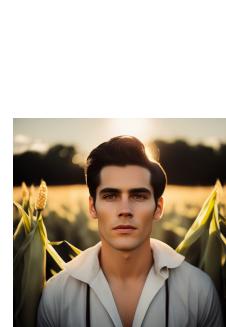
AI tools are increasingly being used to assist with literature review research. These tools can be used to automate tasks such as: Identifying relevant literature, extracting information, and summarizing the content [545], [546]. One such tool is PDFGPT [547], which uses the GPT-3 model to generate responses to user queries.

1) *Fake references:* One of the major drawbacks of using AI tools such as ChatGPT in research is the creation of fake citations and references [548], [549]. Fake citations is an inherent consequence of the generation capabilities of LLMs wherein they may prefer to lean on their generation capabilities rather than on search. When an LLM is asked to provide citations, given its inherent generative nature, it might end up generating text that looks like a research article/source but in reality, might just be a sequence of words that are objectively similar in linguistic/language generation terms [550].

The potential complications that are being created due to the uncontrolled usage of these tools result in many issues among which misleading the scientific community carries vital importance. Fake citations can undermine the credibility of the author [551]. This can lead to inaccurate conclusions and potentially harmful decisions being made based on faulty information. Using fake citations and references can hide the true sources of information used in the research, making it difficult for others to replicate or verify the findings [552]. Recently, WebChatGPT⁶ is an impressive extension that has the potential to address the pervasive issue of fake citations. With the installation of this extension, WebChatGPT becomes equipped with robust capabilities to detect and eliminate fake citations. This advanced tool uses sophisticated algorithms to

⁶<https://tools.zmo.ai/webChatGPT>

TABLE VII: Image generation examples

Prompt:**Negative Prompt:**Different famous personalities in roles other than their original ones
blurry, photorealistic**Generated Images:****Prompt:****Negative Prompt:**Generate an image of Monalisa showing her teeth in a wedding ceremony
blurry, low resolution, artistic**Generated Images:****Prompt:****Negative Prompt:**Area of rocks, deep inside the forest, divine domain, river, sunset, kids playing
artistic, blurry, background**Generated Images:****Prompt:****Negative Prompt:**A man kissing a girl/ Generate an image of a guy/ woman
artistic, blurry, background, young**Generated Images:**

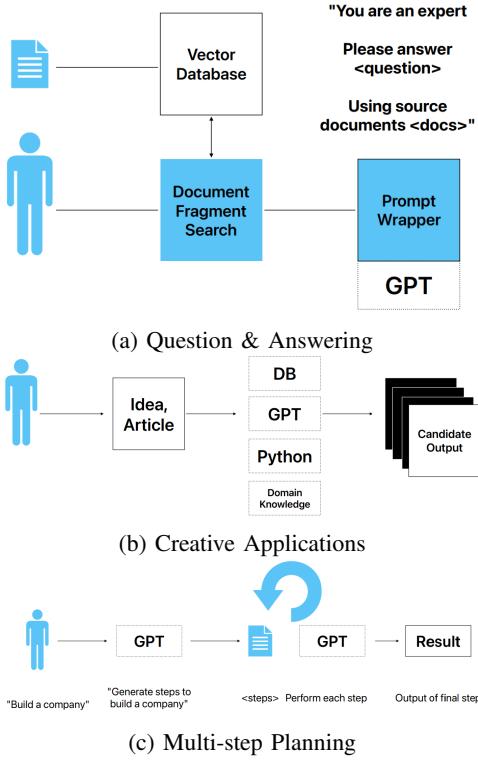


Fig. 11: Templates for LLM-based application development. GPT is taken as an example scenario representing LLMs.

analyze the authenticity and reliability of citations, ensuring that only accurate and legitimate sources are included.

E. AI tools for Coding and Development / CodeGPT

AI tools are increasingly being used to help programmers write code. These tools can be used to automate tasks such as code completion, refactoring, linting, and testing [553]. GitHub Copilot [554] is an AI-powered code completion tool developed by GitHub in collaboration with OpenAI. It utilizes OpenAI's GPT-3 language model to assist developers in writing code more efficiently. Meta also released the CodeLlama [555], a LLM model, that can use text prompts to generate and discuss code. It has the potential to generate clean and robust code with well documentation in Python, c/C++, Java, PHP, Typescript (Javascript), Bash and other programming languages.

Developers can interact with the LLM by providing prompts or queries related to their coding needs, and the model can generate relevant code segments or suggest solutions to programming problems [556]. LLMs have been used to develop applications in three primary categories which include: (a) Question Answering, (b) Creativity (c) Multi-step planning [556]. These template categories are illustrated in Fig. 11. In Table VIII, we present a list of publicly available AI/LLM tools for a variety of applications.

F. Retrieval Augmented Generation

IX. GPT-PLUG-INS

GPT-Plugins are a new way to extend the functionality of ChatGPT. They allow developers to create custom apps that can be integrated into ChatGPT, providing users with new features and capabilities [569]. GPT-Plugins can be used to do things, such as access to external data sources, automate tasks, and enhance user experience [570]. In this Section, we demonstrate several GPT-Plug-ins.

A. ChatGPT Plug-ins

Arguably, the watershed event in the use of ChatGPT was the introduction of plugins by OpenAI. Plugins allow ChatGPT to communicate with third-party sources of data and knowledge bases, thereby providing a platform to extend ChatGPTs capabilities for composition, summarization, nuanced tasks such as sentiment analysis and more to any resource on the internet. Moreover, given that ChatGPT has provided sufficiently acceptable performance for various tasks, plugins allow for ChatGPT to provide answers to queries with updated information from the internet which may not be present in its training dataset [508].

This also has the advantage of providing references for queries to add credibility to answers. For e.g., Bing, the search engine by Microsoft works with OpenAI's ChatGPT through its API to allow its users to ask questions from its Bing search system and get answers with references/sources mentioned. The integration of LLMs in to search engines, thereby allowing users to get answers to human like queries has spearheaded the search engine business in to a new direction [571]. Moreover, this addition of credibility is an important consideration to enable use of ChatGPT and similar LLMs in other critical tasks.

While, at the time of this manuscript, OpenAI still hasn't rolled out plugin development access to all developers, there have been several notable use cases that have already come out. For example, twelve companies have been listed on the OpenAI website ⁷, namely, Expedia, FiscalNote, Instacart, KAYAK, Klarna, Milo, OpenTable, Shopify, Slack, Speak, Wolfram, and Zapier to have created the first plugins. The power that plugins provide in terms of flexibility to develop new applications has drawn a big attention towards plugin development. Apart from the above-mentioned early developers, two notable plugins already made available by OpenAI are the Code interpreter and the knowledge-based retrieval plugin.

- **Code Interpreter:** The Code interpreter is a built-in Python code interpreter which can be used for performing logical calculations as well as writing code. The interpreter can use the language model's understanding of a human language description of a problem and use that as input to develop Python code for the problem's solution.
- **Knowledge-base retrieval:** A knowledge-based retrieval plugin has also been open-sourced⁸ which can be used by developers. This plugin can be used to enable ChatGPT

⁷<https://openai.com/blog/ChatGPT-plugins>

⁸<https://github.com/openai/ChatGPT-retrieval-plugin>

TABLE VIII: Publicly available AI /LLM tools

Tools	Function	Link	Availability
ChatGPT	Conversational AI Chatbot	ChatGPT	Both
RoomGPT	Redesign your room in eight different themes	RoomGPT	Public
HomeGPT	Redesign your home and office	HomeGPT	Subscription based
PDFGPT.IO	Turns PDF into the knowledge base for a ChatGPT type interface	PDFGPT	Subscription based
TexGPT	Harnesses GPT-3's power to help you write in Overleaf	TexGPT	Public
AcademicGPT	An AI tool to write and review scientific papers, critical analysis and explanation of complex concepts	AcademicGPT	Public
DiagramGPT	An AI tool for creating scientific diagrams and flow charts of different processes	DiagramGPT	Public
AutoGPT	Auto-prompting without the user intervention	AutoGPT	Public
HuggingGPT [557]	A framework to connect various AI models to solve AI tasks	HuggingGPT	Public
XrayGPT [558]	Automated analysis of chest radiographs	XrayGPT	Public
Video-ChatGPT	A vision language model for video understanding and conservation about videos	Video-ChatGPT	Public
ClimateGPT	Large language model for a conversation about the climate in English and Arabic	ClimateGPT	Public
CodeGPT	An AI assistant for coding	CodeGPT	Public
Code Llama	Open Foundation Models to generate and discuss code	Code Llama	Public
MiniGPT-4 [335]	Multi-modal model for a number of tasks, including image generation and website development, using prompts	MiniGPT	Public
SearchGPT [559]	Information Retrieval with Advanced Language Models	SearchGPT	Both
RAG [560]	An AI framework that combines generative large language models (LLMs) with traditional information retrieval systems	PostgressML	Both
SoRA [544]	AI based tool for text-to-video generation	SoRA	Both
BiomedGPT [561]	A Unified and Generalist Biomedical Generative Pre-trained Transformer for Vision, Language, and Multi-modal Tasks	BiomedGPT	Public
PatientGPT	An AI engine to transform patient navigation with a seamless and customized experience	PatientGPT	Subscription based
SentimentGPT [562]	Exploiting GPT for sentiment analysis	SentimentGPT	Public
DrugGPT [387]	A GPT based model to design potential ligands, targeting specific proteins	DrugGPT	Public
Elicit [563]	AI research assistant, automated literature reviews	Elicit	Public
Citation AI	AI research assistant to generate real evidence-based answers	Citation AI	Subscription based
Gen-2 [564]	Video generation using text, images, and videos	Gen-2	Public
AI Avatar	Avatar generation	AI Avatar	Public
Langchain [565]	Building applications with LLMs through composability	Langchain	Public

to access data and then use it to gather useful or relevant information from the data. These can be files, emails, notes etc. All this by using queries or questions in normal

human language.

Lastly, third-party plugins are also an option. These can be created and have been created by several entities. The use of

TABLE IX: Comparison of Bard, ChatGPT, Llama-2 and Bing Chat

Feature	ChatGPT (GPT 3.5)	Bard	Bing Chat (GPT - 4)	Llama2
Accuracy	Not as accurate as Bard	Generally more accurate than ChatGPT	Most accurate	least accurate. ⁴
Versatile	Generally more versatile than Bard	Can generate text, translate languages, and write different kinds of creative content	Not as versatile as ChatGPT or Bard	Less than ChatGPT and Bard both better than Bing
Company	OpenAI	Google	Microsoft	Meta
Primary Purpose	Creative text generation	Conversational AI	Information retrieval	Text generation, answer questions, language translation etc
Integration	Standalone model	Standalone model	Integrated with Bing search engine	Standalone model
Easy to use	User-friendly	User friendly	Not as user-friendly as ChatGPT or Bard	User-friendly
Access to online data	No, trained on data available till 2021	Yes	Yes	Yes
Cost	GPT 3.5 free / GPT-4 (20 USD per month)	Free	Free	Free
Availability	Publicly available	Publicly available	Publicly available	Publicly available
Architecture	Generative pre-trained transformer [566]	Pathways Language models (PaLM2) [567]	Next GPT [568]	Transformer
Plagiarism detector	Yes	No	No	Less likely to generate plagiarised text
Limitations	May generate less coherent or incorrect text	Not as creative as ChatGPT	May provide limited information	Trained on a smaller dataset

two third-party plugins, namely ShowMe which can be used to generate diagrams and ScholarAI can be used to access academic journals. Table X provides a list of plugins available for ChatGPT.

X. GUIDELINES FOR EFFECTIVE USE OF LARGE LANGUAGE MODELS

In this section, we provide a list of steps as well as guidelines for effective and responsible use of LLMs [578].

A. Model selection and deployment guidelines

- Identify the task:** Determine the task, LLMs can be used for a wide range of tasks, such as text classification, sentiment analysis, question answering, and text generation [579], [580], [581].
- Choose the right model:** Choose a pre-trained LLM that is suitable for your task. There are several pre-trained LLMs available, such as GPT-3, BERT, and RoBERTa. Each model has different strengths and weaknesses, so it's important to choose the one that best fits your needs [256].
- Fine-tune the model:** Fine-tuning involves adjusting the model's parameters, such as learning rate, batch size, and number of epochs, to optimize its performance on your task [582].

- Evaluate the model:** Evaluate the performance of the model on a test dataset. This involves measuring the accuracy, precision, recall, and F1 score of the model on the test dataset [583].
- Deploy the model:** Deploy the model in your application or system. This involves integrating the model into your application or system and exposing it through an API or user interface. This step also involves setting up monitoring and logging to track the performance of the model in production [584].
- Monitor and retrain the model:** Monitor the performance of the model in production and retrain it as needed. This involves regularly checking the performance of the model and identifying any areas for improvement [501].
- Continuously improve the model:** Continuously improve the model by incorporating user feedback and updating it with new data. This involves collecting feedback from users and incorporating it into the model to improve its performance [585].

B. Ethical guidelines

The following guidelines will help to ensure the responsible development and use of LLMs focusing on user privacy, bias mitigation, ethical considerations, transparency, competition, collaboration, and environmental impact [80].

- Protect User Privacy:** LLMs should uphold user privacy and protect user data. This includes safeguarding user-

TABLE X: Some ChatGPT Plugins. This list is not exhaustive and more and more plugins are being developed.

Name	Task	Example use cases
Language Translation [341]	Translate between languages	This is particularly useful for documents and information from different languages might need to be translated.
Sentiment Analysis [572]	Determine tone of text or conversation	This can be used for the customer analysis and social media monitoring.
Spell Checker [573]	Check and correct spelling mistakes	This service can be useful for formal and informal communication such as emails and also browsing the web.
Question-Answering [574]	Answer questions for a user query	This can find use in education to build learning platforms, search engines, especially when a more 'understandable' response is required. Knowledge graphs can be used for improving on search queries and creating recommendations.
Knowledge Graph [575]	Find and present information from a database	
Speech Recognition [576]	Understand and transcribe speech audio	This service can be used in audio based customer service and also provide services to differently-abled people through audio
Emotion Detection [577]	Detect emotion from text or audio	This service can be used for applications relating to market research using verbal ques, used for healthcare as well as assessing reactions to games and other media

generated content, such as emails, messages, and personal information. Best practices should be followed, such as data minimization, anonymization, and encryption, to ensure user privacy is not compromised [586].

- **Mitigate Bias:** LLMs can inherit and amplify biases present in the data they are trained on. Developers and researchers should actively identify and mitigate bias in their models. This can be achieved through diverse and inclusive training data, bias detection techniques, and evaluation metrics [587].
- **Address Ethical Implications:** LLMs have the potential to be used for harmful purposes, such as spreading disinformation or generating deepfakes [588]. Ethical considerations; including, ensuring accountability, transparency, and responsibility in the development and deployment of models must be taken into account [589].
- **Foster Transparency:** It is crucial that the inner workings of LLMs are transparent and explainable. This can help build user trust and facilitate understanding of the model's behavior. Explainability techniques, such as attention mechanisms and model interpretation tools, can be employed to provide insight into the decision-making process of models [450].
- **Promote Competition:** The development and deployment of LLMs should not be monopolized by a small number of companies or individuals. This can limit innovation and negatively affect competition. Collaboration between academia, industry, and government can foster competition, while also promoting responsible development and use of models [590].
- **Encourage Collaboration:** Collaboration between researchers, developers, and industry should be encouraged to promote the responsible development and use of LLMs. This includes open sourcing models and data, as well as facilitating the sharing of research findings and best practices [591].
- **Minimize Environmental Impact:** Training LLMs can require significant computational resources and energy, which can have negative environmental impacts. Developers should strive to create more energy-efficient models and explore alternative training methods, such as model distillation or transfer learning, to reduce the

environmental footprint of models [592], [593].

- **Optimization is exploitation:** is a statement that holds particular significance in the context of LLMs and AI technologies [594]. While these technologies have the potential to revolutionize the way we live and work, they also have the potential to perpetuate existing inequalities and introduce new forms of exploitation [595]. Therefore, it is important to carefully consider the ethical implications of optimization in the development and deployment of LLMs and AI technologies [77].

XI. CHALLENGES AND LIMITATIONS OF LARGE LANGUAGE MODELS

Although LLMs have made significant contributions to various domains, they have significant limitations and challenges [105], [114]. LLMs are currently perceived as forerunners of Artificial General Intelligence (AGI). However, despite their phenomenal success in conversational tasks, the state-of-the-art LLMs still lack in many aspects that makes them less likely an early manifestation of AGI. We first provide a quick list of the challenges and limitations of LLMs (Fig. 12) and then present a more detailed discussion on a few limitations of critical concerns.

A number of challenges and limitations have been focused on, including biased data, overreliance on surface-level patterns, limited common sense, poor ability to reason and interpret feedback [596], [597]. Other issues include; the need for vast amounts of data and computational resources [598], limited generalizability [599], lack of interpretability [600], difficulty with rare or out-of-vocabulary words, limited understanding of syntax and grammar [601], and limited domain-specific knowledge [602].

The susceptibility to adversarial attacks [248], ethical concerns [75], difficulty with context-dependent language [255], absence of emotion and sentiment analysis [603], limited multilingual capabilities [604], limited memory [605], lack of creativity [499], and restricted real-time capabilities [579] are also critical concerns. High training and maintenance costs, scalability issues, absence of causality, difficulty with multimodal inputs, limited attention span, constrained transfer learning, incomplete world understanding beyond text, insufficient comprehension of human behavior, restricted long-form text generation, limited collaboration, ambiguity handling

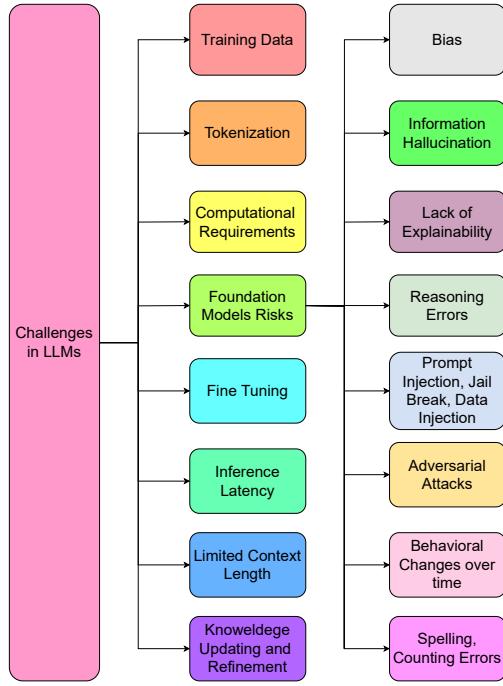


Fig. 12: Challenges in LLMs.

challenges, inadequate understanding of cultural differences, constrained incremental learning, limitations in handling structured data, and difficulty with noise or errors in input data [606], [607], [608], [609], [610], [611], [92] are some of the key challenges in safe, responsible, and efficient deployment of LLMs.

A. Training Data Requirements

Large Language Models (LLMs) require a large corpus of data for pre-training the model. Collecting and curating these datasets can be extremely challenging. The size of the dataset makes it impossible to read or assess the quality of the dataset making it prone to having duplicates, making the model biased and degrading its responses [612] [613]. It also makes it difficult to assess the model as the training data may contain data similar to testing samples leading to incorrect evaluation metrics. Since there is no way of checking the datasets manually, it may contain confidential or personal information too, such as telephone numbers leading to privacy leaks during prompting [596]. Due to the fact that the data distribution and requirements in LLMs is more of a black box, it remains uncertain what amount of data is required for different tasks.

B. Tokenization Problems

LLMs heavily rely on tokenization which consists of breaking down a sequence of words into tokens for the models input. Most LLMs use subword tokenization [56], which is used to create tokens by splitting the words to handle non-familiar vocabulary and at the same time maintaining the computational complexity. However, there are some major drawbacks of

tokenization which includes, different combinations of token can be used to relay the same prompts, which may lead to unfair pricing for the APIs of these LLMs. In a multilingual environment it may cause unexpected model responses due to different spacing in the prompts for languages such Taiwanese mandarin or Chinese mandarin [614].

C. Computational Requirements

Pre-training LLMs requires significant computational costs which can be very expensive, both financially and environmentally. Millions of dollars are spent in training these LLMs with thousands of compute hours and energy consumption. These are classified as Red AI [615], referring to models achieving state of the art results due to vast computation. Scaling these models can also be a challenging task due to the number of resources invested to train these LLMs. The concept of Computer Optimal Training [612] was introduced to address this problem for maximizing the training efficiency with respect to the corpus and model size.

D. Fine-Tuning LLMs

Fine-Tuning LLMs is a useful technique to train LLMs to custom tasks by training further on these task-specific datasets [48]. However, it requires a high amount of memory and large compute resources to store model gradients, parameters and activations, along with storing these fine-tuned models, limiting its access to a few institutions. Parameter-Efficient Fine-Tuning is a technique that can be used to address this problem which consists of updating a subset of model parameters such as prefix fine-tuning [616], prompt-tuning [617] and adapters [618]. Although techniques like Low-Rank Adaptation (LoRA) [619], or LongLora [620], or QLora [621] can be used to optimize the computation cost, but still computational demands remain a significant barrier for Fine-Tuning LLMs.

E. Inference Latency

High inference latency is one of the major challenges of LLMs which is mainly due to large memory footprints and lack of model parallelism. Several techniques can be used to mitigate this problem. Efficient Attention [622] can be used for accelerating attention through sub quadratic approximations such as multi-attention query or flash attention. Quantization [59] can be used to reduce the large memory footprint by reducing the computational precision of activations and weights. Pruning [623] and cascading [624] are some more techniques that can reduce the inference latency drastically for efficient and seamless responses.

F. Limited Context Length

Limited Context Length is one of the crucial aspects of LLMs, as it is extremely useful for interpretation of different prompts and semantic analysis. Without this contextual information, it can drastically degrade the performance of LLMs. There are several strategies that can be used to address

this; Positional Embedding Schemas [625], Efficient Attention [626] and Transformer alternatives. Different Positional Embedding Schemas can help LLMs to generalize well to different prompts which may not exist in the training data. Transient Global [627] and Luna [628] are some efficient attention mechanisms that can process larger context lengths effectively. Recurrent Neural Networks (RNNs) [629] and State Space Models (SSMs) [630] are good alternative for transformer-based approaches and are effective for addressing limiting context length.

G. Knowledge Updating and Refinement

Retraining the models is a costly process and is not sustainable. To address this, approaches such as model editing [631] is a technique which uses non-parametric knowledge resources to alter a model's behavior, and preserving model parameters by feeding new weights to modify the model's behavior can be used. However, these approaches are found to have limited generalizability and may only be applicable to a limited model architecture. On this end, web plugins and access to the web can alleviate the knowledge updating problem.

H. Training and Inference Cost

As LLMs involve billions of parameters, many powerful GPUs are needed for training and inference. According to the Artificial Intelligence Index Report 2024 published by Stanford University AI-index GPT-4 costs US78M and Gemini Ultra costs US191M to train (see Page 64 of the Stanford Report). This is far beyond the affordability of academic institutions and small and medium enterprises (SMEs). The inference latency can be reduced based on several methods, such as weight low-rank approximation, pruning, and quantization, as discussed above. This will effectively produce a small language model (SLM), which is more efficient for training and inference. A language model can also be simplified significantly by “converting textual inputs into cloze questions”, which leads to “three orders of magnitude fewer parameters” [632]. For visual recognition tasks, training can be speeded up for transformers based on the approximation of self-attention matrices using the Nyström approximation [633] and CUR decomposition [634]. In the future, more efficient training and inference algorithms are needed to make LLMs affordable for many researchers, developers, and users in real-world applications. Figure 13 shows a comparison between the input and output token cost, and context size, between benchmark LLMs models.

I. Risks of Foundation models

A foundation model refers to a base or core model that serves as the fundamental architecture for various machine learning tasks. In [138], a careful assessment of the risks and benefits of foundation models is done. A review [579] also highlights the potential threats and benefits of foundation models in health and education.

- **Bias:** Language models have the potential to unintentionally demonstrate bias when the training data used in

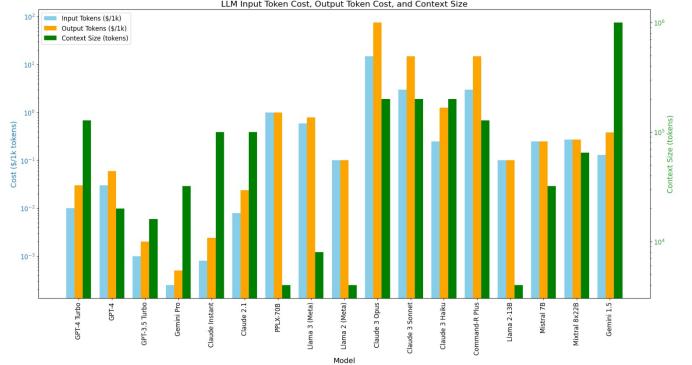


Fig. 13: State-of-the-art LLMs and their comparison in terms of token cost, output token cost, and context size.

their development is biased. According to Schramowski et al. [635], large pre-trained models designed to mimic natural languages can inadvertently perpetuate unfairness and prejudices.

- **Privacy** The manifestations of these biases are as follows:
 - (i) Training data bias: Language models typically rely on extensive datasets of human language for training [636], [637]. If these datasets contain biases related to factors such as race, gender, or socioeconomic status, the model may internalize and reproduce these biases in its responses.
 - (ii) User interaction bias: The responses generated by Chatbots are influenced by the input received from users. If users consistently pose biased or prejudiced questions, the model may learn and perpetuate these biases in its responses.
 - (iii) Algorithmic bias: Biases can also be introduced through the algorithms employed in training and operating language models and Chatbots [638].

- **Information Hallucination:** Hallucination in Natural Language Generation (NLG) is the generation of text that is nonsensical or unfaithful to the provided source content [639]. Hallucinations in LLMs are often the result of the model's attempt to fill in gaps in knowledge or context, with assumptions that are based on the patterns it has learned during training. This can lead to incorrect or misleading outputs, which can be particularly problematic in sensitive applications [475].

The cause of hallucinations in LLMs is an area of active research. Recent advances suggest that it's a complex problem related to the model's training process, dataset, and architectural design [640]. In particular, LLMs might be biased towards producing more "interesting" or fluent outputs, leading to a higher risk of hallucination [641]. There have been several proposed methods to mitigate the issue of hallucinations, one approach is to modify the training process to explicitly penalize hallucinations, such as in the case of "reality grounding" [642]. Another is to provide the model with a larger and more diverse dataset, which might reduce the risk of the model making incorrect assumptions [77]. In addition to this, researchers are exploring the use of "verifiable" or "fact-checkable" data during training, to teach the model to rely more

on facts and less on its own assumptions [643]. This, however, requires careful consideration of the data and metrics used.

- **Lack of Explainability:** No one can explain a model containing 175 billion parameters, the advent of LLMs has ushered in unprecedented advancements in NLP tasks. However, the sheer complexity and scale of these models present challenges in terms of explainability [644], [645]. As LLMs continue to grow in size, with models containing billions of parameters, the ability to comprehensively explain their decision-making processes becomes increasingly elusive [646], [647]. This complexity makes it exceedingly difficult for humans to understand and interpret the decision-making mechanisms employed by the model [85]. The lack of transparency [648] hinders the ability to gain insights into how specific inputs lead to particular outputs [649]. This, in addition to the intricate architecture of LLMs, often consisting of deep neural networks, exacerbates the challenge of explainability [650]. The numerous layers and complex interactions make it challenging to trace the reasoning process of the model. While techniques such as attention mechanisms [651] can provide some insights into the model's focus, they do not provide a comprehensive understanding of how the model arrives at its final output.

Finally, the lack of explainability in LLMs raises concerns regarding accountability, trust, and ethical considerations [652], [648].

- **Reasoning Errors:** LLM can make mistakes in logical reasoning [653], either because of ambiguities in the prompt or inherent limitations in its understanding of complex logical operations. LLMs cannot plan, reason, and have limited knowledge and commonsense [654] about the physical world [655]. From a cognitive science perspective, Auto-regressive LLMs at their best can approximate the Wernicke and Broca areas in the brain [656].
- **Prompt Injection, 'Jail Break' Attacks [657], Data Poisoning Attacks:** GPT-4 is susceptible to various adversarial attacks. For instance, a malicious actor might inject misleading prompts, perform 'jailbreak' attacks to make the model reveal sensitive information, or use data poisoning strategies to manipulate the model's output. Such vulnerabilities have been discussed in [658], [194] through experiments.

- **Adversarial Attacks:** Adversarial attacks on large language models (LLMs) are a type of security threat that can be used to manipulate or control the output of an LLM. These attacks work by deliberately introducing small changes to the input text, which the LLM then misinterprets and produces incorrect or harmful output [659]. One common type of adversarial attack is called a text injection attack. In this type of attack, the attacker introduces carefully crafted text into the input, which the LLM then interprets as a command. For example, the attacker could inject the text "delete all files" into an LLM that is used to control a computer system. The LLM would then

delete all of the files on the system [660]. Visual-prompt based models are also being attacked by these corrupted prompts [661].

- **Behavioral Changes over Time** Chen et. al. [662] investigated the performance of GPT 3.5 and GPT 4 over time, between March 2023 to June 2023, and found that the performance can greatly vary over time. For example, In March, GPT-4 had an accuracy of 84%, but in June, its accuracy dropped to 51%, a decrease of 33%. However, many experts suggest that the performance decrease is due to model drift [663] or prompt drift [664], we need to prompt better for maintaining the performance.
- **Spelling and Counting Errors:** Some specific tasks, like identifying and correcting spelling errors, can be challenging for GPT-4 due to its statistical nature. Another such example are counting errors. Counting error occurs when the model miscounts or misinterprets numerical quantities. For instance, it may provide incorrect calculations or misplaced decimal points when performing arithmetic operations, and counting the number of words or characters in long paragraphs [611], [665].

XII. OPEN QUESTIONS

In this section, we evaluate the open questions that are faced by AI researchers from technical, usage and philosophical standpoints.

A. Environmental and Energy Resources

Studies have revealed that the training process for GPT-3 alone used up 185,000 gallons of water, equivalent to what's needed to fill a cooling tower of a nuclear reactor [666]. This high consumption of water is primarily due to the cooling process of data centers, which necessitates a massive amount of water to regulate the servers' optimal temperature. Moreover, it is expected that the development of newer and advanced version models would need even more significant amounts of water due to their larger data parameters [244]. This concern has been discussed in [667], which presents a method to estimate the water footprint of AI language models and suggests more information transparency in this regard.

Apart from water usage, the training of LLMs demands a considerable amount of electricity. The training of OpenAI's GPT-3 alone resulted in the release of 502 metric tons of carbon, which could provide energy to an average American household for hundreds of years [668]. The amount of energy consumed by AI tools during training can be staggering, with some estimates suggesting that it can take hundreds of thousands or even millions of kWh to train a single large-scale model like GPT-3 [669], [670]. This electricity usage also contributes to indirect water consumption through power generation for data centers located off-site which should be taken into account, leading to carbon emissions [671].

The energy consumption of AI training has significant implications for the environment, particularly in terms of greenhouse gas emissions and climate change [672]. The energy required to train AI models is often generated from fossil fuels, such as coal and natural gas, which emit large

amounts of carbon dioxide and other greenhouse gases into the atmosphere. This can contribute to global warming and other environmental impacts [673]. As AI becomes more pervasive in our daily lives, it is important to consider the energy requirements of these systems and develop strategies to mitigate their impact on the environment [674]. One such solution is for data centers to adopt more eco-friendly cooling systems, such as using recycled water or implementing advanced cooling technologies [675]. Additionally, renewable energy sources, such as solar or wind power, can be utilized to power data centers, thereby reducing carbon emissions. Limiting the size and intricacy of LLMs is another potential solution, as smaller models require less data, resulting in reduced energy consumption [671]. Another study by Chien et. al [676] found that with models like ChatGPT, inference services dominated the power consumption and the power emissions for one year were equivalent to 25 times the training power of GPT3. They suggested the use of request direction approaches as a promising manner of reducing power consumption in LLMs.

B. Ethical Considerations

Inadvertently, LLMs may perpetuate biases inherent in the training data, resulting in outputs that are biased or discriminatory [677] as discussed previously. The challenge lies in identifying and mitigating such biases to ensure fair and equitable treatment across diverse user groups and disciplines [678]. Incorporating robust data authenticity and consent mechanisms, data anonymization techniques, and data retention policies into the development and deployment of LLMs can help ensure the responsible and ethical handling of user data. In [679], presents a comprehensive study of trustworthiness in LLMs, including different dimensions of trustworthiness, established benchmark, evaluation, and analysis.

1) Humans VS LLMs: Human interactions offer a deep level of empathy, emotional intelligence, and the ability to understand complex nuances in everyday life-situations. Humans responses are not only based on the current situation (prompt), but also considers other factors [200].

On the other hand, chatbots powered by AI have their advantages. They can operate 24/7, handle large volumes of inquiries simultaneously, and provide quick and consistent responses [680]. There is also a need to develop new performance metrics for measuring the intelligence of AI systems, as traditional methods of assessing intelligence, such as IQ tests [372], are not well-suited for AI systems, as they are designed to measure human intelligence [681].

2) Copy-right Issues: Training Large Language Models (LLMs) involves using vast amounts of textual data, which often includes copyrighted material [682]. This practice raises significant legal and ethical concerns regarding the unauthorized use of protected content. [682] examines the AI lifecycle from data collection to model deployment, emphasizing that most training materials are copyrighted, creating legal challenges.

3) Training data contamination from AI-generated content: Data sources for training large models are typically scraped from the internet. With the increasing popularity of generative

AI, it is possible that data present on the internet will have a significant component generated by AI models and therefore, reduce the human creativity aspect of the training data. Models, if trained on such data might end up trying to copy the generation aspects of previous AI models rather than humans only. One solution to this could be to use AI detection engines [683] that can determine content generated by AI before passing it through the model during the training process. There is a need to develop a dependable mechanism [587] to perform this task and retain the integrity of data.

4) The future as we perceive it: Large Language has vast potential for practical applications, particularly when combined with human oversight and judgement.

- **Use in Low Stakes Applications, Combine with Human Oversight:** LLMs are best suited for low stakes applications, where errors or inaccuracies can be tolerated. Moreover, combining LLMs with human oversight can significantly mitigate the risk of errors, biases, and other issues [684], [685].
- **Source of Inspiration, Suggestions:** LLMs can serve as an invaluable source of inspiration and suggestions, helping users brainstorm ideas [686], create content [687], and make decisions [688].
- **Copilots Over Autonomous Agents:** Given its limitations, LLMs are better suited as a 'copilot' that provides assistance and suggestions, rather than an autonomous agent that acts without human input or oversight [689], [690].
- **Artificial General Intelligence - AGI** Artificial general intelligence (AGI [691]) is a hypothetical type of artificial intelligence that would have the ability to learn and perform any intellectual task. In [367], GPT-4 is found to have sparks of artificial general intelligence. GPT-4 is able to perform a variety of tasks; such as solving math problems, writing creative contents, writing poems and poetry [692] and answering questions in an informative way. However, in our opinion, realizing the dream of AGI is still far away, despite of the rapid progress in the LLMs development. The key challenges include; understanding natural intelligence [693], developing adaptable fully autonomous models [694], and being safe and reliable with the understanding of the physical world [695], [696].
- **Embodied Artificial Intelligence** Embodied AI [697] bridges digital intelligence and physical action, enabling systems to perform complex tasks, interact naturally with humans, and adapt to dynamic environments, offering transformative potential across various industries. Vision-language models [698] are critical for enhancing the capabilities of Embodied AI by enabling systems to understand and interpret visual data [699] in conjunction with natural language. One such model is EmbodiedGPT [700], which integrates the powerful capabilities of GPT models with vision-language technology [701].
- **Small Language Models** Small language models [632], [702]...., Paradigm shift between accuracy and trade-off, small language models are the future — **need to be extended**

- **Democratizing AI** Democratizing AI [703] is a crucial movement that seeks to make artificial intelligence accessible and inclusive for a wide range of individuals and organizations. By breaking down barriers and providing user-friendly tools, democratization empowers diverse communities to leverage the power of AI to solve problems and drive innovation. It emphasizes the importance of open data, transparency, and accountability, ensuring that AI systems are unbiased, understandable, and ethically grounded.
- **Open-source LLMs** Open-source Large Language Models (LLMs) are AI models that are freely accessible to the public, allowing researchers and developers to study, modify, and improve them [662]. These models provide a transparent and collaborative environment that fosters innovation and rapid development in natural language processing [704]. Open-source LLMs, like GPT-Neo [705], BLOOM [706], and LLaMA [707], are crucial for democratizing AI technology, enabling a diverse community to contribute to advancements, and ensuring that AI development is not solely in the hands of large corporations.

No-code AI platforms [708] may also assist in democratizing AI initiatives [709], by providing a user-friendly interface that allows users to build and deploy ML models without any coding experience [710]. *No-code AI* can be used to leverage machine learning operations (MLOps) [711], to ensure models are deployed and managed effectively in production.

XIII. CONCLUSION

In this survey, we provided a comprehensive exploration of LLMs, their implications, technical concepts, and practical learning and usage. We discussed the potential benefits and risks of LLMs, and explored the different ways in which they can be used. We also provided a number of examples of how LLMs are being used in practice; such as generating images, chatting with pdf files, and also discussed GPT plug-ins. A comparison of popular chatbots; such as, ChatGPT, Bard, and Bing Chat is also provided. Benchmark dataset for LLM training and fine-tuning are also presented.

We particularly explored the applications of LLMs in medicine, engineering, agriculture, education, finance, media, law, and the entertainment industry. A list of popular LLM-based open-source applications for a variety of tasks is also presented. By delving into the technical intricacies, effective utilization, and future potential of LLMs, the survey will contribute to a deeper understanding and usage of these models within the research community. The survey has shed light on the key elements that drive the success of large language models through an examination of their working principles, diverse architectures, guidelines for prompting, AI-enabled tools and plug-ins, optimal strategies for employing LLMs, as well as advancements in pre-training, fine-tuning, and capability evaluation. A thorough comparison between popular chatbots has been provided as well.

Furthermore, the survey has also highlighted the importance of the safe and ethical use of AI tools like ChatGPT and

others. It recognizes the need for developing guidelines and regulations to address concerns related to security, ethics, the economy, and the environment. Ensuring the responsible integration of LLMs in healthcare, academia, and other industries is critical, as it enables these tools to effectively support and enhance human endeavors while upholding the values of integrity, privacy, and fairness. In our opinion, A technology X can replace a technology Y on a task Z , and can also help increase the productivity of humans on several tasks. LLMs have a great potential to transform many fields and bring positive impact on humans and society.

As the field of LLMs continues to evolve and progress, future research and development efforts should focus on improving the accuracy and performance of these models, addressing their limitations, and exploring new ways to use them. By adopting the guidelines presented in this survey, researchers and practitioners can contribute to the ongoing advancement of LLMs and ensure that they are used in a responsible and beneficial manner.

DECLARATION OF INTEREST

The authors have no conflicts of interest to declare.

REFERENCES

- [1] K. S. Jones, “Natural language processing: a historical review,” *Current issues in computational linguistics: in honour of Don Walker*, pp. 3–16, 1994.
- [2] K. Chowdhary and K. Chowdhary, “Natural language processing,” *Fundamentals of artificial intelligence*, pp. 603–649, 2020.
- [3] T. Iqbal and S. Qureshi, “The survey: Text generation models in deep learning,” *Journal of King Saud University-Computer and Information Sciences*, vol. 34, no. 6, pp. 2515–2528, 2022.
- [4] D. Nozza, F. Bianchi, D. Hovy, et al., “Honest: Measuring hurtful sentence completion in language models,” in *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Association for Computational Linguistics, 2021.
- [5] B. Min, H. Ross, E. Sulem, A. P. B. Veyseh, T. H. Nguyen, O. Sainz, E. Agirre, I. Heintz, and D. Roth, “Recent advances in natural language processing via large pre-trained language models: A survey,” *ACM Computing Surveys*, 2021.
- [6] M. Soam and S. Thakur, “Next word prediction using deep learning: A comparative study,” in *2022 12th International Conference on Cloud Computing, Data Science & Engineering (Confluence)*, pp. 653–658, IEEE, 2022.
- [7] S. Diao, R. Xu, H. Su, Y. Jiang, Y. Song, and T. Zhang, “Taming pre-trained language models with n-gram representations for low-resource domain adaptation,” in *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pp. 3336–3349, 2021.
- [8] P. F. Brown, V. J. Della Pietra, P. V. Desouza, J. C. Lai, and R. L. Mercer, “Class-based n-gram models of natural language,” *Computational linguistics*, vol. 18, no. 4, pp. 467–480, 1992.
- [9] N. Omar and Q. Al-Tashi, “Arabic nested noun compound extraction based on linguistic features and statistical measures,” *GEMA Online Journal of Language Studies*, vol. 18, no. 2, pp. 93–107, 2018.
- [10] B. Rawat, A. S. Bist, U. Rahardja, Q. Aini, and Y. P. A. Sanjaya, “Recent deep learning based nlp techniques for chatbot development: An exhaustive survey,” in *2022 10th International Conference on Cyber and IT Service Management (CITSM)*, pp. 1–4, IEEE, 2022.
- [11] Q. Lhoest, A. V. del Moral, Y. Jernite, A. Thakur, P. von Platen, S. Patil, J. Chaumond, M. Drame, J. Plu, L. Tunstall, et al., “Datasets: A community library for natural language processing,” *arXiv preprint arXiv:2109.02846*, 2021.
- [12] O. Sharir, B. Peleg, and Y. Shoham, “The cost of training nlp models: A concise overview,” *arXiv preprint arXiv:2004.08900*, 2020.

- [13] J. Wei, Y. Tay, R. Bommasani, C. Raffel, B. Zoph, S. Borgeaud, D. Yogatama, M. Bosma, D. Zhou, D. Metzler, *et al.*, “Emergent abilities of large language models,” *arXiv preprint arXiv:2206.07682*, 2022.
- [14] A. Srivastava, A. Rastogi, A. Rao, A. A. M. Shoeb, A. Abid, A. Fisch, A. R. Brown, A. Santoro, A. Gupta, A. Garriga-Alonso, *et al.*, “Beyond the imitation game: Quantifying and extrapolating the capabilities of language models,” *arXiv preprint arXiv:2206.04615*, 2022.
- [15] H. Touvron, T. Lavigra, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, *et al.*, “Llama: Open and efficient foundation language models,” *arXiv preprint arXiv:2302.13971*, 2023.
- [16] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [17] D. Luitse and W. Denkera, “The great transformer: Examining the role of large language models in the political economy of ai,” *Big Data & Society*, vol. 8, no. 2, p. 20539517211047734, 2021.
- [18] Z. Dong, T. Tang, L. Li, and W. X. Zhao, “A survey on long text modeling with transformers,” *arXiv preprint arXiv:2302.14502*, 2023.
- [19] K. Adnan and R. Akbar, “An analytical study of information extraction from unstructured and multidimensional big data,” *Journal of Big Data*, vol. 6, no. 1, pp. 1–38, 2019.
- [20] M. Awais, M. Naseer, S. Khan, R. M. Anwer, H. Cholakkal, M. Shah, M.-H. Yang, and F. S. Khan, “Foundational models defining a new era in vision: A survey and outlook,” *arXiv preprint arXiv:2307.13721*, 2023.
- [21] H. Zhang, X. Li, and L. Bing, “Video-llama: An instruction-tuned audio-visual language model for video understanding,” *arXiv preprint arXiv:2306.02858*, 2023.
- [22] A. Rouditchenko, A. Boggust, D. Harwath, B. Chen, D. Joshi, S. Thomas, K. Audhkhasi, H. Kuehne, R. Panda, R. Feris, *et al.*, “Avl-net: Learning audio-visual language representations from instructional videos,” *arXiv preprint arXiv:2006.09199*, 2020.
- [23] Y. Zhao, Z. Lin, D. Zhou, Z. Huang, J. Feng, and B. Kang, “Bubogpt: Enabling visual grounding in multi-modal llms,” *arXiv preprint arXiv:2307.08581*, 2023.
- [24] J. Huang and K. C.-C. Chang, “Towards reasoning in large language models: A survey,” *arXiv preprint arXiv:2212.10403*, 2022.
- [25] N. Pappas and T. Meyer, “A survey on language modeling using neural networks,” tech. rep., Idiap, 2012.
- [26] J. R. Bellegarda, “Statistical language model adaptation: review and perspectives,” *Speech communication*, vol. 42, no. 1, pp. 93–108, 2004.
- [27] J. Lafferty and C. Zhai, “Probabilistic relevance models based on document and query generation,” *Language modeling for information retrieval*, pp. 1–10, 2003.
- [28] V. A. Petrushin, “Hidden markov models: Fundamentals and applications,” in *Online Symposium for Electronics Engineer*, 2000.
- [29] S. Khudanpur and J. Wu, “Maximum entropy techniques for exploiting syntactic, semantic and collocational dependencies in language modeling,” *Computer Speech & Language*, vol. 14, no. 4, pp. 355–372, 2000.
- [30] H. Wang, J. He, X. Zhang, and S. Liu, “A short text classification method based on n-gram and cnn,” *Chinese Journal of Electronics*, vol. 29, no. 2, pp. 248–254, 2020.
- [31] R. Rosenfeld, “Two decades of statistical language modeling: Where do we go from here?,” *Proceedings of the IEEE*, vol. 88, no. 8, pp. 1270–1278, 2000.
- [32] E. Arisoy, T. N. Sainath, B. Kingsbury, and B. Ramabhadran, “Deep neural network language models,” in *Proceedings of the NAACL-HLT 2012 Workshop: Will We Ever Really Replace the N-gram Model? On the Future of Language Modeling for HLT*, pp. 20–28, 2012.
- [33] J. R. Bellegarda, “Exploiting latent semantic information in statistical language modeling,” *Proceedings of the IEEE*, vol. 88, no. 8, pp. 1279–1296, 2000.
- [34] F. Alva-Manchego, C. Scarton, and L. Specia, “Data-driven sentence simplification: Survey and benchmark,” *Computational Linguistics*, vol. 46, no. 1, pp. 135–187, 2020.
- [35] M. Malik, M. K. Malik, K. Mehmood, and I. Makhdoom, “Automatic speech recognition: a survey,” *Multimedia Tools and Applications*, vol. 80, pp. 9411–9457, 2021.
- [36] J. Cervantes, F. Garcia-Lamont, L. Rodríguez-Mazahua, and A. Lopez, “A comprehensive survey on support vector machine classification: Applications, challenges and trends,” *Neurocomputing*, vol. 408, pp. 189–215, 2020.
- [37] M. Crawford, T. M. Khoshgoftaar, J. D. Prusa, A. N. Richter, and H. Al Najada, “Survey of review spam detection using machine learning techniques,” *Journal of Big Data*, vol. 2, no. 1, pp. 1–24, 2015.
- [38] M. Neethu and R. Rajasree, “Sentiment analysis in twitter using machine learning techniques,” in *2013 fourth international conference on computing, communications and networking technologies (ICCCNT)*, pp. 1–5, IEEE, 2013.
- [39] A. Go, L. Huang, R. Bhayani, *et al.*, “Twitter sentiment analysis,” *Entropy*, vol. 17, p. 252, 2009.
- [40] L. Deng and Y. Liu, “A joint introduction to natural language processing and to deep learning,” *Deep learning in natural language processing*, pp. 1–22, 2018.
- [41] W. Yin, K. Kann, M. Yu, and H. Schütze, “Comparative study of cnn and rnn for natural language processing,” *arXiv preprint arXiv:1702.01923*, 2017.
- [42] T. Mikolov, M. Karafiat, L. Burget, J. Cernocky, and S. Khudanpur, “Recurrent neural network based language model,” in *Interspeech*, vol. 2, pp. 1045–1048, Makuhari, 2010.
- [43] S. Hochreiter, “Recurrent neural net learning and vanishing gradient,” *International Journal Of Uncertainty, Fuzziness and Knowledge-Based Systems*, vol. 6, no. 2, pp. 107–116, 1998.
- [44] S. Hihi and Y. Bengio, “Hierarchical recurrent neural networks for long-term dependencies,” *Advances in neural information processing systems*, vol. 8, 1995.
- [45] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” *Advances in neural information processing systems*, vol. 30, 2017.
- [46] P. Shaw, J. Uszkoreit, and A. Vaswani, “Self-attention with relative position representations,” *arXiv preprint arXiv:1803.02155*, 2018.
- [47] B. Ghoghoj and A. Ghodsi, “Attention mechanism, transformers, bert, and gpt: tutorial and survey,” 2020.
- [48] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” *arXiv preprint arXiv:1810.04805*, 2018.
- [49] Q. Liu, M. J. Kusner, and P. Blunsom, “A survey on contextual embeddings,” *arXiv preprint arXiv:2003.07278*, 2020.
- [50] E. Adamopoulou and L. Moussiades, “Chatbots: History, technology, and applications,” *Machine Learning with Applications*, vol. 2, p. 100006, 2020.
- [51] M. Allahyari, S. Pouriyeh, M. Assefi, S. Safaei, E. D. Trippe, J. B. Gutierrez, and K. Kochut, “Text summarization techniques: a brief survey,” *arXiv preprint arXiv:1707.02268*, 2017.
- [52] Y. Ge, W. Hua, J. Ji, J. Tan, S. Xu, and Y. Zhang, “Openagi: When llm meets domain experts,” *arXiv preprint arXiv:2304.04370*, 2023.
- [53] K. I. Roumeliotis, N. D. Tselikas, and D. K. Nasiopoulos, “Llama 2: Early adopters’ utilization of meta’s new open-source pretrained model,” 2023.
- [54] A. Baladn, I. Sastre, L. Chiruzzo, and A. Ros, “Retuyt-inco at bea 2023 shared task: Tuning open-source llms for generating teacher responses,” in *Proceedings of the 18th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2023)*, pp. 756–765, 2023.
- [55] J. J. Nay, “Large language models as fiduciaries: A case study toward robustly communicating with artificial intelligence through legal standards,” 2023.
- [56] R. Sennrich, B. Haddow, and A. Birch, “Neural machine translation of rare words with subword units,” *arXiv preprint arXiv:1508.07909*, 2015.
- [57] T. Y. Zhuo, Z. Li, Y. Huang, Y.-F. Li, W. Wang, G. Haffari, and F. Shiri, “On robustness of prompt-based semantic parsing with large pre-trained language model: An empirical study on codex,” *arXiv preprint arXiv:2301.12868*, 2023.
- [58] W.-L. Chiang, Z. Li, Z. Lin, Y. Sheng, Z. Wu, H. Zhang, L. Zheng, S. Zhuang, Y. Zhuang, J. E. Gonzalez, I. Stoica, and E. P. Xing, “Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality,” March 2023.
- [59] Z. Yao, R. Yazdani Aminabadi, M. Zhang, X. Wu, C. Li, and Y. He, “Zeroquant: Efficient and affordable post-training quantization for large-scale transformers,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 27168–27183, 2022.
- [60] A. Zou, Z. Wang, J. Z. Kolter, and M. Fredrikson, “Universal and transferable adversarial attacks on aligned language models,” *arXiv preprint arXiv:2307.15043*, 2023.
- [61] D. M. Katz, M. J. Bommarito, S. Gao, and P. Arredondo, “GPT-4 Passes the Bar Exam,” March 2023.
- [62] A. Byrd, “Truth-telling: Critical inquiries on llms and the corpus texts that train them,” *Composition Studies*, vol. 51, no. 1, pp. 135–142, 2023.

- [63] X. Zhang, B. Yu, H. Yu, Y. Lv, T. Liu, F. Huang, H. Xu, and Y. Li, “Wider and deeper llm networks are fairer llm evaluators,” *arXiv preprint arXiv:2308.01862*, 2023.
- [64] I. Yildirim and L. Paul, “From task structures to world models: What do llms know?,” *arXiv preprint arXiv:2310.04276*, 2023.
- [65] H. Jin, X. Han, J. Yang, Z. Jiang, C.-Y. Chang, and X. Hu, “Growlength: Accelerating llms pretraining by progressively growing training length,” *arXiv preprint arXiv:2310.00576*, 2023.
- [66] R. V. P. Marcel, B. E. M. Fernando, and Y. V. J. Roberto, “A brief history of the artificial intelligence: chatgpt: The evolution of gpt,” in *2023 18th Iberian Conference on Information Systems and Technologies (CISTI)*, pp. 1–5, IEEE, 2023.
- [67] E. Y. Chang, “Examining gpt-4: Capabilities, implications, and future directions,” 2023.
- [68] M. Zhang and J. Li, “A commentary of gpt-3 in mit technology review 2021,” *Fundamental Research*, vol. 1, no. 6, pp. 831–833, 2021.
- [69] Z. Lan, M. Chen, S. Goodman, K. Gimpel, P. Sharma, and R. Soricut, “Albert: A lite bert for self-supervised learning of language representations,” *arXiv preprint arXiv:1909.11942*, 2019.
- [70] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, “Roberta: A robustly optimized bert pretraining approach,” *arXiv preprint arXiv:1907.11692*, 2019.
- [71] A. J. Thirunavukarasu, D. S. J. Ting, K. Elangovan, L. Gutierrez, T. F. Tan, and D. S. W. Ting, “Large language models in medicine,” *Nature Medicine*, pp. 1–11, 2023.
- [72] H. Wang, T. Fu, Y. Du, W. Gao, K. Huang, Z. Liu, P. Chandak, S. Liu, P. Van Katwyk, A. Deac, et al., “Scientific discovery in the age of artificial intelligence,” *Nature*, vol. 620, no. 7972, pp. 47–60, 2023.
- [73] J. Wang, Y. Huang, C. Chen, Z. Liu, S. Wang, and Q. Wang, “Software testing with large language model: Survey, landscape, and vision,” *arXiv preprint arXiv:2307.07221*, 2023.
- [74] F. F. Xu, U. Alon, G. Neubig, and V. J. Hellendoorn, “A systematic evaluation of large language models of code,” in *Proceedings of the 6th ACM SIGPLAN International Symposium on Machine Programming*, pp. 1–10, 2022.
- [75] J. Cabrera, M. S. Loyola, I. Magaña, and R. Rojas, “Ethical dilemmas, mental health, artificial intelligence, and llm-based chatbots,” in *International Work-Conference on Bioinformatics and Biomedical Engineering*, pp. 313–326, Springer, 2023.
- [76] A. Creswell, M. Shanahan, and I. Higgins, “Selection-inference: Exploiting large language models for interpretable logical reasoning,” *arXiv preprint arXiv:2205.09712*, 2022.
- [77] E. Ferrara, “Should chatgpt be biased? challenges and risks of bias in large language models,” *arXiv preprint arXiv:2304.03738*, 2023.
- [78] K. Tirumala, D. Simig, A. Aghajanyan, and A. S. Morcos, “D4: Improving llm pretraining via document de-duplication and diversification,” *arXiv preprint arXiv:2308.12284*, 2023.
- [79] J. White, Q. Fu, S. Hays, M. Sandborn, C. Olea, H. Gilbert, A. El-nashar, J. Spencer-Smith, and D. C. Schmidt, “A prompt pattern catalog to enhance prompt engineering with chatgpt,” *arXiv preprint arXiv:2302.11382*, 2023.
- [80] P. P. Ray, “Chatgpt: A comprehensive review on background, applications, key challenges, bias, ethics, limitations and future scope,” *Internet of Things and Cyber-Physical Systems*, 2023.
- [81] A. Sudmann, “On the media-political dimension of artificial intelligence: Deep learning as a black box and openai,” *Digital Culture & Society*, vol. 4, no. 1, pp. 181–200, 2018.
- [82] A. Koubaa, “Gpt-4 vs. gpt-3.5: A concise showdown,” 2023.
- [83] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, et al., “Training language models to follow instructions with human feedback,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 27730–27744, 2022.
- [84] S. Huang, L. Dong, W. Wang, Y. Hao, S. Singhal, S. Ma, T. Lv, L. Cui, O. K. Mohammed, Q. Liu, et al., “Language is not all you need: Aligning perception with language models,” *arXiv preprint arXiv:2302.14045*, 2023.
- [85] W. X. Zhao, K. Zhou, J. Li, T. Tang, X. Wang, Y. Hou, Y. Min, B. Zhang, J. Zhang, Z. Dong, et al., “A survey of large language models,” *arXiv preprint arXiv:2303.18223*, 2023.
- [86] Y. Du, Z. Liu, J. Li, and W. X. Zhao, “A survey of vision-language pre-trained models,” *arXiv preprint arXiv:2202.10936*, year=2022.
- [87] G. Mialon, R. Dessì, M. Lomeli, C. Nalmpantis, R. Pasunuru, R. Raileanu, B. Rozière, T. Schick, J. Dwivedi-Yu, A. Celikyilmaz, et al., “Augmented language models: a survey,” *arXiv preprint arXiv:2302.07842*, 2023.
- [88] R. Qureshi, M. Irfan, H. Ali, A. Khan, A. S. Nittala, S. Ali, A. Shah, T. M. Gondal, F. Sadak, Z. Shah, et al., “Artificial intelligence and biosensors in healthcare and its clinical relevance: A review,” *IEEE Access*, 2023.
- [89] Q. Al-Tashi, M. B. Saad, A. Sheshadri, C. C. Wu, J. Y. Chang, B. Al-Lazikani, C. Gibbons, N. I. Vokes, J. Zhang, J. J. Lee, et al., “Swarmdeepsurv: swarm intelligence advances deep survival network for prognostic radiomics signatures in four solid cancers,” *Patterns*.
- [90] S. Mohamadi, G. Mujtaba, N. Le, G. Doretto, and D. A. Adjeroh, “Chatgpt in the age of generative ai and large language models: A concise survey,” *arXiv preprint arXiv:2307.04251*, 2023.
- [91] H. Naveed, A. U. Khan, S. Qiu, M. Saqib, S. Anwar, M. Usman, N. Barnes, and A. Mian, “A comprehensive overview of large language models,” *arXiv preprint arXiv:2307.06435*, 2023.
- [92] E. Kasneci, K. Seßler, S. Küchemann, M. Bannert, D. Dementieva, F. Fischer, U. Gasser, G. Groh, S. Günemann, E. Hüllermeier, et al., “Chatgpt for good? on opportunities and challenges of large language models for education,” *Learning and Individual Differences*, vol. 103, p. 102274, 2023.
- [93] M. Sallam, “The utility of chatgpt as an example of large language models in healthcare education, research and practice: Systematic review on the future perspectives and potential limitations,” *medRxiv*, pp. 2023–02, 2023.
- [94] Z. Lin, H. Akin, R. Rao, B. Hie, Z. Zhu, W. Lu, A. dos Santos Costa, M. Fazel-Zarandi, T. Sercu, S. Candido, et al., “Language models of protein sequences at the scale of evolution enable accurate structure prediction,” *BioRxiv*, vol. 2022, p. 500902, 2022.
- [95] A. Madani, B. McCann, N. Naik, N. S. Keskar, N. Anand, R. R. Eguchi, P.-S. Huang, and R. Socher, “Progen: Language modeling for protein generation,” *arXiv preprint arXiv:2004.03497*, 2020.
- [96] Y. Cao, S. Li, Y. Liu, Z. Yan, Y. Dai, P. S. Yu, and L. Sun, “A comprehensive survey of ai-generated content (aigc): A history of generative ai from gan to chatgpt,” *arXiv preprint arXiv:2303.04226*, 2023.
- [97] I. Beltagy, K. Lo, and A. Cohan, “Scibert: A pretrained language model for scientific text,” *arXiv preprint arXiv:1903.10676*, 2019.
- [98] J. Li, T. Tang, W. X. Zhao, J.-Y. Nie, and J.-R. Wen, “Pretrained language models for text generation: A survey,” *arXiv preprint arXiv:2201.05273*, 2022.
- [99] S. Wu, O. Irsay, S. Lu, V. Dabrowski, M. Dredze, S. Gehrmann, P. Kambadur, D. Rosenberg, and G. Mann, “Bloomberggpt: A large language model for finance,” *arXiv preprint arXiv:2303.17564*, 2023.
- [100] T. Eloundou, S. Manning, P. Mishkin, and D. Rock, “Gpts are gpts: An early look at the labor market impact potential of large language models,” *arXiv preprint arXiv:2303.10130*, 2023.
- [101] B. Li, K. Mellou, B. Zhang, J. Pathuri, and I. Menache, “Large language models for supply chain optimization,” *arXiv preprint arXiv:2307.03875*, 2023.
- [102] L. Baria, Q. Zhao, H. Zou, Y. Tian, F. Bader, and M. Debbah, “Large language models for telecom: The next big thing?,” *arXiv preprint arXiv:2306.10249*, 2023.
- [103] M. Chen, J. Tworek, H. Jun, Q. Yuan, et al., “Evaluating large language models trained on code,” *arXiv preprint arXiv:2107.03374*, 2021.
- [104] S. Salman, J. A. Shamsi, and R. Qureshi, “Deep fake generation and detection: Issues, challenges, and solutions,” *IT Professional*, vol. 25, no. 1, pp. 52–59, 2023.
- [105] Z. Sun, “A short survey of viewing large language models in legal aspect,” *arXiv preprint arXiv:2303.09136*, year=2023.
- [106] R. Qureshi, M. Irfan, T. M. Gondal, S. Khan, J. Wu, M. U. Hadi, J. Heymach, X. Le, H. Yan, and T. Alam, “Ai in drug discovery and its clinical relevance,” *Heliyon*, 2023.
- [107] O. B. Shoham and N. Rappoport, “Cpllm: Clinical prediction with large language models,” *arXiv preprint arXiv:2309.11295*, 2023.
- [108] Q. Al-Tashi, M. B. Saad, A. Muneer, R. Qureshi, S. Mirjalili, A. Sheshadri, X. Le, N. I. Vokes, J. Zhang, and J. Wu, “Machine learning models for the identification of prognostic and predictive cancer biomarkers: A systematic review,” *International journal of molecular sciences*, vol. 24, no. 9, p. 7781, 2023.
- [109] A. Holzinger, K. Keiblinger, P. Holub, K. Zatloukal, and H. Müller, “Ai for life: Trends in artificial intelligence for biotechnology,” *New Biotechnology*, vol. 74, pp. 16–24, 2023.
- [110] L. Wang, C. Ma, X. Feng, Z. Zhang, H. Yang, J. Zhang, Z. Chen, J. Tang, X. Chen, Y. Lin, et al., “A survey on large language model based autonomous agents,” *arXiv preprint arXiv:2308.11432*, 2023.
- [111] Y. Zhu, X. Wang, J. Chen, S. Qiao, Y. Ou, Y. Yao, S. Deng, H. Chen, and N. Zhang, “Llms for knowledge graph construction and

- reasoning: Recent capabilities and future opportunities,” *arXiv preprint arXiv:2305.13168*, 2023.
- [112] L. Huynh, J. Hong, A. Mian, H. Suzuki, Y. Wu, and S. Camtepe, “Quantum-inspired machine learning: a survey,” *arXiv preprint arXiv:2308.11269*, 2023.
- [113] E. Brynjolfsson, D. Li, and L. R. Raymond, “Generative ai at work,” tech. rep., National Bureau of Economic Research, 2023.
- [114] P. Samuelson, “Generative ai meets copyright,” *Science*, vol. 381, no. 6654, pp. 158–161, 2023.
- [115] I. Chiang, *Unleashing the Power of Generative AI: The Race for Advancement and the Global Ramifications*. PhD thesis, Massachusetts Institute of Technology, 2023.
- [116] S. Wang, S. Menon, T. Long, K. Henderson, D. Li, K. Crowston, M. Hansen, J. V. Nickerson, and L. B. Chilton, “Reelframer: Co-creating news reels on social media with generative ai,” *arXiv preprint arXiv:2304.09653*, 2023.
- [117] S. Mayahi and M. Vidrih, “The impact of generative ai on the future of visual content marketing,” *arXiv preprint arXiv:2211.12660*, 2022.
- [118] S.-C. Chen, “Multimedia research toward the metaverse,” *IEEE Multi-Media*, vol. 29, no. 1, pp. 125–127, 2022.
- [119] A. Zentner, “Applied innovation: Artificial intelligence in higher education,” *Available at SSRN 4314180*, 2022.
- [120] J. Sun, Q. V. Liao, M. Muller, M. Agarwal, S. Houde, K. Talamadupula, and J. D. Weisz, “Investigating explainability of generative ai for code through scenario-based design,” in *27th International Conference on Intelligent User Interfaces*, pp. 212–228, 2022.
- [121] J. Morley, N. J. DeVito, and J. Zhang, “Generative ai for medical research,” 2023.
- [122] P. Ghimire, K. Kim, and M. Acharya, “Generative ai in the construction industry: Opportunities & challenges,” *arXiv preprint arXiv:2310.04427*, 2023.
- [123] H. Cui, C. Wang, H. Maan, K. Pang, F. Luo, and B. Wang, “scgpt: Towards building a foundation model for single-cell multi-omics using generative ai,” *bioRxiv*, pp. 2023–04, 2023.
- [124] S. B. Kotsiantis, I. D. Zaharakis, and P. E. Pintelas, “Machine learning: a review of classification and combining techniques,” *Artificial Intelligence Review*, vol. 26, pp. 159–190, 2006.
- [125] A. Pérez-Suárez, J. F. Martínez-Trinidad, and J. A. Carrasco-Ochoa, “A review of conceptual clustering algorithms,” *Artificial Intelligence Review*, vol. 52, pp. 1267–1296, 2019.
- [126] Z.-Q. Zhao, P. Zheng, S.-t. Xu, and X. Wu, “Object detection with deep learning: A review,” *IEEE transactions on neural networks and learning systems*, vol. 30, no. 11, pp. 3212–3232, 2019.
- [127] C. Chen, C. Qin, H. Qiu, G. Tarroni, J. Duan, W. Bai, and D. Rueckert, “Deep learning for cardiac image segmentation: a review,” *Frontiers in Cardiovascular Medicine*, vol. 7, p. 25, 2020.
- [128] A. A. de Hond, A. M. Leeuwenberg, L. Hooft, I. M. Kant, S. W. Nijman, H. J. van Os, J. J. Aardoom, T. P. Debray, E. Schuit, M. van Smeden, et al., “Guidelines and quality criteria for artificial intelligence-based prediction models in healthcare: a scoping review,” *NPJ digital medicine*, vol. 5, no. 1, p. 2, 2022.
- [129] C. Zhang, C. Zhang, S. Zheng, Y. Qiao, C. Li, M. Zhang, S. K. Dam, C. M. Thwal, Y. L. Tun, L. L. Huy, et al., “A complete survey on generative ai (aige): Is chatgpt from gpt-4 to gpt-5 all you need?,” *arXiv preprint arXiv:2303.11717*, 2023.
- [130] C. Zhang, C. Zhang, S. Zheng, M. Zhang, M. Qamar, S.-H. Bae, and I. S. Kweon, “A survey on audio diffusion models: Text to speech synthesis and enhancement in generative ai,” *arXiv preprint arXiv:2303.13336*, vol. 2, 2023.
- [131] L. Wang, W. Chen, W. Yang, F. Bi, and F. R. Yu, “A state-of-the-art review on image synthesis with generative adversarial networks,” *IEEE Access*, vol. 8, pp. 63514–63537, 2020.
- [132] N. Alidausari, A. Sowmya, N. Marcus, and G. Mohammadi, “Video generative adversarial networks: a review,” *ACM Computing Surveys (CSUR)*, vol. 55, no. 2, pp. 1–25, 2022.
- [133] S. Barke, M. B. James, and N. Polikarpova, “Grounded copilot: How programmers interact with code-generating models,” *Proceedings of the ACM on Programming Languages*, vol. 7, no. OOPSLA1, pp. 85–111, 2023.
- [134] D. Zhang, S. Li, X. Zhang, J. Zhan, P. Wang, Y. Zhou, and X. Qiu, “Speechgpt: Empowering large language models with intrinsic cross-modal conversational abilities,” *arXiv preprint arXiv:2305.11000*, 2023.
- [135] S. Hong, J. Seo, S. Hong, H. Shin, and S. Kim, “Large language models are frame-level directors for zero-shot text-to-video generation,” *arXiv preprint arXiv:2305.14330*, 2023.
- [136] Ö. AYDIN and E. KARAARSLAN, “Is chatgpt leading generative ai? what is beyond expectations,” *What is Beyond Expectations*, 2023.
- [137] B. Kim, H. Kim, S.-W. Lee, G. Lee, D. Kwak, D. H. Jeon, S. Park, S. Kim, S. Kim, D. Seo, et al., “What changes can large-scale language models bring? intensive study on hyperclova: Billions-scale korean generative pretrained transformers,” *arXiv preprint arXiv:2109.04650*, 2021.
- [138] R. Bommasani, D. A. Hudson, E. Adeli, R. Altman, S. Arora, S. von Arx, M. S. Bernstein, J. Bohg, A. Bosselut, E. Brunskill, et al., “On the opportunities and risks of foundation models,” *arXiv preprint arXiv:2108.07258*, 2021.
- [139] Y. Yuan, “On the power of foundation models.” in *International Conference on Machine Learning*, pp. 40519–40530, PMLR, 2023.
- [140] V. Sanh, L. Debut, J. Chaumond, and T. Wolf, “Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter,” *arXiv preprint arXiv:1910.01108*, 2019.
- [141] A. Jo, “The promise and peril of generative ai,” *Nature*, vol. 614, no. 1, pp. 214–216, 2023.
- [142] H. Bansal, K. Gopalakrishnan, S. Dingliwal, S. Bodapati, K. Kirchhoff, and D. Roth, “Rethinking the role of scale for in-context learning: An interpretability-based case study at 66 billion scale,” *arXiv preprint arXiv:2212.09095*, 2022.
- [143] M. Mariani, “Generative artificial intelligence and innovation: Conceptual foundations,” *Available at SSRN 4249382*, 2022.
- [144] W. Zeng, X. Ren, T. Su, H. Wang, Y. Liao, Z. Wang, X. Jiang, Z. Yang, K. Wang, X. Zhang, et al., “Pangu-alpha: Large-scale autoregressive pretrained chinese language models with auto-parallel computation,” *arXiv preprint arXiv:2104.12369*, 2021.
- [145] L. Mescheder, S. Nowozin, and A. Geiger, “Adversarial variational bayes: Unifying variational autoencoders and generative adversarial networks,” in *International conference on machine learning*, pp. 2391–2400, PMLR, 2017.
- [146] A. Creswell, T. White, V. Dumoulin, K. Arulkumaran, B. Bengio, and A. A. Bharath, “Generative adversarial networks: An overview,” *IEEE signal processing magazine*, vol. 35, no. 1, pp. 53–65, 2018.
- [147] F.-A. Croitoru, V. Hundru, R. T. Ionescu, and M. Shah, “Diffusion models in vision: A survey,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- [148] Y. Bai, S. Kadavath, S. Kundu, A. Askell, et al., “Training a helpful and harmless assistant with reinforcement learning from human feedback,” in *arXiv preprint arXiv:2204.05862*, 2022.
- [149] A. Muneer and S. M. Fati, “A comparative analysis of machine learning techniques for cyberbullying detection on twitter,” *Future Internet*, vol. 12, no. 11, p. 187, 2020.
- [150] H. Hassani and E. S. Silva, “The role of chatgpt in data science: how ai-assisted conversational interfaces are revolutionizing the field,” *Big data and cognitive computing*, vol. 7, no. 2, p. 62, 2023.
- [151] K. Wang, C. Gou, Y. Duan, Y. Lin, X. Zheng, and F.-Y. Wang, “Generative adversarial networks: introduction and outlook,” *IEEE/CAA Journal of Automatica Sinica*, vol. 4, no. 4, pp. 588–598, 2017.
- [152] J. N. Kather, N. Ghaffari Laleh, S. Foersch, and D. Truhn, “Medical domain knowledge in domain-agnostic generative ai,” *NPJ digital medicine*, vol. 5, no. 1, p. 90, 2022.
- [153] I. Solaiman, “The gradient of generative ai release: Methods and considerations,” in *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*, pp. 111–122, 2023.
- [154] M. Shoman, T. Ghoul, G. Lanzaro, T. Alsharif, S. Gargoum, and T. Sayed, “Enforcing traffic safety: A deep learning approach for detecting motorcyclists’ helmet violations using yolov8 and deep convolutional generative adversarial network-generated images,” *Algorithms*, vol. 17, no. 5, 2024.
- [155] M. Shoman, D. Wang, A. Aboah, and M. Abdel-Aty, “Enhancing traffic safety with parallel dense video captioning for end-to-end event analysis,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, pp. 7125–7133, June 2024.
- [156] S. Tan, Y. Shen, and B. Zhou, “Improving the fairness of deep generative models without retraining,” *arXiv preprint arXiv:2012.04842*, 2020.
- [157] K. Wach, C. D. Duong, J. Ejdys, R. Kazlauskaitė, P. Korzynski, G. Mazurek, J. Palisziewicz, and E. Ziembia, “The dark side of generative artificial intelligence: A critical analysis of controversies and risks of chatgpt,” *Entrepreneurial Business and Economics Review*, vol. 11, no. 2, pp. 7–24, 2023.
- [158] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative adversarial net-

- works,” *Communications of the ACM*, vol. 63, no. 11, pp. 139–144, 2020.
- [159] Z. Che, Y. Cheng, S. Zhai, Z. Sun, and Y. Liu, “Boosting deep learning risk prediction with generative adversarial networks for electronic health records,” in *2017 IEEE International Conference on Data Mining (ICDM)*, pp. 787–792, IEEE, 2017.
- [160] A. Shafahi, M. Najibi, M. A. Ghiasi, Z. Xu, J. Dickerson, C. Studer, L. S. Davis, G. Taylor, and T. Goldstein, “Adversarial training for free!” *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [161] S. Mukherjee, H. Asnani, E. Lin, and S. Kannan, “Clustergan: Latent space clustering in generative adversarial networks,” in *Proceedings of the AAAI conference on artificial intelligence*, vol. 33, pp. 4610–4617, 2019.
- [162] M. Morrison, R. Kumar, K. Kumar, P. Seetharaman, A. Courville, and Y. Bengio, “Chunked autoregressive gan for conditional waveform synthesis,” *arXiv preprint arXiv:2110.10139*, 2021.
- [163] S. Kaushik, A. Choudhury, S. Natarajan, L. A. Pickett, and V. Dutt, “Medicine expenditure prediction via a variance-based generative adversarial network,” *IEEE Access*, vol. 8, pp. 110947–110958, 2020.
- [164] L.-C. Yang and A. Lerch, “On the evaluation of generative models in music,” *Neural Computing and Applications*, vol. 32, no. 9, pp. 4773–4784, 2020.
- [165] N. Geneva and N. Zabaras, “Multi-fidelity generative deep learning turbulent flows,” *arXiv preprint arXiv:2006.04731*, 2020.
- [166] J. D. Cohen, S. M. McClure, and A. J. Yu, “Should i stay or should i go? how the human brain manages the trade-off between exploitation and exploration,” *Philosophical Transactions of the Royal Society B: Biological Sciences*, vol. 362, no. 1481, pp. 933–942, 2007.
- [167] D. Xu, F. Zhu, Q. Liu, and P. Zhao, “Improving exploration efficiency of deep reinforcement learning through samples produced by generative model,” *Expert Systems with Applications*, vol. 185, p. 115680, 2021.
- [168] A. Aboah, M. Shoman, V. Mandal, S. Davami, Y. Adu-Gyamfi, and A. Sharma, “A vision-based system for traffic anomaly detection using deep learning and decision trees,” in *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pp. 4202–4207, 2021.
- [169] P. Dhoni, “Exploring the synergy between generative ai, data and analytics in the modern age,” 2023.
- [170] G. Vigliensoni, P. Perry, R. Fiebrink, et al., “A small-data mindset for generative ai creative work,” 2022.
- [171] M. Shoman, A. Aboah, A. Morehead, Y. Duan, A. Daud, and Y. Adu-Gyamfi, “A region-based deep learning approach to automated retail checkout,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, pp. 3210–3215, June 2022.
- [172] Y. Kossale, M. Airaj, and A. Darouichi, “Mode collapse in generative adversarial networks: An overview,” in *2022 8th International Conference on Optimization and Applications (ICOA)*, pp. 1–6, IEEE, 2022.
- [173] Y. Ding, N. Mishra, and H. Hoffmann, “Generative and multi-phase learning for computer systems optimization,” in *Proceedings of the 46th International Symposium on Computer Architecture*, pp. 39–52, 2019.
- [174] A. Bandi, P. V. S. R. Adapa, and Y. E. V. P. K. Kuchi, “The power of generative ai: A review of requirements, models, input–output formats, evaluation metrics, and challenges,” *Future Internet*, vol. 15, no. 8, p. 260, 2023.
- [175] D. Q. Tran, A. Aboah, Y. Jeon, M. Shoman, M. Park, and S. Park, “Low-light image enhancement framework for improved object detection in fisheye lens datasets,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, pp. 7056–7065, June 2024.
- [176] N. Chen, A. Klushyn, R. Kurle, X. Jiang, J. Bayer, and P. Smagt, “Metrics for deep generative models,” in *International Conference on Artificial Intelligence and Statistics*, pp. 1540–1550, PMLR, 2018.
- [177] S. Barratt and R. Sharma, “A note on the inception score,” *arXiv preprint arXiv:1801.01973*, 2018.
- [178] A. Obukhov and M. Krasnyanskiy, “Quality assessment method for gan based on modified metrics inception score and fréchet inception distance,” in *Software Engineering Perspectives in Intelligent Systems: Proceedings of 4th Computational Methods in Systems and Software 2020, Vol. 1 4*, pp. 102–114, Springer, 2020.
- [179] A. Verine, B. Negrevergne, M. S. Pydi, and Y. Chevaleyre, “Precision-recall divergence optimization for generative modeling with gans and normalizing flows,” *arXiv preprint arXiv:2305.18910*, 2023.
- [180] R. Kansal, J. Duarte, H. Su, B. Orzari, T. Tomei, M. Pierini, M. Touranakou, D. Gunopoulos, et al., “Particle cloud generation with message passing generative adversarial networks,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 23858–23871, 2021.
- [181] A. Borji, “Pros and cons of gan evaluation measures: New developments,” *Computer Vision and Image Understanding*, vol. 215, p. 103329, 2022.
- [182] DeepLearning.AI, “Generative AI with LLMs.” <https://www.deeplearning.ai/courses/generative-ai-with-lm/>, n.d. Coursera.
- [183] A. Schmidt, “Speeding up the engineering of interactive systems with generative ai,” in *Companion Proceedings of the 2023 ACM SIGCHI Symposium on Engineering Interactive Computing Systems*, pp. 7–8, 2023.
- [184] H. Muse, S. Bulathwela, and E. Yilmaz, “Pre-training with scientific text improves educational question generation (student abstract)” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, pp. 16288–16289, 2023.
- [185] D. Foster, *Generative deep learning*. ” O’Reilly Media, Inc.”, 2022.
- [186] Y. Wang, W. Zhong, L. Li, F. Mi, X. Zeng, W. Huang, L. Shang, X. Jiang, and Q. Liu, “Aligning large language models with human: A survey,” *arXiv preprint arXiv:2307.12966*, 2023.
- [187] R. Zhong, K. Lee, Z. Zhang, and D. Klein, “Adapting language models for zero-shot learning by meta-tuning on dataset and prompt collections,” *arXiv preprint arXiv:2104.04670*, 2021.
- [188] H. Dang, L. Mecke, F. Lehmann, S. Goller, and D. Buschek, “How to prompt? opportunities and challenges of zero-and few-shot learning for human-ai interaction in creative applications of generative models,” *arXiv preprint arXiv:2209.01390*, 2022.
- [189] F. Song, B. Yu, M. Li, H. Yu, F. Huang, Y. Li, and H. Wang, “Preference ranking optimization for human alignment,” *arXiv preprint arXiv:2306.17492*, 2023.
- [190] H. Liu, Z. Teng, L. Cui, C. Zhang, Q. Zhou, and Y. Zhang, “Logicot: Logical chain-of-thought instruction-tuning data collection with gpt-4,” *arXiv preprint arXiv:2305.12147*, 2023.
- [191] J. Oppenlaender, “Prompt engineering for text-based generative art,” *arXiv preprint arXiv:2204.13988*, 2022.
- [192] N. Ratner, Y. Levine, Y. Belinkov, O. Ram, I. Magar, O. Abend, E. Karpas, A. Shashua, K. Leyton-Brown, and Y. Shoham, “Parallel context windows for large language models,” in *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 6383–6402, 2023.
- [193] T. Xue, Z. Wang, and H. Ji, “Parameter-efficient tuning helps language model alignment,” *arXiv preprint arXiv:2310.00819*, 2023.
- [194] Y. Liu, G. Deng, Z. Xu, Y. Li, Y. Zheng, Y. Zhang, L. Zhao, T. Zhang, and Y. Liu, “Jailbreaking chatgpt via prompt engineering: An empirical study,” *arXiv preprint arXiv:2305.13860*, 2023.
- [195] K. Yang, S. Ji, T. Zhang, Q. Xie, Z. Kuang, and S. Ananiadou, “Towards interpretable mental health analysis with chatgpt,” 2023.
- [196] J. White, S. Hays, Q. Fu, J. Spencer-Smith, and D. C. Schmidt, “Chatgpt prompt patterns for improving code quality, refactoring, requirements elicitation, and software design,” *arXiv preprint arXiv:2303.07839*, 2023.
- [197] C. Liu, X. Bao, H. Zhang, N. Zhang, H. Hu, X. Zhang, and M. Yan, “Improving chatgpt prompt for code generation,” *arXiv preprint arXiv:2305.08360*, 2023.
- [198] I. Singh, V. Blukis, A. Mousavian, A. Goyal, D. Xu, J. Tremblay, D. Fox, J. Thomason, and A. Garg, “Progrompt: Generating situated robot task plans using large language models,” in *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 11523–11530, IEEE, 2023.
- [199] S. K. K. Santu and D. Feng, “Teler: A general taxonomy of llm prompts for benchmarking complex tasks,” *arXiv preprint arXiv:2305.11430*, 2023.
- [200] J. Zamfirescu-Pereira, R. Y. Wong, B. Hartmann, and Q. Yang, “Why johnny can’t prompt: how non-ai experts try (and fail) to design llm prompts,” in *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, pp. 1–21, 2023.
- [201] J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou, et al., “Chain-of-thought prompting elicits reasoning in large language models,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 24824–24837, 2022.
- [202] S. Diao, P. Wang, Y. Lin, and T. Zhang, “Active prompting with chain-of-thought for large language models,” *arXiv preprint arXiv:2302.12246*, 2023.
- [203] P. Denny, J. Leinonen, J. Prather, A. Luxton-Reilly, T. Amarouche, B. A. Becker, and B. N. Reeves, “Promptly: Using prompt problems to teach learners how to effectively utilize ai code generators,” *arXiv preprint arXiv:2307.16364*, 2023.

- [204] S. S. Raman, V. Cohen, E. Rosen, I. Idrees, D. Paulius, and S. Tellex, “Planning with large language models via corrective re-prompting,” *arXiv preprint arXiv:2211.09935*, 2022.
- [205] Q. Dong, L. Li, D. Dai, C. Zheng, Z. Wu, B. Chang, X. Sun, J. Xu, and Z. Sui, “A survey for in-context learning,” *arXiv preprint arXiv:2301.00234*, 2022.
- [206] T. Kojima, S. S. Gu, M. Reid, Y. Matsuo, and Y. Iwasawa, “Large language models are zero-shot reasoners,” *Advances in neural information processing systems*, vol. 35, pp. 22199–22213, 2022.
- [207] F. Liu, J. M. Eisenschlos, F. Piccinno, S. Krichene, C. Pang, K. Lee, M. Joshi, W. Chen, N. Collier, and Y. Altun, “Deplot: One-shot visual language reasoning by plot-to-table translation,” *arXiv preprint arXiv:2212.10505*, 2022.
- [208] X. Liu, D. McDuff, G. Kovacs, I. Galatzer-Levy, J. Sunshine, J. Zhan, M.-Z. Poh, S. Liao, P. Di Achille, and S. Patel, “Large language models are few-shot health learners,” *arXiv preprint arXiv:2305.15525*, 2023.
- [209] Z. Hu, Y. Lan, L. Wang, W. Xu, E.-P. Lim, R. K.-W. Lee, L. Bing, and S. Poria, “Llm-adapters: An adapter family for parameter-efficient fine-tuning of large language models,” *arXiv preprint arXiv:2304.01933*, 2023.
- [210] D. Miyake, A. Iohara, Y. Saito, and T. Tanaka, “Negative-prompt inversion: Fast image inversion for editing with text-guided diffusion models,” *arXiv preprint arXiv:2305.16807*, 2023.
- [211] N. Liu, S. Li, Y. Du, A. Torralba, and J. B. Tenenbaum, “Compositional visual generation with composable diffusion models,” in *European Conference on Computer Vision*, pp. 423–439, Springer, 2022.
- [212] AUTOMATIC1111, “Negative-prompt.” <https://github.com/AUTOMATIC1111/stable-diffusion-webui/wiki/Negative-prompt>, 2022. Accessed on August 1, 2023.
- [213] F. Ma, C. Zhang, L. Ren, J. Wang, Q. Wang, W. Wu, X. Quan, and D. Song, “Xprompt: Exploring the extreme of prompt tuning,” *arXiv preprint arXiv:2210.04457*, 2022.
- [214] N. Tumanyan, M. Geyer, S. Bagor, and T. Dekel, “Plug-and-play diffusion features for text-driven image-to-image translation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 1921–1930, 2023.
- [215] H. Chang, H. Zhang, J. Barber, A. Maschinot, J. Lezama, L. Jiang, M.-H. Yang, K. Murphy, W. T. Freeman, M. Rubinstein, *et al.*, “Muse: Text-to-image generation via masked generative transformers,” *arXiv preprint arXiv:2301.00704*, 2023.
- [216] A. Chen, Y. Yao, P.-Y. Chen, Y. Zhang, and S. Liu, “Understanding and improving visual prompting: A label-mapping perspective,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 19133–19143, 2023.
- [217] A. Bar, Y. Gandelsman, T. Darrell, A. Globerson, and A. Efros, “Visual prompting via image inpainting,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 25005–25017, 2022.
- [218] M. Jia, L. Tang, B.-C. Chen, C. Cardie, S. Belongie, B. Hariharan, and S.-N. Lim, “Visual prompt tuning,” in *European Conference on Computer Vision*, pp. 709–727, Springer, 2022.
- [219] T. Chakrabarty, A. Saakyan, O. Winn, A. Panagopoulou, Y. Yang, M. Apidianaki, and S. Muresan, “I spy a metaphor: Large language models and diffusion models co-create visual metaphors,” *arXiv preprint arXiv:2305.14724*, 2023.
- [220] R. Volum, S. Rao, M. Xu, G. A. DesGrennes, C. Brockett, B. Van Durme, O. Deng, A. Malhotra, and B. Dolan, “Craft an iron sword: Dynamically generating interactive game characters by prompting large language models tuned on code,” in *The Third Wordplay: When Language Meets Games Workshop*, 2022.
- [221] D. Hegde, J. M. J. Valanarasu, and V. M. Patel, “Clip goes 3d: Leveraging prompt tuning for language grounded 3d recognition,” *arXiv preprint arXiv:2303.11313*, 2023.
- [222] J. Hoffmann, S. Borgeaud, A. Mensch, E. Buchatskaya, T. Cai, E. Rutherford, D. d. L. Casas, L. A. Hendricks, J. Welbl, A. Clark, *et al.*, “Training compute-optimal large language models,” *arXiv preprint arXiv:2203.15556*, 2022.
- [223] N. Chomsky, “Syntactic structures. mouton de gruyter,” *Mouton de Gruyter*, 2002.
- [224] L. Pan, A. Albalak, X. Wang, and W. Y. Wang, “Logic-lm: Empowering large language models with symbolic solvers for faithful logical reasoning,” *arXiv preprint arXiv:2305.12295*, 2023.
- [225] M. Du, F. He, N. Zou, D. Tao, and X. Hu, “Shortcut learning of large language models in natural language understanding: A survey,” *arXiv preprint arXiv:2208.11857*, 2022.
- [226] B. D. Lund, T. Wang, N. R. Mannuru, B. Nie, S. Shimray, and Z. Wang, “Chatgpt and a new academic reality: Artificial intelligence-written research papers and the ethics of the large language models in scholarly publishing,” *Journal of the Association for Information Science and Technology*, vol. 74, no. 5, pp. 570–581, 2023.
- [227] P. F. Brown, J. Cocke, S. A. Della Pietra, V. J. Della Pietra, F. Jelinek, J. Lafferty, R. L. Mercer, and P. S. Roossin, “A statistical approach to machine translation,” *Computational linguistics*, vol. 16, no. 2, pp. 79–85, 1990.
- [228] M. Setnes, R. Babuska, and H. B. Verbruggen, “Rule-based modeling: Precision and transparency,” *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, vol. 28, no. 1, pp. 165–169, 1998.
- [229] M. Wu, F. Liu, and T. Cohn, “Evaluating the utility of hand-crafted features in sequence labelling,” *arXiv preprint arXiv:1808.09075*, 2018.
- [230] E. D. Liddy, “Natural language processing,” 2001.
- [231] T. Hofmann, “Unsupervised learning by probabilistic latent semantic analysis,” *Machine learning*, vol. 42, pp. 177–196, 2001.
- [232] X. Liu and W. B. Croft, “Statistical language modeling,” *Annual Review of Information Science and Technology*, vol. 39, p. 1, 2004.
- [233] B.-H. Juang and L. R. Rabiner, “Automatic speech recognition—a brief history of the technology development,” *Georgia Institute of Technology. Atlanta Rutgers University and the University of California. Santa Barbara*, vol. 1, p. 67, 2005.
- [234] P. Azurin, *Transfer learning for natural language processing*. Simon and Schuster, 2021.
- [235] A. Kovačević and D. Kečo, “Bidirectional lstm networks for abstractive text summarization,” in *Advanced Technologies, Systems, and Applications VI: Proceedings of the International Symposium on Innovative and Interdisciplinary Applications of Advanced Technologies (IAT) 2021*, pp. 281–293, Springer, 2022.
- [236] Y. Wu, M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, W. Macherey, M. Krikun, Y. Cao, Q. Gao, K. Macherey, *et al.*, “Google’s neural machine translation system: Bridging the gap between human and machine translation,” *arXiv preprint arXiv:1609.08144*, 2016.
- [237] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, “End-to-end object detection with transformers,” in *European conference on computer vision*, pp. 213–229, Springer, 2020.
- [238] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz, *et al.*, “Transformers: State-of-the-art natural language processing,” in *Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations*, pp. 38–45, 2020.
- [239] A. Radford, K. Narasimhan, T. Salimans, I. Sutskever, *et al.*, “Improving language understanding by generative pre-training,” 2018.
- [240] N. A. Akbar, I. Darmayanti, S. M. Fati, and A. Muneer, “Deep learning of a pre-trained language model’s joke classifier using gpt-2,” *Journal of Hunan University Natural Sciences*, vol. 48, no. 8, 2021.
- [241] R. Dale, “Gpt-3: What’s it good for?,” *Natural Language Engineering*, vol. 27, no. 1, pp. 113–118, 2021.
- [242] T. B. Brown, B. Mann, N. Ryder, *et al.*, “Language models are few-shot learners,” in *Proceedings of the 34th Conference on Neural Information Processing Systems (NeurIPS)*, 2020.
- [243] L. Floridi and M. Chiriaci, “Gpt-3: Its nature, scope, limits, and consequences,” *Minds and Machines*, vol. 30, pp. 681–694, 2020.
- [244]
- [245] M. R. King, “Can bard, google’s experimental chatbot based on the lamda large language model, help to analyze the gender and racial diversity of authors in your cited scientific references?,” *Cellular and Molecular Bioengineering*, vol. 16, no. 2, pp. 175–179, 2023.
- [246] C. Khatri, B. Hedayatnia, A. Venkatesh, J. Nunn, Y. Pan, Q. Liu, H. Song, A. Gottardi, S. Kwatra, S. Pancholi, *et al.*, “Advancing the state of the art in open domain dialog systems through the alexa prize,” *arXiv preprint arXiv:1812.10757*, 2018.
- [247] A. Karpathy, “State of GPT.” <https://www.youtube.com/watch?v=bZQun8Y4L2A>, 2023.
- [248] X. Liu, H. Cheng, P. He, W. Chen, Y. Wang, H. Poon, and J. Gao, “Adversarial training for large neural language models,” *arXiv preprint arXiv:2004.08994*, 2020.
- [249] C. Chelba, T. Mikolov, M. Schuster, Q. Ge, T. Brants, P. Koehn, and T. Robinson, “One billion word benchmark for measuring progress in statistical language modeling,” *arXiv preprint arXiv:1312.3005*, 2013.
- [250] S. Biderman and S. et. al., “Pythia: A suite for analyzing large language models across training and scaling,” *arXiv preprint arXiv:2304.01373*, 2023.
- [251] D. Hernandez, T. Brown, T. Conerly, N. DasSarma, D. Drain, S. ElShowk, N. Elhage, Z. Hatfield-Dodds, T. Henighan, T. Hume, *et al.*, “Scaling laws and interpretability of learning from repeated data,” *arXiv preprint arXiv:2205.10487*, year=2022.

- [252] J. W. Rae, S. Borgeaud, T. Cai, K. Millican, J. Hoffmann, F. Song, J. Aslanides, S. Henderson, R. Ring, S. Young, *et al.*, “Scaling language models: Methods, analysis & insights from training gopher,” *arXiv preprint arXiv:2112.11446* , year=2021.
- [253] N. Carlini, D. Ippolito, M. Jagielski, K. Lee, F. Tramer, and C. Zhang, “Quantifying memorization across neural language models,” *arXiv preprint arXiv:2202.07646*, 2022.
- [254] P. Banerjee and H. Han, “Language modeling approaches to information retrieval,” 2009.
- [255] Z. Dai, Z. Yang, Y. Yang, J. Carbonell, Q. V. Le, and R. Salakhutdinov, “Transformer-xl: Attentive language models beyond a fixed-length context,” *arXiv preprint arXiv:1901.02860*, 2019.
- [256] I. Dergaa, K. Chamari, P. Zmijewski, and H. B. Saad, “From human writing to artificial intelligence generated text: examining the prospects and potential threats of chatgpt in academic writing,” *Biology of Sport*, vol. 40, no. 2, pp. 615–622, 2023.
- [257] L. Bottou, “Stochastic gradient descent tricks,” in *Neural Networks: Tricks of the Trade: Second Edition*, pp. 421–436, Springer, 2012.
- [258] P. J. Werbos, “Backpropagation through time: what it does and how to do it,” *Proceedings of the IEEE*, vol. 78, no. 10, pp. 1550–1560, 1990.
- [259] S. Praveen and V. Vajroboi, “Understanding the perceptions of healthcare researchers regarding chatgpt: a study based on bidirectional encoder representation from transformers (bert) sentiment analysis and topic modeling,” *Annals of Biomedical Engineering*, pp. 1–3, 2023.
- [260] J. Salazar, D. Liang, T. Q. Nguyen, and K. Kirchhoff, “Masked language model scoring,” *arXiv preprint arXiv:1910.14659*, 2019.
- [261] Y. Sun, Y. Zheng, C. Hao, and H. Qiu, “Nsp-bert: A prompt-based zero-shot learner through an original pre-training task–next sentence prediction,” *arXiv e-prints*, pp. arXiv–2109, 2021.
- [262] W. Zhao, H. Hu, W. Zhou, J. Shi, and H. Li, “Best: Bert pre-training for sign language recognition with coupling tokenization,” *arXiv preprint arXiv:2302.05075* , year=2023.
- [263] C. Sun, X. Qiu, Y. Xu, and X. Huang, “How to fine-tune bert for text classification?,” in *Chinese Computational Linguistics: 18th China National Conference, CCL 2019, Kunming, China, October 18–20, 2019, Proceedings 18*, pp. 194–206, Springer, 2019.
- [264] L. Jiarong, X. Hong, J. Wenchao, Y. Jianren, and W. Tao, “Knowledge enhanced bert based on corpus associate generation,” in *Machine Learning for Cyber Security: 4th International Conference, ML4CS 2022, Guangzhou, China, December 2–4, 2022, Proceedings, Part III*, pp. 533–547, Springer, 2023.
- [265] M. Irfan, A. I. Sanka, Z. Ullah, and R. C. Cheung, “Reconfigurable content-addressable memory (CAM) on FPGAs: A tutorial and survey,” *Future Generation Computer Systems*, vol. 128, pp. 451–465, 2022.
- [266] L. Fan, L. Li, Z. Ma, S. Lee, H. Yu, and L. Hemphill, “A bibliometric review of large language models research from 2017 to 2023,” *arXiv preprint arXiv:2304.02020*, 2023.
- [267] T. Dettmers, M. Lewis, Y. Belkada, and L. Zettlemoyer, “Llm. int8 () : 8-bit matrix multiplication for transformers at scale,” *arXiv preprint arXiv:2208.07339*, 2022.
- [268] Z. Liu, J. Wang, T. Dao, T. Zhou, B. Yuan, Z. Song, A. Shrivastava, C. Zhang, Y. Tian, C. Re, *et al.*, “Deja vu: Contextual sparsity for efficient llms at inference time,” in *International Conference on Machine Learning*, pp. 22137–22176, PMLR, 2023.
- [269] Z. Yang, Z. Dai, Y. Yang, J. Carbonell, R. R. Salakhutdinov, and Q. V. Le, “Xlnet: Generalized autoregressive pretraining for language understanding,” *Advances in neural information processing systems*, vol. 32, 2019.
- [270] H. Ye, Z. Chen, D.-H. Wang, and B. Davison, “Pretrained generalized autoregressive model with adaptive probabilistic label clusters for extreme multi-label text classification,” in *International Conference on Machine Learning*, pp. 10809–10819, PMLR, 2020.
- [271] J. Su, S. Yu, and D. Luo, “Enhancing aspect-based sentiment analysis with capsule network,” *IEEE Access*, vol. 8, pp. 100551–100561, 2020.
- [272] M. Kolbæk, D. Yu, Z.-H. Tan, and J. Jensen, “Multitalker speech separation with utterance-level permutation invariant training of deep recurrent neural networks,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 25, no. 10, pp. 1901–1913, 2017.
- [273] M. Hobbhahn, T. Lieberum, and D. Seiler, “Investigating causal understanding in llms,” in *NeurIPS ML Safety Workshop*, 2022.
- [274] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu, “Exploring the limits of transfer learning with a unified text-to-text transformer,” *Journal of Machine Learning Research*, vol. 21, no. 140, pp. 1–67, 2020.
- [275] J. Cho, J. Lei, H. Tan, and M. Bansal, “Unifying vision-and-language tasks via text generation,” in *International Conference on Machine Learning*, pp. 1931–1942, PMLR, 2021.
- [276] N. S. Keskar, B. McCann, L. R. Varshney, C. Xiong, and R. Socher, “Ctrl: A conditional transformer language model for controllable generation,” *arXiv preprint arXiv:1909.05858*, 2019.
- [277] S.-A. Rebuffi, H. Bilen, and A. Vedaldi, “Learning multiple visual domains with residual adapters,” *Advances in neural information processing systems*, vol. 30, 2017.
- [278] S. Smith *et al.*, “Common crawl corpus,” <https://commoncrawl.org/>, 2013.
- [279] Y. Zhu, R. Kiros, R. S. Zemel, R. Salakhutdinov, R. Urtasun, A. Torralba, and S. Fidler, “Aligning books and movies: Towards story-like visual explanations by watching movies and reading books,” *Proceedings of the IEEE International Conference on Computer Vision*, pp. 19–27, 2015.
- [280] L. Gao, S. Biderman, S. Black, *et al.*, “The pile: An 800gb dataset of diverse text for language modeling,” *arXiv preprint arXiv:2101.00027*, 2020.
- [281] P. Koehn, “Europarl: A parallel corpus for statistical machine translation,” in *MT summit*, vol. 5, pp. 79–86, 2005.
- [282] J. Baumgartner, S. Zannettou, B. Keegan, M. Squire, and J. Blackburn, “The pushshift reddit dataset,” *Proceedings of the International AAAI Conference on Web and Social Media*, vol. 14, pp. 830–839, 2020.
- [283] W. contributors, “Wikipedia dumps.” <https://dumps.wikimedia.org/>, 2023.
- [284] P. J. Ortiz Suárez, B. Sagot, and L. Romary, “Asynchronous pipeline for processing huge corpora on medium to low resource infrastructures,” in *Proceedings of the Workshop on Challenges in the Management of Large Corpora (CMLC-7) 2019: 7th Workshop on Challenges in the Management of Large Corpora (CMLC-7)*, pp. 9–16, 2019.
- [285] D. Araci, “Finbert: A pretrained language model for financial communications,” in *arXiv preprint arXiv:1908.10063*, 2019.
- [286] Y. Zhang, M. Ghaly, and A. Sarker, “Bert-based clinical natural language processing for medication-related information extraction,” *AMIA Summits on Translational Science Proceedings*, vol. 2020, p. 516, 2020.
- [287] “Pubmed.” <https://pubmed.ncbi.nlm.nih.gov/>, 2020.
- [288] B. Workshop *et al.*, “The bigscience roots corpus: A 1.6tb composite multilingual dataset,” *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (ACL)*, 2023.
- [289] S. CRFM, “Stanford alpaca: An instruction-following llama model,” 2023. <https://crfm.stanford.edu/2023/03/13/alpaca.html>.
- [290] Y. Wang, Y. Kordi, S. Mishra, A. Liu, N. A. Smith, H. Hajishirzi, *et al.*, “Self-instruct: Aligning language model with self generated instructions,” *arXiv preprint arXiv:2212.10560*, 2022.
- [291] D. Hendrycks, C. Burns, S. Basart, A. Zou, M. Mazeika, D. Song, and J. Steinhardt, “Measuring massive multitask language understanding,” in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2020.
- [292] K. Cobbe, V. Kosaraju, M. Bavarian, M. Chen, , *et al.*, “Training verifiers to solve math word problems,” in *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2021.
- [293] Z. Xie, G. Lai, Z. Dai, E. Hovy, A. Yates, *et al.*, “Cloth: A large-scale dataset for cloze test in high school english examinations,” in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2018.
- [294] S. Reddy, D. Chen, and C. D. Manning, “Coqa: A conversational question answering challenge,” in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2019.
- [295] P. Rajpurkar, J. Zhang, K. Lopyrev, and P. Liang, “Squad: 100,000+ questions for machine comprehension of text,” in *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2016.
- [296] B. Y. Lin, W. Zhou, M. Shen, P. Zhou, *et al.*, “Commongen: A constrained text generation challenge for generative commonsense reasoning,” in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2020.
- [297] A. See, P. J. Liu, and C. D. Manning, “Get to the point: Summarization with pointer-generator networks,” *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (ACL)*, 2017.
- [298] M. Sap, H. Rashkin, D. Chen, R. Le Bras, *et al.*, “Social iq: Commonsense reasoning about social interactions,” *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, 2019.
- [299] O. Bojar, R. Chatterjee, C. Federmann, B. Haddow, *et al.*, “Findings of the 2014 workshop on statistical machine translation,” in *Proceedings*

- of the Ninth Workshop on Statistical Machine Translation*, pp. 12–58, 2014.
- [300] S. Zhang, S. Roller, N. Goyal, *et al.*, “Opt: Open pre-trained transformer language models,” in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (ACL)*, 2022.
- [301] J. Lee, W. Yoon, S. Kim, D. Kim, *et al.*, “Biobert: a pre-trained biomedical language representation model for biomedical text mining,” *Bioinformatics*, vol. 36, no. 4, pp. 1234–1240, 2020.
- [302] I. Chalkidis, M. Fergadiotis, N. Nikiforos, *et al.*, “Legal-bert: The muppets straight out of law school,” in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: Findings*, 2020.
- [303] L. Ouyang, J. Wu, X. Jiang, C. Almeida, *et al.*, “Training language models to follow instructions with human feedback,” *arXiv preprint arXiv:2203.02155*, 2022.
- [304] C. Schuhmann, R. Beaumont, R. Vencu, *et al.*, “Laion-5b: An open large-scale dataset for training next generation image-text models,” *arXiv preprint arXiv:2210.08402*, 2022.
- [305] Y. Bai, S. Kadavath, S. Kundu, *et al.*, “Constitutional ai: Harmlessness from ai feedback,” *arXiv preprint arXiv:2212.08073*, 2022.
- [306] H. Zhong, A. Madaan, R. Kumar, *et al.*, “Agieval: A human-centric benchmark for evaluating foundation models,” in *arXiv preprint arXiv:2304.06364*, 2023.
- [307] A. Wang, Y. Prusachatkun, N. Nangia, *et al.*, “Superglue: A stickier benchmark for general-purpose language understanding systems,” *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [308] R. Zellers, A. Holtzman, Y. Bisk, *et al.*, “Hellaswag: Can a machine really finish your sentence?,” in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 2019.
- [309] P. Clark, W. Cowell, M. Gardner, *et al.*, “Think you have solved question answering? try arc, the ai2 reasoning challenge,” in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2018.
- [310] D. Patil, S. Iyer, *et al.*, “Gorilla: Large language models connected with massive apis,” *arXiv preprint arXiv:2305.15334*, 2023.
- [311] I. Chalkidis, N. Aletras, *et al.*, “Lexglue: A benchmark dataset for legal language understanding in english,” *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2021.
- [312] D. Hendrycks, C. Burns, *et al.*, “Aligning ai with shared human values,” in *Proceedings of the 2021 Conference on Neural Information Processing Systems (NeurIPS)*, 2021.
- [313] M. Geva, D. Khashabi, *et al.*, “Did aristotle use a laptop? a question answering benchmark with implicit reasoning strategies,” *Transactions of the Association for Computational Linguistics*, vol. 9, pp. 346–361, 2021.
- [314] A. Srivastava, H. He, D. Ippolito, *et al.*, “Big-bench: Beyond the imitation game,” in *Proceedings of the 2022 Conference on Neural Information Processing Systems (NeurIPS)*, 2022.
- [315] M. Joshi, E. Choi, D. Weld, and L. Zettlemoyer, “Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension,” in *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics*, 2017.
- [316] S. Pradhan, A. Moschitti, N. Xue, *et al.*, “Conll-2012 shared task: Modeling multilingual unrestricted coreference in ontonotes,” *Proceedings of the Sixteenth Conference on Computational Natural Language Learning (CoNLL)*, 2012.
- [317] A. L. Maas, R. E. Daly, P. T. Pham, D. Huang, *et al.*, “Learning word vectors for sentiment analysis,” *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics (ACL)*, 2011.
- [318] R. Socher, A. Perelygin, J. Wu, J. Chuang, C. D. Manning, A. Y. Ng, and C. Potts, “Recursive deep models for semantic compositionality over a sentiment treebank,” *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2013.
- [319] W. B. Dolan and C. Brockett, “Automatically constructing a corpus of sentential paraphrases,” in *Proceedings of the Third International Workshop on Paraphrasing (IWP)*, 2005.
- [320] D. Cer, M. Diab, E. Agirre, I. Lopez-Gazpio, and L. Specia, “Semeval-2017 task 1: Semantic textual similarity multilingual and crosslingual focused evaluation,” in *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, 2017.
- [321] S. Narayan, S. B. Cohen, and M. Lapata, “Don’t give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization,” in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2018.
- [322] X. Zhang, J. Zhao, and Y. LeCun, “Character-level convolutional networks for text classification,” in *Proceedings of the 28th International Conference on Neural Information Processing Systems (NeurIPS)*, 2015.
- [323] E. M. Voorhees *et al.*, “The trec-8 question answering track report,” in *Proceedings of the Eighth Text Retrieval Conference (TREC-8)*, 1999.
- [324] A. Warstadt, A. Singh, and S. R. Bowman, “Neural network acceptability judgments,” *Transactions of the Association for Computational Linguistics*, vol. 7, pp. 625–641, 2019.
- [325] J. Austin, A. Odena, *et al.*, “Program synthesis with large language models,” *arXiv preprint arXiv:2108.07732*, 2021.
- [326] E. F. Tjong Kim Sang and F. De Meulder, “Introduction to the conll-2003 shared task: Language-independent named entity recognition,” in *Proceedings of the Seventh Conference on Natural Language Learning at HLT-NAACL 2003*, pp. 142–147, 2003.
- [327] Q. Zhang, P. Wang, W. Wei, and T. Zhou, “Position-aware attention and supervised data improve slot filling,” in *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2017.
- [328] A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, and S. R. Bowman, “Glue: A multi-task benchmark and analysis platform for natural language understanding,” *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2018.
- [329] D. Su, Y. Xu, G. I. Winata, P. Xu, H. Kim, Z. Liu, and P. Fung, “Generalizing question answering system with pre-trained language model fine-tuning,” in *Proceedings of the 2nd Workshop on Machine Reading for Question Answering*, pp. 203–211, 2019.
- [330] N. Pawłowski, J. Vaughan, J. Jennings, and C. Zhang, “Answering causal questions with augmented llms,” 2023.
- [331] Y. Zhuang, Y. Yu, K. Wang, H. Sun, and C. Zhang, “Toolqa: A dataset for llm question answering with external tools,” *arXiv preprint arXiv:2306.13304*, 2023.
- [332] K. Singhal, S. Azizi, T. Tu, S. S. Mahdavi, J. Wei, H. W. Chung, N. Scales, A. Tanwani, H. Cole-Lewis, S. Pfahl, *et al.*, “Large language models encode clinical knowledge,” *Nature*, pp. 1–9, 2023.
- [333] A. Piñeiro-Martín, C. García-Mateo, L. Docío-Fernández, and M. d. C. López-Pérez, “Ethical challenges in the development of virtual assistants powered by large language models,” *Electronics*, vol. 12, no. 14, p. 3170, 2023.
- [334] A. Lazaridou, E. Gribovskaya, W. Stokowiec, and N. Grigorev, “Internet-augmented language models through few-shot prompting for open-domain question answering,” *arXiv preprint arXiv:2203.05115*, 2022.
- [335] D. Zhu, J. Chen, X. Shen, X. Li, and M. Elhoseiny, “Minigpt-4: Enhancing vision-language understanding with advanced large language models,” *arXiv preprint arXiv:2304.10592*, 2023.
- [336] A. Celikyilmaz, E. Clark, and J. Gao, “Evaluation of text generation: A survey,” *arXiv preprint arXiv:2006.14799*, 2020.
- [337] S. Dathathri, A. Madotto, J. Lan, J. Hung, E. Frank, P. Molino, J. Yosinski, and R. Liu, “Plug and play language models: A simple approach to controlled text generation,” *arXiv preprint arXiv:1912.02164*, 2019.
- [338] Y. Li, Q. Pan, S. Wang, T. Yang, and E. Cambria, “A generative model for category text generation,” *Information Sciences*, vol. 450, pp. 301–315, 2018.
- [339] L. Wang, C. Lyu, T. Ji, Z. Zhang, D. Yu, S. Shi, and Z. Tu, “Document-level machine translation with large language models,” *arXiv preprint arXiv:2304.02210*, 2023.
- [340] H. Huang, S. Wu, X. Liang, B. Wang, Y. Shi, P. Wu, M. Yang, and T. Zhao, “Towards making the most of llm for translation quality estimation,” in *CCF International Conference on Natural Language Processing and Chinese Computing*, pp. 375–386, Springer, 2023.
- [341] C. Lyu, J. Xu, and L. Wang, “New trends in machine translation using large language models: Case examples with chatgpt,” *arXiv preprint arXiv:2305.01181*, 2023.
- [342] B. Zhang, B. Haddow, and A. Birch, “Prompting large language model for machine translation: A case study,” *arXiv preprint arXiv:2301.07069*, 2023.
- [343] X. Sun, X. Li, J. Li, F. Wu, S. Guo, T. Zhang, and G. Wang, “Text classification via large language models,” *arXiv preprint arXiv:2305.08377*, 2023.
- [344] R. Song, Z. Liu, X. Chen, H. An, Z. Zhang, X. Wang, and H. Xu, “Label prompt for multi-label text classification,” *Applied Intelligence*, vol. 53, no. 8, pp. 8761–8775, 2023.
- [345] W. Zhang, Y. Deng, B. Liu, S. J. Pan, and L. Bing, “Sentiment analysis in the era of large language models: A reality check,” *arXiv preprint arXiv:2305.15005*, 2023.

- [346] M. Labonne and S. Moran, "Spam-t5: Benchmarking large language models for few-shot email spam detection," *arXiv preprint arXiv:2304.01238*, 2023.
- [347] S. S. Mullick, M. Bhamhani, S. Sinha, A. Mathur, S. Gupta, and J. Shah, "Content moderation for evolving policies using binary question answering," in *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 5: Industry Track)*, pp. 561–573, 2023.
- [348] B. Peng, M. Galley, P. He, H. Cheng, Y. Xie, Y. Hu, Q. Huang, L. Liden, Z. Yu, W. Chen, et al., "Check your facts and try again: Improving large language models with external knowledge and automated feedback," *arXiv preprint arXiv:2302.12813*, 2023.
- [349] R. C. Fernandez, A. J. Elmore, M. J. Franklin, S. Krishnan, and C. Tan, "How large language models will disrupt data management," *Proceedings of the VLDB Endowment*, vol. 16, no. 11, pp. 3302–3309, 2023.
- [350] S. Yu, C. Fang, Y. Ling, C. Wu, and Z. Chen, "Llm for test script generation and migration: Challenges, capabilities, and opportunities," *arXiv preprint arXiv:2309.13574*, 2023.
- [351] T. Zhang, F. Ladhak, E. Durmus, P. Liang, K. McKeown, and T. B. Hashimoto, "Benchmarking large language models for news summarization," *arXiv preprint arXiv:2301.13848*, 2023.
- [352] C. Shen, L. Cheng, Y. You, and L. Bing, "Are large language models good evaluators for abstractive summarization?," *arXiv preprint arXiv:2305.13091*, 2023.
- [353] J. Wu, R. Antonova, A. Kan, M. Lepert, A. Zeng, S. Song, J. Bohg, S. Rusinkiewicz, and T. Funkhouser, "Tidybot: Personalized robot assistance with large language models," *arXiv preprint arXiv:2305.05658*, 2023.
- [354] R. Luo, Z. Zhao, M. Yang, J. Dong, M. Qiu, P. Lu, T. Wang, and Z. Wei, "Valley: Video assistant with large language model enhanced ability," *arXiv preprint arXiv:2306.07207*, 2023.
- [355] S. Wadhwa, S. Amir, and B. C. Wallace, "Revisiting relation extraction in the era of large language models," *arXiv preprint arXiv:2305.05003*, 2023.
- [356] X. Wei, X. Cui, N. Cheng, X. Wang, X. Zhang, S. Huang, P. Xie, J. Xu, Y. Chen, M. Zhang, et al., "Zero-shot information extraction via chatting with chatgpt," *arXiv preprint arXiv:2302.10205*, 2023.
- [357] C. Li, X. Zhang, D. Chrysostomou, and H. Yang, "Tod4ir: A humanised task-oriented dialogue system for industrial robots," *IEEE Access*, vol. 10, pp. 91631–91649, 2022.
- [358] J. Deng, H. Sun, Z. Zhang, J. Cheng, and M. Huang, "Recent advances towards safe, responsible, and moral dialogue systems: A survey," *arXiv preprint arXiv:2302.09270*, 2023.
- [359] J. Wei, J. Wei, Y. Tay, D. Tran, A. Webson, Y. Lu, X. Chen, H. Liu, D. Huang, D. Zhou, et al., "Larger language models do in-context learning differently," *arXiv preprint arXiv:2303.03846*, 2023.
- [360] V. Bhat and P. Bhattacharyya, "Survey: Automatic speech recognition for indian languages,"
- [361] A. Radford, J. W. Kim, T. Xu, G. Brockman, C. McLeavey, and I. Sutskever, "Robust speech recognition via large-scale weak supervision," in *International Conference on Machine Learning*, pp. 28492–28518, PMLR, 2023.
- [362] H. W. Chung, L. Hou, S. Longpre, B. Zoph, Y. Tay, W. Fedus, E. Li, X. Wang, M. Dehghani, S. Brahma, et al., "Scaling instruction-finetuned language models," *arXiv preprint arXiv:2210.11416*, 2022.
- [363] M. Suzgun, N. Scales, N. Schärli, S. Gehrmann, Y. Tay, H. W. Chung, A. Chowdhery, Q. V. Le, E. H. Chi, D. Zhou, et al., "Challenging big-bench tasks and whether chain-of-thought can solve them," *arXiv preprint arXiv:2210.09261*, 2022.
- [364] M. Chen, J. Tworek, H. Jun, Q. Yuan, H. P. de Oliveira Pinto, J. Kaplan, H. Edwards, Y. Burda, N. Joseph, G. Brockman, A. Ray, R. Puri, G. Krueger, M. Petrov, H. Khlaaf, G. Sastry, P. Mishkin, B. Chan, S. Gray, N. Ryder, M. Pavlov, A. Power, L. Kaiser, M. Bavarian, C. Winter, P. Tillet, F. P. Such, D. Cummings, M. Plappert, F. Chantzis, E. Barnes, A. Herbert-Voss, W. H. Guss, A. Nichol, A. Paino, N. Tezak, J. Tang, I. Babuschkin, S. Balaji, S. Jain, W. Saunders, C. Hesse, A. N. Carr, J. Leike, J. Achiam, V. Misra, E. Morikawa, A. Radford, M. Knight, M. Brundage, M. Murati, K. Mayer, P. Welinder, B. McGrew, D. Amodei, S. McCandlish, I. Sutskever, and W. Zaremba, "Evaluating large language models trained on code," 2021.
- [365] Y. Huang, Y. Bai, Z. Zhu, J. Zhang, J. Zhang, T. Su, J. Liu, C. Lv, Y. Zhang, J. Lei, et al., "C-eval: A multi-level multi-discipline chinese evaluation suite for foundation models," *arXiv preprint arXiv:2305.08322*, 2023.
- [366] W. Chen and E. W. X. M. J. X. T. X. X. W. P. L. Ming Yin, Max Ku, "Theoremqa: A theorem-driven question answering dataset," *arXiv preprint arXiv:2305.12524*, 2023.
- [367] S. Bubeck, V. Chandrasekaran, R. Eldan, J. Gehrke, E. Horvitz, E. Kamar, P. Lee, Y. T. Lee, Y. Li, S. Lundberg, et al., "Sparks of artificial general intelligence: Early experiments with gpt-4," *arXiv preprint arXiv:2303.12712*, 2023.
- [368] R. Anil, A. M. Dai, O. Firat, M. Johnson, D. Lepikhin, A. Passos, S. Shakeri, E. Taropa, P. Bailey, Z. Chen, et al., "Palm 2 technical report," *arXiv preprint arXiv:2305.10403*, 2023.
- [369] Y. Wang, H. Le, A. D. Gotmare, N. D. Bui, J. Li, and S. C. Hoi, "Codet5+: Open code large language models for code understanding and generation," *arXiv preprint arXiv:2305.07922*, 2023.
- [370] S. A. Basit, R. Qureshi, S. Musleh, R. Guler, M. S. Rahman, K. H. Biswas, and T. Alam, "Covid-19base v3: Update of the knowledgebase for drugs and biomedical entities linked to covid-19," *Frontiers in Public Health*, vol. 11, p. 1125917, 2023.
- [371] F. C. Kitamura, "Chatgpt is shaping the future of medical writing but still requires human judgment," 2023.
- [372] T. H. Kung, M. Cheatham, A. Medenilla, C. Sillos, L. De Leon, C. Elepaño, M. Madriaga, R. Aggabao, G. Diaz-Candido, J. Maningo, et al., "Performance of chatgpt on usmle: Potential for ai-assisted medical education using large language models," *PLoS digital health*, vol. 2, no. 2, p. e0000198, 2023.
- [373] M. Sallam, "Chatgpt utility in healthcare education, research, and practice: systematic review on the promising perspectives and valid concerns," in *Healthcare*, vol. 11, p. 887, MDPI, 2023.
- [374] A. Gilson, C. W. Safranek, T. Huang, V. Socrates, L. Chi, R. A. Taylor, D. Chartash, 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*, vol. 9, no. 1, p. e45312, 2023.
- [375] M. Cascella, J. Montomoli, V. Bellini, and E. Bignami, "Evaluating the feasibility of chatgpt in healthcare: an analysis of multiple clinical and research scenarios," *Journal of Medical Systems*, vol. 47, no. 1, p. 33, 2023.
- [376] S. K. Karn, R. Ghosh, O. Farri, et al., "shs-nlp at radsum23: Domain-adaptive pre-training of instruction-tuned llms for radiology report impression generation," *arXiv preprint arXiv:2306.03264*, 2023.
- [377] A. Rao, J. Kim, M. Kamineni, M. Pang, W. Lie, and M. D. Succi, "Evaluating chatgpt as an adjunct for radiologic decision-making," *medRxiv*, pp. 2023–02, 2023.
- [378] D. Duong and B. D. Solomon, "Analysis of large-language model versus human performance for genetics questions," *medRxiv*, pp. 2023–01, 2023.
- [379] N. Fijačko, L. Gosak, G. Štiglic, C. T. Picard, and M. J. Douma, "Can chatgpt pass the life support exams without entering the american heart association course?," *Resuscitation*, vol. 185, 2023.
- [380] M. F. Romano, L. C. Shih, I. C. Paschalidis, R. Au, and V. B. Kolachalam, "Large language models in neurology research and future practice," *Neurology*, 2023.
- [381] M. R. Haque and S. Rubya, "An overview of chatbot-based mobile mental health apps: Insights from app description and user reviews," *JMIR mHealth and uHealth*, vol. 11, no. 1, p. e44838, 2023.
- [382] S. M. Jungmann, T. Klan, S. Kuhn, and F. Jungmann, "Accuracy of a chatbot (ada) in the diagnosis of mental disorders: comparative case study with lay and expert users," *JMIR formative research*, vol. 3, no. 4, p. e13863, 2019.
- [383] D. Magalhaes Azevedo and S. Kieffer, "User reception of ai-enabled mhealth apps: The case of babylon health.."
- [384] P. Malik, M. Pathania, V. K. Rathaur, et al., "Overview of artificial intelligence in medicine," *Journal of family medicine and primary care*, vol. 8, no. 7, p. 2328, 2019.
- [385] mbzuai oryx, "Xraygpt: Chest radiographs summarization using medical vision-language models," 2023.
- [386] J. Ma and B. Wang, "Segment anything in medical images," *arXiv preprint arXiv:2304.12306*, 2023.
- [387] Y. Li, C. Gao, X. Song, X. Wang, Y. Xu, and S. Han, "Druggpt: A gpt-based strategy for designing potential ligands targeting specific proteins," *bioRxiv*, pp. 2023–06, 2023.
- [388] M. Moor, O. Banerjee, Z. S. H. Abad, H. M. Krumholz, J. Leskovec, E. J. Topol, and P. Rajpurkar, "Foundation models for generalist medical artificial intelligence," *Nature*, vol. 616, no. 7956, pp. 259–265, 2023.
- [389] C. Cohn, N. Hutchins, and G. Biswas, "Towards a formative feedback generation agent: Leveraging a human-in-the-loop, chain-of-thought

- prompting approach with llms to evaluate formative assessment responses in k-12 science.,” 2023.
- [390] Y. Wu, S. Y. Min, Y. Bisk, R. Salakhutdinov, A. Azaria, Y. Li, T. Mitchell, and S. Prabhumoye, “Plan, eliminate, and track–language models are good teachers for embodied agents,” *arXiv preprint arXiv:2305.02412*, 2023.
- [391] S. Doveh, A. Arbelae, S. Harary, E. Schwartz, R. Herzig, R. Giryes, R. Feris, R. Panda, S. Ullman, and L. Karlinsky, “Teaching structured vision & language concepts to vision & language models,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2657–2668, 2023.
- [392] S. Saha, P. Hase, and M. Bansal, “Can language models teach weaker agents? teacher explanations improve students via theory of mind,” *arXiv preprint arXiv:2306.09299*, 2023.
- [393] J. S. () and W. Y. (), “Unlocking the power of chatgpt: A framework for applying generative ai in education,” *ECNU Review of Education*, vol. 0, no. 0, p. 20965311231168423, 0.
- [394] “An era of chatgpt as a significant futuristic support tool: A study on features, abilities, and challenges,” *BenchCouncil Transactions on Benchmarks, Standards and Evaluations*, vol. 2, no. 4, p. 100089, 2022.
- [395] H. Crompton and D. Burke, “Artificial intelligence in higher education: the state of the field,” *International Journal of Educational Technology in Higher Education*, vol. 20, no. 1, p. 22, 2023.
- [396] L. Zhu, W. Mou, T. Yang, and R. Chen, “Chatgpt can pass the aha exams: Open-ended questions outperform multiple-choice format,” *Resuscitation*, vol. 188, p. 109783, 2023.
- [397] E. Kasneci, K. Sessler, S. Küchemann, M. Bannert, D. Dementieva, F. Fischer, U. Gasser, G. Groh, S. Günemann, E. Hüllermeier, S. Krusche, G. Kutyniok, T. Michaeli, C. Nerdel, J. Pfeffer, O. Poquet, M. Sailer, A. Schmidt, T. Seidel, M. Stadler, J. Weller, J. Kuhn, and G. Kasneci, “Chatgpt for good? on opportunities and challenges of large language models for education,” *Learning and Individual Differences*, vol. 103, p. 102274, 2023.
- [398] B.-C. Kuo, F. T. Chang, and Z.-E. Bai, “Leveraging llms for adaptive testing and learning in taiwan adaptive learning platform (talp),” 2023.
- [399] “Khan academy explores the potential for gpt-4 in a limited pilot program,” 2023.
- [400] “Harnessing gpt-4 so that all students benefit. a nonprofit approach for equal access,” 2023.
- [401] T. Soubhari, S. S. Nanda, T. A. Lone, and P. S. Beegam, “Digital hacks, creativity shacks, and academic menace: The ai effect,” in *Sustainable Development Goal Advancement Through Digital Innovation in the Service Sector*, pp. 208–232, IGI Global, 2023.
- [402] J. Prather, P. Denny, J. Leinonen, B. A. Becker, I. Albluwi, M. Craig, H. Keuning, N. Kiesler, T. Kohn, A. Luxton-Reilly, et al., “The robots are here: Navigating the generative ai revolution in computing education,” *arXiv preprint arXiv:2310.00658*, 2023.
- [403] H. H. Thorp, “Chatgpt is fun, but not an author,” 2023.
- [404] C. Stokel-Walker, “Chatgpt listed as author on research papers: many scientists disapprove,” *Nature*, vol. 613, no. 7945, pp. 620–621, 2023.
- [405] O. Buruk, “Academic writing with gpt-3.5: Reflections on practices, efficacy and transparency,” *arXiv preprint arXiv:2304.11079*, 2023.
- [406] E. Hannan and S. Liu, “Ai: new source of competitiveness in higher education,” *Competitiveness Review: An International Business Journal*, vol. 33, no. 2, pp. 265–279, 2023.
- [407] J. Li, M. Xu, L. Xiang, D. Chen, W. Zhuang, X. Yin, and Z. Li, “Foundation models in smart agriculture: Basics, opportunities, and challenges,” *Computers and Electronics in Agriculture*, vol. 222, p. 109032, 2024.
- [408] Z. Rui, Z. Zhang, M. Zhang, A. Azizi, C. Igathinathane, H. Cen, S. Vougioukas, H. Li, J. Zhang, Y. Jiang, et al., “High-throughput proximal ground crop phenotyping systems—a comprehensive review,” *Computers and Electronics in Agriculture*, vol. 224, p. 109108, 2024.
- [409] J. Roberts, T. Lüddecke, R. Sheikh, K. Han, and S. Albanie, “Charting new territories: Exploring the geographic and geospatial capabilities of multimodal llms,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 554–563, 2024.
- [410] X. Yang, J. Gao, W. Xue, and E. Alexandersson, “Pllama: An open-source large language model for plant science,” *arXiv preprint arXiv:2401.01600*, 2024.
- [411] F. Stella, C. Della Santina, and J. Hughes, “How can llms transform the robotic design process?,” *Nature machine intelligence*, vol. 5, no. 6, pp. 561–564, 2023.
- [412] K. Tan, “Large language models for crop yield prediction,” 2024.
- [413] J. M. Shutske, “Harnessing the power of large language models in agricultural safety & health,” *Journal of Agricultural Safety and Health*, p. 0, 2023.
- [414] B. Arora, D. S. Chaudhary, M. Satsangi, M. Yadav, L. Singh, and P. S. Sudhish, “Agribot: a natural language generative neural networks engine for agricultural applications,” in *2020 International Conference on Contemporary Computing and Applications (IC3A)*, pp. 28–33, IEEE, 2020.
- [415] J. Qing, X. Deng, Y. Lan, and Z. Li, “Gpt-aided diagnosis on agricultural image based on a new light yolop,” *Computers and Electronics in Agriculture*, vol. 213, p. 108168, 2023.
- [416] M. T. Kuska, M. Wahabzada, and S. Paulus, “Ai-chatbots for agriculture-where can large language models provide substantial value?,” *Available at SSRN 4685971*, 2024.
- [417] C. Dhavale, T. Pawar, A. Singh, S. Pole, and K. Sabat, “Revolutionizing farming: Gan-enhanced imaging, cnn disease detection, and llm farmer assistant,” in *2024 2nd International Conference on Computer, Communication and Control (IC4)*, pp. 1–6, IEEE, 2024.
- [418] Y. Cao, L. Chen, Y. Yuan, and G. Sun, “Cucumber disease recognition with small samples using image-text-label-based multi-modal language model,” *Computers and electronics in agriculture*, vol. 211, p. 107993, 2023.
- [419] G. Lu, S. Li, G. Mai, J. Sun, D. Zhu, L. Chai, H. Sun, X. Wang, H. Dai, N. Liu, et al., “Agi for agriculture,” *arXiv preprint arXiv:2304.06136*, 2023.
- [420] A. Tzachor, M. Devare, C. Richards, P. Pypers, A. Ghosh, J. Koo, S. Johal, and B. King, “Large language models and agricultural extension services,” *Nature food*, vol. 4, no. 11, pp. 941–948, 2023.
- [421] T. R. Wanasinghe, R. G. Gosine, O. De Silva, G. K. Mann, L. A. James, and P. Warrian, “Unmanned aerial systems for the oil and gas industry: Overview, applications, and challenges,” *IEEE access*, vol. 8, pp. 166980–166997, 2020.
- [422] S. Z. Hassan, P. Sun, M. Gokgoz, J. Chen, D. R. Reinhart, and S. Gustitus-Graham, “Uav-based approach for municipal solid waste landfill monitoring and water ponding issue detection using sensor fusion,” *Journal of Hydroinformatics*, vol. 25, no. 6, pp. 2107–2127, 2023.
- [423] S. Hassan, M. Gokgoz, P. Sun, J. Chen, and B. Nam, “A cost-effective uav-based sensing system for waste landfill management,” in *The 9th International Conference on Water Resources and Environment Research, Virtual Conference*, 2022.
- [424] C. Ju and H. I. Son, “Multiple uav systems for agricultural applications: Control, implementation, and evaluation,” *Electronics*, vol. 7, no. 9, p. 162, 2018.
- [425] E. Alvarez-Vanhard, T. Corpetti, and T. Houet, “Uav & satellite synergies for optical remote sensing applications: A literature review,” *Science of remote sensing*, vol. 3, p. 100019, 2021.
- [426] S. Khanal, K. Kc, J. P. Fulton, S. Shearer, and E. Ozkan, “Remote sensing in agriculture—accomplishments, limitations, and opportunities,” *Remote Sensing*, vol. 12, no. 22, p. 3783, 2020.
- [427] S. Asadzadeh, W. J. de Oliveira, and C. R. de Souza Filho, “Uav-based remote sensing for the petroleum industry and environmental monitoring: State-of-the-art and perspectives,” *Journal of Petroleum Science and Engineering*, vol. 208, p. 109633, 2022.
- [428] P. Lottes, R. Khanna, J. Pfeifer, R. Siegwart, and C. Stachniss, “Uav-based crop and weed classification for smart farming,” in *2017 IEEE international conference on robotics and automation (ICRA)*, pp. 3024–3031, IEEE, 2017.
- [429] J. Anderegg, F. Tschurr, N. Kirchgessner, S. Treier, M. Schmucki, B. Streit, and A. Walter, “On-farm evaluation of uav-based aerial imagery for season-long weed monitoring under contrasting management and pedoclimatic conditions in wheat,” *Computers and Electronics in Agriculture*, vol. 204, p. 107558, 2023.
- [430] Y. Huang, K. N. Reddy, R. S. Fletcher, and D. Pennington, “Uav low-altitude remote sensing for precision weed management,” *Weed technology*, vol. 32, no. 1, pp. 2–6, 2018.
- [431] K. Ivushkin, H. Bartholomaeus, A. K. Bregt, A. Pulatov, M. H. Franceschini, H. Kramer, E. N. van Loo, V. J. Roman, and R. Finkers, “Uav based soil salinity assessment of cropland,” *Geoderma*, vol. 338, pp. 502–512, 2019.
- [432] G. Sona, D. Passoni, L. Pinto, D. Pagliari, D. Masseroni, B. Ortuan, and A. Facchi, “Uav multispectral survey to map soil and crop for precision farming applications,” *The international archives of the photogrammetry, remote sensing and spatial information sciences*, vol. 41, pp. 1023–1029, 2016.
- [433] K. Neupane and F. Baysal-Gurel, “Automatic identification and monitoring of plant diseases using unmanned aerial vehicles: A review,” *Remote Sensing*, vol. 13, no. 19, p. 3841, 2021.

- [434] D. Gao, Q. Sun, B. Hu, and S. Zhang, "A framework for agricultural pest and disease monitoring based on internet-of-things and unmanned aerial vehicles," *Sensors*, vol. 20, no. 5, p. 1487, 2020.
- [435] T. B. Shahi, C.-Y. Xu, A. Neupane, and W. Guo, "Recent advances in crop disease detection using uav and deep learning techniques," *Remote Sensing*, vol. 15, no. 9, p. 2450, 2023.
- [436] S. Gokool, M. Mahomed, R. Kunz, A. Clulow, M. Sibanda, V. Naiken, K. Chetty, and T. Mabhaudhi, "Crop monitoring in smallholder farms using unmanned aerial vehicles to facilitate precision agriculture practices: a scoping review and bibliometric analysis," *Sustainability*, vol. 15, no. 4, p. 3557, 2023.
- [437] M. Schirrmann, A. Giebel, F. Gleiniger, M. Pflanz, J. Lentschke, and K.-H. Dammer, "Monitoring agronomic parameters of winter wheat crops with low-cost uav imagery," *Remote Sensing*, vol. 8, no. 9, p. 706, 2016.
- [438] J. Bendig, A. Bolten, and G. Bareth, "Uav-based imaging for multi-temporal, very high resolution crop surface models to monitor crop growth variability, photogramm. fernerkun., 6, 551–562," 2013.
- [439] F. Toscano, C. Fiorentino, N. Capece, U. Erra, D. Travascia, A. Scopa, M. Drosos, and P. D'Antonio, "Unmanned aerial vehicle for precision agriculture: A review," *IEEE Access*, 2024.
- [440] M. Dowling and B. Lucey, "Chatgpt for (finance) research: The bananarama conjecture," *Finance Research Letters*, vol. 53, p. 103662, 2023.
- [441] X.-Y. Liu, G. Wang, and D. Zha, "Fingpt: Democratizing internet-scale data for financial large language models," *arXiv preprint arXiv:2307.10485*, 2023.
- [442] A. Zaremba and E. Demir, "Chatgpt: Unlocking the future of nlp in finance," Available at SSRN 4323643, 2023.
- [443] A. Lopez-Lira and Y. Tang, "Can chatgpt forecast stock price movements? return predictability and large language models," *arXiv preprint arXiv:2304.07619*, 2023.
- [444] Y. Yang, M. C. S. Uy, and A. Huang, "Finbert: A pretrained language model for financial communications," *arXiv preprint arXiv:2006.08097*, 2020.
- [445] D. Peskoff and B. M. Stewart, "Credible without credit: Domain experts assess generative language models," in *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pp. 427–438, 2023.
- [446] S. Agarwal, S. Alok, P. Ghosh, and S. Gupta, "Financial inclusion and alternate credit scoring for the millennials: role of big data and machine learning in fintech," *Business School, National University of Singapore Working Paper, SSRN*, vol. 3507827, 2020.
- [447] K. B. Hansen, "The virtue of simplicity: On machine learning models in algorithmic trading," *Big Data & Society*, vol. 7, no. 1, p. 2053951720926558, 2020.
- [448] Z. Lin, Z. Song, Z. Dai, and Q. V. Le, "Fingpt: Open-source financial large language models," *arXiv preprint arXiv:2306.06031*, 2023.
- [449] X. Hou, Y. Zhao, Y. Liu, Z. Yang, K. Wang, L. Li, X. Luo, D. Lo, J. Grundy, and H. Wang, "Large language models for software engineering: A systematic literature review," *arXiv preprint arXiv:2308.10620*, 2023.
- [450] M. Fraiwan and N. Khasawneh, "A review of chatgpt applications in education, marketing, software engineering, and healthcare: Benefits, drawbacks, and research directions," *arXiv preprint arXiv:2305.00237*, 2023.
- [451] D. Tiro, "The possibility of applying chatgpt (ai) for calculations in mechanical engineering," in *New Technologies, Development and Application VI: Volume 1*, pp. 313–320, Springer, 2023.
- [452] X. Wang, N. Anwer, Y. Dai, and A. Liu, "Chatgpt for design, manufacturing, and education," 2023.
- [453] S. Badini, S. Regondi, E. Frontoni, and R. Pugliese, "Assessing the capabilities of chatgpt to improve additive manufacturing troubleshooting," *Advanced Industrial and Engineering Polymer Research*, 2023.
- [454] J. V. Pavlik, "Collaborating with chatgpt: Considering the implications of generative artificial intelligence for journalism and media education," *Journalism & Mass Communication Educator*, vol. 78, no. 1, pp. 84–93, 2023.
- [455] L. Chan, L. Hogaboam, and R. Cao, "Ai in media and entertainment," in *Applied Artificial Intelligence in Business: Concepts and Cases*, pp. 305–324, Springer, 2022.
- [456] R. Lachman and M. Joffe, "Applications of artificial intelligence in media and entertainment," in *Analyzing future applications of AI, sensors, and robotics in society*, pp. 201–220, IGI Global, 2021.
- [457] J. Kirchenbauer, J. Geiping, Y. Wen, J. Katz, I. Miers, and T. Goldstein, "A watermark for large language models," *arXiv preprint arXiv:2301.10226*, 2023.
- [458] Z. Wang, "Mediagpt: A large language model target chinese media," *arXiv preprint arXiv:2307.10930*, 2023.
- [459] J. M. Pérez, D. A. Furman, L. A. Alemany, and F. Luque, "Robertuito: a pre-trained language model for social media text in spanish," *arXiv preprint arXiv:2111.09453*, 2021.
- [460] M. Abdulhai, C. Crepy, D. Valter, J. Canny, and N. Jaques, "Moral foundations of large language models," in *AAAI 2023 Workshop on Representation Learning for Responsible Human-Centric AI*, 2022.
- [461] H. Steck, L. Baltrunas, E. Elahi, D. Liang, Y. Raimond, and J. Basilio, "Deep learning for recommender systems: A netflix case study," *AI Magazine*, vol. 42, no. 3, pp. 7–18, 2021.
- [462] "Databricks - Media Entertainment Solutions." <https://www.databricks.com/solutions/industries/media-and-entertainment>. Accessed: Insert date accessed.
- [463] J. Kim, K. Xu, and K. Merrill Jr, "Man vs. machine: Human responses to an ai newscaster and the role of social presence," *The Social Science Journal*, pp. 1–13, 2022.
- [464] M. Feng, "The development of " ai" synthetic anchor in the context of artificial intelligence," *Highlights in Art and Design*, vol. 2, no. 1, pp. 38–40, 2023.
- [465] A. of the article, "Ai-generated news presenter appears in kuwait," *Al Jazeera*, April 2023.
- [466] A. of the article, "This is how ai could change the future of journalism," *Sky News*, 2023.
- [467] K. Y. Iu and V. M.-Y. Wong, "ChatGPT by OpenAI: The End of Litigation Lawyers?", 2023. Available at SSRN.
- [468] M. Ajevski, K. Barker, A. Gilbert, L. Hardie, and F. Ryan, "Chatgpt and the future of legal education and practice," *The Law Teacher*, vol. 0, no. 0, pp. 1–13, 2023.
- [469] "The legal ai you've been waiting for." <https://casetext.com/cocounsel/>. 31 July 2023].
- [470] J. Cui, Z. Li, Y. Yan, B. Chen, and L. Yuan, "Chatlaw: Open-source legal large language model with integrated external knowledge bases," *arXiv preprint arXiv:2306.16092*, 2023.
- [471] J. J. Nay, "Law Informs Code: A Legal Informatics Approach to Aligning Artificial Intelligence with Humans," *arXiv preprint*, 2022.
- [472] D. Trautmann, A. Petrova, and F. Schilder, "Legal Prompt Engineering for Multilingual Legal Judgement Prediction," *arXiv preprint*, 2022.
- [473] J. H. Choi, K. E. Hickman, A. Monahan, and D. B. Schwarcz, "ChatGPT Goes to Law School," *Journal of Legal Education*, January 2023. Forthcoming.
- [474] F. Yu, L. Quarley, and F. Schilder, "Legal prompting: Teaching a language model to think like a lawyer," 2022.
- [475] J. H. Choi, K. E. Hickman, A. Monahan, and D. B. Schwarcz, "Supra Note 7," 2023. Reference to a previously cited work.
- [476] The Guardian, "Two US Lawyers Fined for Submitting Fake Court Citations by ChatGPT," June 2023.
- [477] ABC News, "US Lawyer Uses ChatGPT to Research Case with Embarrassing Result," June 2023.
- [478] Business Standard, "US Judge Orders Lawyers Not to Use ChatGPT-drafted Content in Court," May 2023.
- [479] P. Rivas and L. Zhao, "Marketing with ChatGPT: Navigating the Ethical Terrain of GPT-Based Chatbot Technology," *AI*, vol. 4, pp. 375–384, 2023.
- [480] A. K. Kushwaha and A. K. Kar, "Language model-driven chatbot for business to address marketing and selection of products," in *Re-imaging Diffusion and Adoption of Information Technology and Systems: A Continuing Conversation: IFIP WG 8.6 International Conference on Transfer and Diffusion of IT, TDIT 2020, Tiruchirappalli, India, December 18–19, 2020, Proceedings, Part I*, pp. 16–28, Springer, 2020.
- [481] S. Verma, R. Sharma, S. Deb, and D. Maitra, "Artificial intelligence in marketing: Systematic review and future research direction," *International Journal of Information Management Data Insights*, vol. 1, no. 1, p. 100002, 2021.
- [482] C. Zielinski, M. Winker, R. Aggarwal, L. Ferris, M. Heinemann, J. Lapeña, S. Pai, and L. Citrome, "Chatbots, ChatGPT, and Scholarly Manuscripts - WAME Recommendations on ChatGPT and Chatbots in Relation to Scholarly Publications," *Afro-Egypt. J. Infect. Endem. Dis.*, vol. 13, pp. 75–79, 2023.
- [483] A. F.-B. Sun, Grace H. DNP and R.-B.-C. C. F. Hoelscher, Stephanie H. DNP, "The ChatGPT Storm and What Faculty Can Do," *Nurse Educator*, vol. 48, pp. 119–124, May/June 2023.
- [484] L. Ma and B. Sun, "Machine learning and ai in marketing – connecting computing power to human insights," *International Journal of Research in Marketing*, vol. 37, no. 3, pp. 481–504, 2020.

- [485] O. Yara, A. Brazheyev, L. Golovko, and V. Bashkatova, "Legal regulation of the use of artificial intelligence: Problems and development prospects," *European Journal of Sustainable Development*, vol. 10, p. 281, Feb. 2021.
- [486] M. Stone, E. Aravopoulou, Y. Ekinci, G. Evans, M. Hobbs, A. Labib, P. Laughlin, J. Machtynger, and L. Machtynger, "Artificial Intelligence (AI) in Strategic Marketing Decision-Making: A research agenda," *Bottom Line*, vol. 33, pp. 183–200, 2020.
- [487] E. Hermann, "Leveraging Artificial Intelligence in Marketing for Social Good—An Ethical Perspective," *J Bus Ethics*, vol. 179, pp. 43–61, 2022.
- [488] K. Jarek and G. Mazurek, "Marketing and artificial intelligence," *Central European Business Review*, vol. 8, no. 2, 2019.
- [489] Z. Liu, X. Yu, L. Zhang, Z. Wu, C. Cao, H. Dai, L. Zhao, W. Liu, D. Shen, Q. Li, et al., "Deid-gpt: Zero-shot medical text de-identification by gpt-4," *arXiv preprint arXiv:2303.11032*, 2023.
- [490] A. D. Subagja, A. M. A. Ausat, A. R. Sari, M. I. Wanof, and S. Suherlan, "Improving customer service quality in msmses through the use of chatgpt," *Jurnal Minfo Polgan*, vol. 12, no. 2, pp. 380–386, 2023.
- [491] S. Makridakis, F. Petropoulos, and Y. Kang, "Large language models: Their success and impact," *Forecasting*, vol. 5, no. 3, pp. 536–549, 2023.
- [492] K. Howell, G. Christian, P. Fomitchov, G. Kehat, J. Marzulla, L. Rolston, J. Tredup, I. Zimmerman, E. Selfridge, and J. Bradley, "The economic trade-offs of large language models: A case study," *arXiv preprint arXiv:2306.07402*, 2023.
- [493] J. Potts, D. W. Allen, C. Berg, and N. Ilyushina, "Large language models reduce agency costs," *Available at SSRN*, 2023.
- [494] P. A. Olujimi and A. Ade-Ibijola, "Nlp techniques for automating responses to customer queries: a systematic review," *Discover Artificial Intelligence*, vol. 3, no. 1, p. 20, 2023.
- [495] O. Topsakal and T. C. Akinci, "Creating large language model applications utilizing langchain: A primer on developing llm apps fast," in *International Conference on Applied Engineering and Natural Sciences*, vol. 1, pp. 1050–1056, 2023.
- [496] Q. Wu, G. Bansal, J. Zhang, Y. Wu, S. Zhang, E. Zhu, B. Li, L. Jiang, X. Zhang, and C. Wang, "Autogen: Enabling next-gen llm applications via multi-agent conversation framework," *arXiv preprint arXiv:2308.08155*, 2023.
- [497] N. Morrical, J. Tremblay, Y. Lin, S. Tyree, S. Birchfield, V. Pascucci, and I. Wald, "Nvisii: A scriptable tool for photorealistic image generation," *arXiv preprint arXiv:2105.13962*, 2021.
- [498] Y. Li, M. Min, D. Shen, D. Carlson, and L. Carin, "Video generation from text," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 32, 2018.
- [499] G. Franceschelli and M. Musolesi, "On the creativity of large language models," *arXiv preprint arXiv:2304.00008*, 2023.
- [500] A. Conneau and G. Lample, "Cross-lingual language model pretraining," *Advances in neural information processing systems*, vol. 32, 2019.
- [501] M. J. Ali, "Chatgpt and lacrimal drainage disorders: performance and scope of improvement," *Ophthalmic Plastic and Reconstructive Surgery*, vol. 39, no. 3, p. 221, 2023.
- [502] X.-Q. Dao, "Performance comparison of large language models on vnhsgc english dataset: Openai chatgpt, microsoft bing chat, and google bard," *arXiv preprint arXiv:2307.02288*, 2023.
- [503] J. Rudolph, S. Tan, and S. Tan, "War of the chatbots: Bard, bing chat, chatgpt, ernie and beyond. the new ai gold rush and its impact on higher education," *Journal of Applied Learning and Teaching*, vol. 6, no. 1, 2023.
- [504] I. Ahmed, M. Kajol, U. Hasan, P. P. Datta, A. Roy, and M. R. Reza, "Chatgpt vs. bard: A comparative study," *UMBC Student Collection*, 2023.
- [505] R. Thoppilan, D. De Freitas, J. Hall, N. Shazeer, A. Kulshreshtha, H.-T. Cheng, A. Jin, T. Bos, L. Baker, Y. Du, et al., "Lamda: Language models for dialog applications," *arXiv preprint arXiv:2201.08239*, 2022.
- [506] X. Amatriain, "Transformer models: an introduction and catalog," *arXiv preprint arXiv:2302.07730*, 2023.
- [507] Y. Hu, I. Ameer, X. Zuo, X. Peng, Y. Zhou, Z. Li, Y. Li, J. Li, X. Jiang, and H. Xu, "Zero-shot clinical entity recognition using chatgpt," *arXiv preprint arXiv:2303.16416*, 2023.
- [508] T. Wu, S. He, J. Liu, S. Sun, K. Liu, Q.-L. Han, and Y. Tang, "A brief overview of chatgpt: The history, status quo and potential future development," *IEEE/CAA Journal of Automatica Sinica*, vol. 10, no. 5, pp. 1122–1136, 2023.
- [509] M. Abdullah, A. Madain, and Y. Jararweh, "Chatgpt: Fundamentals, applications and social impacts," in *2022 Ninth International Conference on Social Networks Analysis, Management and Security (SNAMS)*, pp. 1–8, IEEE, 2022.
- [510] K. Arulkumaran, M. P. Deisenroth, M. Brundage, and A. A. Bharath, "Deep reinforcement learning: A brief survey," *IEEE Signal Processing Magazine*, vol. 34, no. 6, pp. 26–38, 2017.
- [511] M. Kiran, P. Murphy, I. Monga, J. Dugan, and S. S. Baveja, "Lambda architecture for cost-effective batch and speed big data processing," in *2015 IEEE international conference on big data (big data)*, pp. 2785–2792, IEEE, 2015.
- [512] D. Xuan-Quy, L. Ngoc-Bich, V. The-Duy, P. Xuan-Dung, N. Bac-Bien, N. Van-Tien, N. Thi-My-Thanh, and N. Hong-Phuoc, "Vnhsgc: Vietnamese high school graduation examination dataset for large language models," *arXiv preprint arXiv:2305.12199*, 2023.
- [513] H. Trng, "Chatgpt in education-a global and vietnamese research overview," 2023.
- [514] S. Bhardwaz and J. Kumar, "An extensive comparative analysis of chatbot technologies - chatgpt, google bard and microsoft bing," in *2023 2nd International Conference on Applied Artificial Intelligence and Computing (ICAAIC)*, pp. 673–679, 2023.
- [515] B. Campello de Souza, A. Serrano de Andrade Neto, and A. Roazzi, "Are the new ais smart enough to steal your job? iq scores for chatgpt, microsoft bing, google bard and quora poe," *IQ Scores for ChatGPT, Microsoft Bing, Google Bard and Quora Poe (April 7, 2023)*, 2023.
- [516] K. M. Caramancion, "News verifiers showdown: A comparative performance evaluation of chatgpt 3.5, chatgpt 4.0, bing ai, and bard in news fact-checking," *arXiv preprint arXiv:2306.17176*, 2023.
- [517] L. Zhang, "Four tax questions for chatgpt and other language models," 2023.
- [518] M. King, "Gpt-4 aligns with the new liberal party, while other large language models refuse to answer political questions."
- [519] L. Yang, Z. Zhang, Y. Song, S. Hong, R. Xu, Y. Zhao, Y. Shao, W. Zhang, B. Cui, and M.-H. Yang, "Diffusion models: A comprehensive survey of methods and applications," *arXiv preprint arXiv:2209.00796*, 2022.
- [520] J. Ho, C. Saharia, W. Chan, D. J. Fleet, M. Norouzi, and T. Salimans, "Cascaded diffusion models for high fidelity image generation," *The Journal of Machine Learning Research*, vol. 23, no. 1, pp. 2249–2281, 2022.
- [521] W. Wang, R. Metzler, and A. G. Cherstvy, "Anomalous diffusion, aging, and nonergodicity of scaled brownian motion with fractional gaussian noise: Overview of related experimental observations and models," *Physical Chemistry Chemical Physics*, vol. 24, no. 31, pp. 18482–18504, 2022.
- [522] F.-A. Croitoru, V. Hondru, R. T. Ionescu, and M. Shah, "Diffusion models in vision: A survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–20, 2023.
- [523] J. Huggins and J. Zou, "Quantifying the accuracy of approximate diffusions and markov chains," in *Artificial Intelligence and Statistics*, pp. 382–391, PMLR, 2017.
- [524] W. Feller, "On the theory of stochastic processes, with particular reference to applications, p 403–432," 1949.
- [525] J. R. Hershey and P. A. Olsen, "Approximating the kullback leibler divergence between gaussian mixture models," in *2007 IEEE International Conference on Acoustics, Speech and Signal Processing-ICASSP'07*, vol. 4, pp. IV-317, IEEE, 2007.
- [526] J. Kukačka, V. Golkov, and D. Cremers, "Regularization for deep learning: A taxonomy," *arXiv preprint arXiv:1710.10686*, 2017.
- [527] A. Sedghi, L. J. O'Donnell, T. Kapur, E. Learned-Miller, P. Mousavi, and W. M. Wells III, "Image registration: Maximum likelihood, minimum entropy and deep learning," *Medical image analysis*, vol. 69, p. 101939, 2021.
- [528] X. Pan, A. Tewari, T. Leimkühler, L. Liu, A. Meka, and C. Theobalt, "Drag your gan: Interactive point-based manipulation on the generative image manifold," in *ACM SIGGRAPH 2023 Conference Proceedings*, pp. 1–11, 2023.
- [529] P. Dhariwal and A. Nichol, "Diffusion models beat gans on image synthesis," *Advances in neural information processing systems*, vol. 34, pp. 8780–8794, 2021.
- [530] A. Nichol, P. Dhariwal, A. Ramesh, P. Shyam, P. Mishkin, B. McGrew, I. Sutskever, and M. Chen, "Glide: Towards photorealistic image generation and editing with text-guided diffusion models," *arXiv preprint arXiv:2112.10741*, 2021.
- [531] C. Saharia, W. Chan, H. Chang, C. Lee, J. Ho, T. Salimans, D. Fleet, and M. Norouzi, "Palette: Image-to-image diffusion models," in *ACM SIGGRAPH 2022 Conference Proceedings*, pp. 1–10, 2022.

- [532] G. Kim, T. Kwon, and J. C. Ye, "Diffusionclip: Text-guided diffusion models for robust image manipulation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2426–2435, 2022.
- [533] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models," *Advances in neural information processing systems*, vol. 33, pp. 6840–6851, 2020.
- [534] S. AI, "Stablediffusion2.1 release." <https://stability.ai/blog/stablediffusion2-1-release7-dec-2022>, 2022. Accessed on August 3, 2023.
- [535] C. Ju, T. Han, K. Zheng, Y. Zhang, and W. Xie, "Prompting visual-language models for efficient video understanding," in *European Conference on Computer Vision*, pp. 105–124, Springer, 2022.
- [536] H. Lin, A. Zala, J. Cho, and M. Bansal, "Videodirectorgpt: Consistent multi-scene video generation via llm-guided planning," *arXiv preprint arXiv:2309.15091*, 2023.
- [537] A. Munoz, M. Zolfaghari, M. Argus, and T. Brox, "Temporal shift gan for large scale video generation," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 3179–3188, 2021.
- [538] M.-Y. Liu, X. Huang, J. Yu, T.-C. Wang, and A. Mallya, "Generative adversarial networks for image and video synthesis: Algorithms and applications," *Proceedings of the IEEE*, vol. 109, no. 5, pp. 839–862, 2021.
- [539] U. Singer, A. Polyak, T. Hayes, X. Yin, J. An, S. Zhang, Q. Hu, H. Yang, O. Ashual, O. Gafni, et al., "Make-a-video: Text-to-video generation without text-video data," *arXiv preprint arXiv:2209.14792*, 2022.
- [540] T. Reddy, R. Williams, and C. Breazeal, "Text classification for ai education.," in *SIGCSE*, p. 1381, 2021.
- [541] E. Loper and S. Bird, "Nltk: The natural language toolkit," *arXiv preprint cs/0205028*, 2002.
- [542] Y. Vasiliev, *Natural language processing with Python and spaCy: A practical introduction*. No Starch Press, 2020.
- [543] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz, et al., "Huggingface's transformers: State-of-the-art natural language processing," *arXiv preprint arXiv:1910.03771*, 2019.
- [544] Y. Liu, K. Zhang, Y. Li, Z. Yan, C. Gao, R. Chen, Z. Yuan, Y. Huang, H. Sun, J. Gao, et al., "Sora: A review on background, technology, limitations, and opportunities of large vision models," *arXiv preprint arXiv:2402.17177*, 2024.
- [545] J. Pachouly, S. Ahirrao, K. Kotecha, G. Selvachandran, and A. Abraham, "A systematic literature review on software defect prediction using artificial intelligence: Datasets, data validation methods, approaches, and tools," *Engineering Applications of Artificial Intelligence*, vol. 111, p. 104773, 2022.
- [546] E. A. Van Dis, J. Bollen, W. Zuidema, R. van Rooij, and C. L. Bockting, "Chatgpt: five priorities for research," *Nature*, vol. 614, no. 7947, pp. 224–226, 2023.
- [547] K. Nguyen-Trung, A. K. Saeri, and S. Kaufman, "Applying chatgpt and ai-powered tools to accelerate evidence reviews," 2023.
- [548] N. Gleason, "Chatgpt and the rise of ai writers: How should higher education respond?," *Times Higher Education*, 2022.
- [549] G. Cooper, "Examining science education in chatgpt: An exploratory study of generative artificial intelligence," *Journal of Science Education and Technology*, vol. 32, pp. 444–452, 2023.
- [550] Z. Epstein, A. Hertzmann, I. of Human Creativity, M. Akten, H. Farid, J. Fjeld, M. R. Frank, M. Groh, L. Herman, N. Leach, et al., "Art and the science of generative ai," *Science*, vol. 380, no. 6650, pp. 1110–1111, 2023.
- [551] L. Skavronskaia, A. H. Hadinejad, and D. Cotterell, "Reversing the threat of artificial intelligence to opportunity: a discussion of chatgpt in tourism education," *Journal of Teaching in Travel & Tourism*, vol. 23, no. 2, pp. 253–258, 2023.
- [552] J. Huang and K. C.-C. Chang, "Citation: A key to building responsible and accountable large language models," *arXiv preprint arXiv:2307.02185*, 2023.
- [553] B. Yetişiren, I. Özsoy, M. Ayerdem, and E. Tüzün, "Evaluating the code quality of ai-assisted code generation tools: An empirical study on github copilot, amazon codewhisperer, and chatgpt," *arXiv preprint arXiv:2304.10778*, 2023.
- [554] A. M. Dakhel, V. Majdinasab, A. Nikanjam, F. Khomh, M. C. Desmarais, and Z. M. J. Jiang, "Github copilot ai pair programmer: Asset or liability?," *Journal of Systems and Software*, vol. 203, p. 111734, 2023.
- [555] F. G. e. a. Baptiste Rozière, Jonas Gehring, "Code llama: Open foundation models for code," *arXiv preprint*, 2023.
- [556] T. Calò and L. De Russis, "Leveraging large language models for end-user website generation," in *International Symposium on End User Development*, pp. 52–61, Springer, 2023.
- [557] Y. Shen, K. Song, X. Tan, D. Li, W. Lu, and Y. Zhuang, "Hugginggpt: Solving ai tasks with chatgpt and its friends in huggingface," *arXiv preprint arXiv:2303.17580*, 2023.
- [558] O. Thawkar, A. Shaker, S. S. Mullappilly, H. Cholakkal, R. M. Anwer, S. Khan, J. Laaksonen, and F. S. Khan, "Xraygpt: Chest radiographs summarization using medical vision-language models," *arXiv preprint arXiv:2306.07971*, 2023.
- [559] A. Rayhan and D. Gross, "Searchgpt: Revolutionizing information retrieval with advanced language models."
- [560] Y. Ding, W. Fan, L. Ning, S. Wang, H. Li, D. Yin, T.-S. Chua, and Q. Li, "A survey on rag meets llms: Towards retrieval-augmented large language models," *arXiv preprint arXiv:2405.06211*, 2024.
- [561] Y. Luo, J. Zhang, S. Fan, K. Yang, Y. Wu, M. Qiao, and Z. Nie, "Biomedgpt: Open multimodal generative pre-trained transformer for biomedicine," *arXiv preprint arXiv:2308.09442*, 2023.
- [562] K. Kheiri and H. Karimi, "Sentimentgpt: Exploiting gpt for advanced sentiment analysis and its departure from current machine learning," *arXiv preprint arXiv:2307.10234*, 2023.
- [563] S. Whitfield and M. A. Hofmann, "Elicit: Ai literature review research assistant," *Public Services Quarterly*, vol. 19, no. 3, pp. 201–207, 2023.
- [564] L. Cao, V. Buchner, Z. Senane, and F. Yang, "Introducing genception for multimodal llm benchmarking: You may bypass annotations," in *Proceedings of the 4th Workshop on Trustworthy Natural Language Processing (TrustNLP 2024)*, pp. 196–201, 2024.
- [565] H. Chase, "Langchain, 10 2022," URL <https://github.com/hwchase17/langchain>.
- [566] M. R. Chavez, T. S. Butler, P. Rekawek, H. Heo, and W. L. Kinzler, "Chat generative pre-trained transformer: why we should embrace this technology," *American Journal of Obstetrics and Gynecology*, 2023.
- [567] A. Chowdhery, S. Narang, J. Devlin, M. Bosma, G. Mishra, A. Roberts, P. Barham, H. W. Chung, C. Sutton, S. Gehrmann, et al., "Palm: Scaling language modeling with pathways," *arXiv preprint arXiv:2204.02311*, year=2022.
- [568] T. Teubner, C. M. Flath, C. Weinhardt, W. van der Aalst, and O. Hinz, "Welcome to the era of chatgpt et al. the prospects of large language models," *Business & Information Systems Engineering*, pp. 1–7, 2023.
- [569] I. Poola and V. Božić, "These plug-ins revolutionize chatgpt into an everything app-chatgpt's plugin store," 2023.
- [570] C. Xu, Y. Xu, S. Wang, Y. Liu, C. Zhu, and J. McAuley, "Small models are valuable plug-ins for large language models," *arXiv preprint arXiv:2305.08848*, 2023.
- [571] H. Gimpel, K. Hall, S. Decker, T. Eymann, L. Lämmermann, A. Mädche, M. Röglinger, C. Ruiner, M. Schoch, M. Schoop, et al., "Unlocking the power of generative ai models and systems such as gpt-4 and chatgpt for higher education: A guide for students and lecturers," tech. rep., Hohenheim Discussion Papers in Business, Economics and Social Sciences, 2023.
- [572] K.-L. Liu, W.-J. Li, and M. Guo, "Emoticon smoothed language models for twitter sentiment analysis," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 26, pp. 1678–1684, 2012.
- [573] M. M. Al-Jefri and S. A. Mahmoud, "Context-sensitive arabic spell checker using context words and n-gram language models," in *2013 Taibah University International Conference on Advances in Information Technology for the Holy Quran and Its Sciences*, pp. 258–263, IEEE, 2013.
- [574] Z. Jiang, J. Araki, H. Ding, and G. Neubig, "How can we know when language models know? on the calibration of language models for question answering," *Transactions of the Association for Computational Linguistics*, vol. 9, pp. 962–977, 2021.
- [575] L. Wang, W. Zhao, Z. Wei, and J. Liu, "Simkgc: Simple contrastive knowledge graph completion with pre-trained language models," *arXiv preprint arXiv:2203.02167*, 2022.
- [576] H. J. Dolfig and I. L. Hetherington, "Incremental language models for speech recognition using finite-state transducers," in *IEEE Workshop on Automatic Speech Recognition and Understanding, 2001. ASRU'01.*, pp. 194–197, IEEE, 2001.
- [577] R. Mao, Q. Liu, K. He, W. Li, and E. Cambria, "The biases of pre-trained language models: An empirical study on prompt-based sentiment analysis and emotion detection," *IEEE Transactions on Affective Computing*, 2022.
- [578] J. Zhang, R. Xie, Y. Hou, W. X. Zhao, L. Lin, and J.-R. Wen, "Recommendation as instruction following: A large language model empowered recommendation approach," *arXiv preprint arXiv:2305.07001*, 2023.

- [579] Y. Liu, T. Han, S. Ma, J. Zhang, Y. Yang, J. Tian, H. He, A. Li, M. He, Z. Liu, et al., “Summary of chatgpt/gpt-4 research and perspective towards the future of large language models,” *arXiv preprint arXiv:2304.01852*, 2023.
- [580] X. He, X. Shen, Z. Chen, M. Backes, and Y. Zhang, “Mgtbench: Benchmarking machine-generated text detection,” *arXiv preprint arXiv:2303.14822*, 2023.
- [581] M. T. I. Khondaker, A. Waheed, E. M. B. Nagoudi, and M. Abdul-Mageed, “Gptaraeval: A comprehensive evaluation of chatgpt on arabic nlp,” *arXiv preprint arXiv:2305.14976*, 2023.
- [582] J. Kim, J. H. Lee, S. Kim, J. Park, K. M. Yoo, S. J. Kwon, and D. Lee, “Memory-efficient fine-tuning of compressed large language models via sub-4-bit integer quantization,” *arXiv preprint arXiv:2305.14152*, 2023.
- [583] S. Arora, B. Yang, S. Eyuboglu, A. Narayan, A. Hojel, I. Trummer, and C. Ré, “Language models enable simple systems for generating structured views of heterogeneous data lakes,” *arXiv preprint arXiv:2304.09433*, 2023.
- [584] S. R. Bowman, “Eight things to know about large language models,” *arXiv preprint arXiv:2304.00612*, 2023.
- [585] Y. Fu, H. Peng, T. Khot, and M. Lapata, “Improving language model negotiation with self-play and in-context learning from ai feedback,” *arXiv preprint arXiv:2305.10142*, 2023.
- [586] H. Matsumi, D. Hallinan, D. Dimitrova, E. Kosta, and P. De Hert, *Data Protection and Privacy, Volume 15: In Transitional Times*. Bloomsbury Publishing, 2023.
- [587] P. Hacker, A. Engel, and M. Mauer, “Regulating chatgpt and other large generative ai models,” in *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*, pp. 1112–1123, 2023.
- [588] A. Gokul, “Llms and ai: Understanding its reach and impact,” 2023.
- [589] S. A. Khowaja, P. Khuwaja, and K. Dev, “Chatgpt needs spade (sustainability, privacy, digital divide, and ethics) evaluation: A review,” *arXiv preprint arXiv:2305.03123*, 2023.
- [590] A. Chan, H. Bradley, and N. Rajkumar, “Reclaiming the digital commons: A public data trust for training data,” *arXiv preprint arXiv:2303.09001*, 2023.
- [591] W. H. Deng, B. Guo, A. Devrio, H. Shen, M. Eslami, and K. Holstein, “Understanding practices, challenges, and opportunities for user-engaged algorithm auditing in industry practice,” in *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, pp. 1–18, 2023.
- [592] M. Kraus, J. A. Bingler, M. Leippold, T. Schimanski, C. C. Senni, D. Stammbach, S. A. Vaghefi, and N. Webersinke, “Enhancing large language models with climate resources,” *arXiv preprint arXiv:2304.00116*, 2023.
- [593] E. Agathokleous, C. J. Saitanis, C. Fang, and Z. Yu, “Use of chatgpt: What does it mean for biology and environmental science?,” *Science of The Total Environment*, p. 164154, 2023.
- [594] Y. Shen, L. Heacock, J. Elias, K. D. Hentel, B. Reig, G. Shih, and L. Moy, “Chatgpt and other large language models are double-edged swords,” 2023.
- [595] Z. Sun, Y. Shen, Q. Zhou, H. Zhang, Z. Chen, D. Cox, Y. Yang, and C. Gan, “Principle-driven self-alignment of language models from scratch with minimal human supervision,” *arXiv preprint arXiv:2305.03047*, 2023.
- [596] J. Kaddour, “The minipile challenge for data-efficient language models,” *arXiv preprint arXiv:2304.08442*, 2023.
- [597] J. Yang, H. Jin, R. Tang, X. Han, Q. Feng, H. Jiang, B. Yin, and X. Hu, “Harnessing the power of llms in practice: A survey on chatgpt and beyond,” *arXiv preprint arXiv:2304.13712*, 2023.
- [598] Z. Chen, L. Cao, S. Madden, J. Fan, N. Tang, Z. Gu, Z. Shang, C. Liu, M. Cafarella, and T. Kraska, “Seed: Simple, efficient, and effective data management via large language models,” *arXiv preprint arXiv:2310.00749*, 2023.
- [599] J. Yuan, R. Tang, X. Jiang, and X. Hu, “Llm for patient-trial matching: Privacy-aware data augmentation towards better performance and generalizability,” *arXiv preprint arXiv:2303.16756*, 2023.
- [600] Y. Wang, “Deciphering the enigma: A deep dive into understanding and interpreting llm outputs,” 2023.
- [601] J. V. Lochter, R. M. Silva, and T. A. Almeida, “Deep learning models for representing out-of-vocabulary words,” in *Brazilian Conference on Intelligent Systems*, pp. 418–434, Springer, 2020.
- [602] Y. Gu, R. Tinn, H. Cheng, M. Lucas, N. Usuyama, X. Liu, T. Naumann, J. Gao, and H. Poon, “Domain-specific language model pretraining for biomedical natural language processing,” *ACM Transactions on Computing for Healthcare (HEALTH)*, vol. 3, no. 1, pp. 1–23, 2021.
- [603] X. Wang, X. Li, Z. Yin, Y. Wu, and L. Jia, “Emotional intelligence of large language models,” *arXiv preprint arXiv:2307.09042*, 2023.
- [604] S. Doddapaneni, G. Ramesh, M. M. Khapra, A. Kunchukuttan, and P. Kumar, “A primer on pretrained multilingual language models,” *arXiv preprint arXiv:2107.00676*, 2021.
- [605] D. Schuurmans, “Memory augmented large language models are computationally universal,” *arXiv preprint arXiv:2301.04589*, 2023.
- [606] Y. Wolf, N. Wies, Y. Levine, and A. Shashua, “Fundamental limitations of alignment in large language models,” *arXiv preprint arXiv:2304.11082*, 2023.
- [607] R. Tang, Y.-N. Chuang, and X. Hu, “The science of detecting llm-generated texts,” *arXiv preprint arXiv:2303.07205*, 2023.
- [608] F. Ufuk, “The role and limitations of large language models such as chatgpt in clinical settings and medical journalism,” *Radiology*, vol. 307, no. 3, p. e230276, 2023.
- [609] R. Bhayana, S. Krishna, and R. R. Bleakney, “Performance of chatgpt on a radiology board-style examination: Insights into current strengths and limitations,” *Radiology*, p. 230582, 2023.
- [610] C.-H. Chiang and H.-y. Lee, “Can large language models be an alternative to human evaluations?,” *arXiv preprint arXiv:2305.01937*, 2023.
- [611] X. Yang, Y. Li, X. Zhang, H. Chen, and W. Cheng, “Exploring the limits of chatgpt for query or aspect-based text summarization,” *arXiv preprint arXiv:2302.08081*, 2023.
- [612] J. Kaplan, S. McCandlish, T. Henighan, T. B. Brown, B. Chess, R. Child, S. Gray, A. Radford, J. Wu, and D. Amodei, “Scaling laws for neural language models,” *arXiv preprint arXiv:2001.08361*, 2020.
- [613] J. Kaddour, J. Harris, M. Mozes, H. Bradley, R. Raileanu, and R. McHardy, “Challenges and applications of large language models,” *arXiv preprint arXiv:2307.10169*, 2023.
- [614] T. Fujii, K. Shibata, A. Yamaguchi, T. Morishita, and Y. Sogawa, “How do different tokenizers perform on downstream tasks in scriptio continua languages?: A case study in japanese,” *arXiv preprint arXiv:2306.09572*, 2023.
- [615] R. Schwartz, J. Dodge, N. A. Smith, and O. Etzioni, “Green ai,” *Communications of the ACM*, vol. 63, no. 12, pp. 54–63, 2020.
- [616] X. L. Li and P. Liang, “Prefix-tuning: Optimizing continuous prompts for generation,” *arXiv preprint arXiv:2101.00190*, 2021.
- [617] H. Laurençon, L. Saulnier, T. Wang, C. Akiki, A. Villanova del Moral, T. Le Scao, L. Von Werra, C. Mou, E. González Ponferrada, H. Nguyen, et al., “The bigscience roots corpus: A 1.6 tb composite multilingual dataset,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 31809–31826, 2022.
- [618] N. Houlsby, A. Giurgiu, S. Jastrzebski, B. Morrone, Q. De Laroussilhe, A. Gesmundo, M. Attaryan, and S. Gelly, “Parameter-efficient transfer learning for nlp,” in *International Conference on Machine Learning*, pp. 2790–2799, PMLR, 2019.
- [619] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen, “Lora: Low-rank adaptation of large language models,” *arXiv preprint arXiv:2106.09685*, 2021.
- [620] Y. Chen, S. Qian, H. Tang, X. Lai, Z. Liu, S. Han, and J. Jia, “Longlora: Efficient fine-tuning of long-context large language models,” *arXiv preprint arXiv:2309.12307*, 2023.
- [621] T. Dettmers, A. Pagnoni, A. Holtzman, and L. Zettlemoyer, “Qlora: Efficient finetuning of quantized llms,” *arXiv preprint arXiv:2305.14314*, 2023.
- [622] M. Pagliardini, D. Paliotta, M. Jaggi, and F. Fleuret, “Faster causal attention over large sequences through sparse flash attention,” *arXiv preprint arXiv:2306.01160*, 2023.
- [623] S. Liu and Z. Wang, “Ten lessons we have learned in the new “sparseland”: A short handbook for sparse neural network researchers,” *arXiv preprint arXiv:2302.02596*, 2023.
- [624] L. Chen, M. Zaharia, and J. Zou, “Frugalgpt: How to use large language models while reducing cost and improving performance,” *arXiv preprint arXiv:2305.05176*, 2023.
- [625] S. Chen, S. Wong, L. Chen, and Y. Tian, “Extending context window of large language models via positional interpolation,” *arXiv preprint arXiv:2306.15595*, 2023.
- [626] R. Li, J. Su, C. Duan, and S. Zheng, “Linear attention mechanism: An efficient attention for semantic segmentation,” *arXiv preprint arXiv:2007.14902*, 2020.
- [627] M. Guo, J. Ainslie, D. Uthus, S. Ontanon, J. Ni, Y.-H. Sung, and Y. Yang, “Longt5: Efficient text-to-text transformer for long sequences,” *arXiv preprint arXiv:2112.07916*, 2021.
- [628] X. Ma, X. Kong, S. Wang, C. Zhou, J. May, H. Ma, and L. Zettlemoyer, “Luna: Linear unified nested attention,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 2441–2453, 2021.

- [629] B. Peng, E. Alcaide, Q. Anthony, A. Albalak, S. Arcadinho, H. Cao, X. Cheng, M. Chung, M. Grella, K. K. GV, et al., “Rwkv: Reinventing rnns for the transformer era,” *arXiv preprint arXiv:2305.13048*, 2023.
- [630] T. Dao, D. Y. Fu, K. K. Saab, A. W. Thomas, A. Rudra, and C. Ré, “Hungry hungry hippos: Towards language modeling with state space models,” *arXiv preprint arXiv:2212.14052*, 2022.
- [631] Y. Yao, P. Wang, B. Tian, S. Cheng, Z. Li, S. Deng, H. Chen, and N. Zhang, “Editing large language models: Problems, methods, and opportunities,” *arXiv preprint arXiv:2305.13172*, 2023.
- [632] T. Schick and H. Schütze, “It’s not just size that matters: Small language models are also few-shot learners,” *arXiv preprint arXiv:2009.07118*, 2020.
- [633] Y. Xiong, Z. Zeng, R. Chakraborty, M. Tan, G. Fung, Y. Li, and V. Singh, “Nyströmformer: A nyström-based algorithm for approximating self-attention,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, pp. 14138–14148, 2021.
- [634] C. Wu, M. Che, and H. Yan, “The CUR Decomposition of Self-attention Matrices in Vision Transformers,” *TechRxiv*, August 01 2024.
- [635] P. Schramowski, C. Turan, N. Andersen, C. A. Rothkopf, and K. Kersting, “Large pre-trained language models contain human-like biases of what is right and wrong to do,” *Nature Machine Intelligence*, vol. 4, no. 3, pp. 258–268, 2022.
- [636] Y. Yu, Y. Zhuang, J. Zhang, Y. Meng, A. Ratner, R. Krishna, J. Shen, and C. Zhang, “Large language model as attributed training data generator: A tale of diversity and bias,” *arXiv preprint arXiv:2306.15895*, 2023.
- [637] D. Oba, M. Kaneko, and D. Bollegala, “In-contextual bias suppression for large language models,” *arXiv preprint arXiv:2309.07251*, 2023.
- [638] N. Gillibrand and C. Draper, “Informational sovereignty: A new framework for ai regulation,” *Gillibrand, Nicky, and Chris Draper. “Informational Sovereignty: A New Framework For AI Regulation” (July 17, 2023)*, 2023.
- [639] Z. Ji, N. Lee, R. Frieske, T. Yu, D. Su, Y. Xu, E. Ishii, Y. J. Bang, A. Madotto, and P. Fung, “Survey of hallucination in natural language generation,” *ACM Computing Surveys*, vol. 55, no. 12, pp. 1–38, 2023.
- [640] J. Greene, “Will ChatGPT Make Lawyers Obsolete? (Hint: Be Afraid),” *Reuters*, December 2022.
- [641] T. McCoy, E. Pavlick, and T. Linzen, “Right for the wrong reasons: Diagnosing syntactic heuristics in natural language inference,” in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, (Florence, Italy), pp. 3428–3448, Association for Computational Linguistics, July 2019.
- [642] J. Weston, E. Dinan, and A. Miller, “Retrieve and refine: Improved sequence generation models for dialogue,” in *Proceedings of the 2018 EMNLP Workshop SCAI: The 2nd International Workshop on Search-Oriented Conversational AI*, (Brussels, Belgium), pp. 87–92, Association for Computational Linguistics, Oct. 2018.
- [643] J. Thorne, A. VLachos, C. Christodoulopoulos, and A. Mittal, “FEVER: a large-scale dataset for fact extraction and VERification,” in *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, (New Orleans, Louisiana), pp. 809–819, Association for Computational Linguistics, June 2018.
- [644] D. V. Hada and S. K. Shevade, “Replug: Explainable recommendation using plug-and-play language model,” in *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 81–91, 2021.
- [645] Y. Gao, T. Sheng, Y. Xiang, Y. Xiong, H. Wang, and J. Zhang, “Chatrec: Towards interactive and explainable llms-augmented recommender system,” *arXiv preprint arXiv:2303.14524*, year=2023.
- [646] A. Uchendu, *REVERSE TURING TEST IN THE AGE OF DEEPFAKE TEXTS*. PhD thesis, The Pennsylvania State University, 2023.
- [647] E. M. Bonsu and D. Baffour-Koduah, “From the consumers’ side: Determining students’ perception and intention to use chatgpt in ghanaiian higher education,” *Journal of Education, Society & Multiculturalism*, vol. 4, no. 1, pp. 1–29, 2023.
- [648] Q. V. Liao and J. W. Vaughan, “AI Transparency in the Age of LLMs: A human-centered research roadmap,” *arXiv preprint arXiv:2306.01941*, year=2023.
- [649] G. Vilone and L. Longo, “Notions of explainability and evaluation approaches for explainable artificial intelligence,” *Information Fusion*, vol. 76, pp. 89–106, 2021.
- [650] N. M. Deshpande, S. Gite, B. Pradhan, and M. E. Assiri, “Explainable artificial intelligence—a new step towards the trust in medical diagnosis with ai frameworks: A review,” *Comput. Model. Eng. Sci.*, vol. 133, pp. 1–30, 2022.
- [651] H. Zhao, J. Jia, and V. Koltun, “Exploring self-attention for image recognition,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10076–10085, 2020.
- [652] D. Shin, “The effects of explainability and causability on perception, trust, and acceptance: Implications for explainable ai,” *International Journal of Human-Computer Studies*, vol. 146, p. 102551, 2021.
- [653] K. Valmeeakam, A. Olmo, S. Sreedharan, and S. Kambhampati, “Large language models still can’t plan (a benchmark for llms on planning and reasoning about change),” *arXiv preprint arXiv:2206.10498*, 2022.
- [654] N. Bian, X. Han, L. Sun, H. Lin, Y. Lu, and B. He, “Chatgpt is a knowledgeable but inexperienced solver: An investigation of commonsense problem in large language models,” *arXiv preprint arXiv:2303.16421*, 2023.
- [655] Y. LeCun, “A path towards autonomous machine intelligence version 0.9. 2, 2022-06-27,” *Open Review*, vol. 62, 2022.
- [656] M. Hardy, I. Sucholutsky, B. Thompson, and T. Griffiths, “Large language models meet cognitive science: Llms as tools, models, and participants,” in *Proceedings of the annual meeting of the cognitive science society*, 2023.
- [657] A. Wan, E. Wallace, S. Shen, and D. Klein, “Poisoning language models during instruction tuning,” *arXiv preprint arXiv:2305.00944*, 2023.
- [658] H. Li, D. Guo, W. Fan, M. Xu, and Y. Song, “Multi-step jailbreaking privacy attacks on chatgpt,” *arXiv preprint arXiv:2304.05197*, 2023.
- [659] F. Perez and I. Ribeiro, “Ignore previous prompt: Attack techniques for language models,” *arXiv preprint arXiv:2211.09527*, 2022.
- [660] Y. Liu, G. Deng, Y. Li, K. Wang, T. Zhang, Y. Liu, H. Wang, Y. Zheng, and Y. Liu, “Prompt injection attack against llm-integrated applications,” *arXiv preprint arXiv:2306.05499*, 2023.
- [661] C. Zhang, C. Zhang, T. Kang, D. Kim, S.-H. Bae, and I. S. Kweon, “Attack-sam: Towards evaluating adversarial robustness of segment anything model,” *arXiv preprint arXiv:2305.00866*, 2023.
- [662] H. Chen, F. Jiao, X. Li, C. Qin, M. Ravaut, R. Zhao, C. Xiong, and S. Joty, “Chatgpt’s one-year anniversary: are open-source large language models catching up?,” *arXiv preprint arXiv:2311.16989*, 2023.
- [663] K. Nelson, G. Corbin, M. Anania, M. Kovacs, J. Tobias, and M. Blowers, “Evaluating model drift in machine learning algorithms,” in *2015 IEEE Symposium on Computational Intelligence for Security and Defense Applications (CISDA)*, pp. 1–8, IEEE, 2015.
- [664] Y. Zhou, A. I. Muresanu, Z. Han, K. Paster, S. Pitis, H. Chan, and J. Ba, “Large language models are human-level prompt engineers,” *arXiv preprint arXiv:2211.01910*, 2022.
- [665] F. Huang, H. Kwak, and J. An, “Is chatgpt better than human annotators? potential and limitations of chatgpt in explaining implicit hate speech,” *arXiv preprint arXiv:2302.07736*, 2023.
- [666] G. Fergusson, C. Fitzgerald, C. Frascella, M. Iorio, T. McBrien, C. Schroeder, B. Winters, and E. Zhou, “Contributions by,”
- [667] P. Li, J. Yang, M. A. Islam, and S. Ren, “Making ai less” thirsty”: Uncovering and addressing the secret water footprint of ai models,” *arXiv preprint arXiv:2304.03271*, 2023.
- [668] D. Patterson, J. Gonzalez, Q. Le, C. Liang, L.-M. Munguia, D. Rothchild, D. So, M. Texier, and J. Dean, “Carbon emissions and large neural network training,” 2021.
- [669] S. Biswas, “Potential use of chat gpt in global warming,” *Ann Biomed Eng*, vol. 51, pp. 1126–1127, 2023.
- [670] Z. Yao, Y. Lum, A. Johnston, and et al., “Machine learning for a sustainable energy future,” *Nat Rev Mater*, vol. 8, pp. 202–215, 2023.
- [671] M. C. Rillig, M. Ågerstrand, M. Bi, K. A. Gould, and U. Sauerland, “Risks and benefits of large language models for the environment,” *Environmental Science & Technology*, vol. 57, no. 9, pp. 3464–3466, 2023.
- [672] X. Zhi and J. Wang, “Editorial: Ai-based prediction of high-impact weather and climate extremes under global warming: A perspective from the large-scale circulations and teleconnections,” *Frontiers in Earth Science*, vol. 11, 2023.
- [673] J. Zhong, Y. Zhong, M. Han, T. Yang, and Q. Zhang, “The impact of ai on carbon emissions: evidence from 66 countries,” *Applied Economics*, vol. 0, no. 0, pp. 1–15, 2023.
- [674] M. A. Habila, M. Ouladsmane, and Z. A. Alothman, “Chapter 21 - role of artificial intelligence in environmental sustainability,” in *Visualization Techniques for Climate Change with Machine Learning and Artificial Intelligence* (A. Srivastav, A. Dubey, A. Kumar, S. Kumar Narang, and M. Ali Khan, eds.), pp. 449–469, Elsevier, 2023.
- [675] A. S. George, A. H. George, and A. G. Martin, “The environmental impact of ai: A case study of water consumption by chat gpt,” *Partners*

- Universal International Innovation Journal*, vol. 1, no. 2, pp. 97–104, 2023.
- [676] A. A. Chien, L. Lin, H. Nguyen, V. Rao, T. Sharma, and R. Wijayawardana, “Reducing the carbon impact of generative ai inference (today and in 2035),” *ACM Hot Carbon 2023*, 2023.
- [677] L. Weidinger, J. Mellor, M. Rauh, C. Griffin, J. Uesato, P.-S. Huang, M. Cheng, M. Glaese, B. Balle, A. Kasirzadeh, *et al.*, “Ethical and social risks of harm from language models,” *arXiv preprint arXiv:2112.04359*, 2021.
- [678] G. Pistilli, “What lies behind agi: ethical concerns related to llms,” *Revue Ethique et Numérique*, 2022.
- [679] L. Sun, Y. Huang, H. Wang, S. Wu, Q. Zhang, C. Gao, Y. Huang, W. Lyu, Y. Zhang, X. Li, *et al.*, “Trustllm: Trustworthiness in large language models,” *arXiv preprint arXiv:2401.05561*, 2024.
- [680] J. Hill, W. R. Ford, and I. G. Farreras, “Real conversations with artificial intelligence: A comparison between human-human online conversations and human-chatbot conversations,” *Computers in human behavior*, vol. 49, pp. 245–250, 2015.
- [681] M. Mitchell, “How do we know how smart ai systems are?,” 2023.
- [682] T. Chu, Z. Song, and C. Yang, “How to protect copyright data in optimization of large language models?,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, pp. 17871–17879, 2024.
- [683] W. J. Dupp Jr., “Artificial intelligence and academic publishing,” *Journal of Cataract & Refractive Surgery*, vol. 49, no. 7, pp. 655–656, 2023.
- [684] C. Novelli, F. Casolari, A. Rotolo, M. Taddeo, and L. Floridi, “Taking ai risks seriously: a new assessment model for the ai act,” *AI & SOCIETY*, pp. 1–5, 2023.
- [685] U. M. Fayyad, “From stochastic parrots to intelligent assistants—the secrets of data and human interventions,” *IEEE Intelligent Systems*, vol. 38, no. 3, pp. 63–67, 2023.
- [686] V. Salikutuk, D. Koert, and F. Jäkel, “Interacting with large language models: A case study on ai-aided brainstorming for guesstimation problems,” in *HHAI 2023: Augmenting Human Intellect*, pp. 153–167, IOS Press, 2023.
- [687] S. Moore, R. Tong, A. Singh, Z. Liu, X. Hu, Y. Lu, J. Liang, C. Cao, H. Khosravi, P. Denny, *et al.*, “Empowering education with llms—the next-gen interface and content generation,” in *International Conference on Artificial Intelligence in Education*, pp. 32–37, Springer, 2023.
- [688] K. Nottingham, P. Ammanabrolu, A. Suhr, Y. Choi, H. Hajishirzi, S. Singh, and R. Fox, “Do embodied agents dream of pixelated sheep?: Embodied decision making using language guided world modelling,” *arXiv preprint arXiv:2301.12050*, 2023.
- [689] K. Greshake, S. Abdelnabi, S. Mishra, C. Endres, T. Holz, and M. Fritz, “More than you’ve asked for: A comprehensive analysis of novel prompt injection threats to application-integrated large language models,” *arXiv preprint arXiv:2302.12173*, 2023.
- [690] S. Kang, J. Yoon, and S. Yoo, “Large language models are few-shot testers: Exploring llm-based general bug reproduction,” in *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*, pp. 2312–2323, IEEE, 2023.
- [691] B. Goertzel, “Artificial general intelligence: concept, state of the art, and future prospects,” *Journal of Artificial General Intelligence*, vol. 5, no. 1, p. 1, 2014.
- [692] T. Giannini and J. P. Bowen, “Generative art and computational imagination: Integrating poetry and art,” in *Proceedings of EVA London 2023*, pp. 211–219, BCS Learning & Development, 2023.
- [693] C. Summerfield, *Natural General Intelligence: How understanding the brain can help us build AI*. Oxford University Press, 2022.
- [694] N. Nascimento, P. Alencar, and D. Cowan, “Self-adaptive large language model (llm)-based multiagent systems,” *arXiv preprint arXiv:2307.06187*, 2023.
- [695] J. Pei, L. Deng, S. Song, M. Zhao, Y. Zhang, S. Wu, G. Wang, Z. Zou, Z. Wu, W. He, *et al.*, “Towards artificial general intelligence with hybrid tianjic chip architecture,” *Nature*, vol. 572, no. 7767, pp. 106–111, 2019.
- [696] T. Everitt, *Towards safe artificial general intelligence*. PhD thesis, The Australian National University (Australia), 2019.
- [697] J. Hughes, A. Abdulali, R. Hashem, and F. Iida, “Embodied artificial intelligence: Enabling the next intelligence revolution,” in *IOP Conference Series: Materials Science and Engineering*, vol. 1261, p. 012001, IOP Publishing, 2022.
- [698] Y. Ma, Z. Song, Y. Zhuang, J. Hao, and I. King, “A survey on vision-language-action models for embodied ai,” *arXiv preprint arXiv:2405.14093*, 2024.
- [699] R. Sapkota, R. Qureshi, M. F. Calero, M. Hussain, C. Badjugar, U. Nepal, A. Poulose, P. Zeno, U. B. P. Vaddevolu, H. Yan, *et al.*, “Yolov10 to its genesis: A decadal and comprehensive review of the you only look once series,” *arXiv preprint arXiv:2406.19407*, 2024.
- [700] Y. Mu, Q. Zhang, M. Hu, W. Wang, M. Ding, J. Jin, B. Wang, J. Dai, Y. Qiao, and P. Luo, “Embodiedgpt: Vision-language pre-training via embodied chain of thought,” *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [701] P. Gao, P. Wang, F. Gao, F. Wang, and R. Yuan, “Vision-language navigation with embodied intelligence: A survey,” *arXiv preprint arXiv:2402.14304*, 2024.
- [702] L. C. Magister, J. Mallinson, J. Adamek, E. Malmi, and A. Severyn, “Teaching small language models to reason,” *arXiv preprint arXiv:2212.08410*, 2022.
- [703] C. T. Wolf, “Democratizing ai? experience and accessibility in the age of artificial intelligence,” *XRDS: Crossroads, The ACM Magazine for Students*, vol. 26, no. 4, pp. 12–15, 2020.
- [704] Z. Liu, A. Qiao, W. Neiswanger, H. Wang, B. Tan, T. Tao, J. Li, Y. Wang, S. Sun, O. Pangarkar, *et al.*, “Llm360: Towards fully transparent open-source llms,” *arXiv preprint arXiv:2312.06550*, 2023.
- [705] S.-h. Huang and C.-y. Chen, “Combining lora to gpt-neo to reduce large language model hallucination,” 2024.
- [706] T. Le Scao, A. Fan, C. Akiki, E. Pavlick, S. Ilić, D. Hesslow, R. Castagné, A. S. Lucchioni, F. Yvon, M. Gallé, *et al.*, “Bloom: A 176b-parameter open-access multilingual language model,” 2023.
- [707] H. Inan, K. Upasani, J. Chi, R. Runpta, K. Iyer, Y. Mao, M. Tontchev, Q. Hu, B. Fuller, D. Testuggine, *et al.*, “Llama guard: Llm-based input-output safeguard for human-ai conversations,” *arXiv preprint arXiv:2312.06674*, 2023.
- [708] C. K. Iyer, F. Hou, H. Wang, Y. Wang, K. Oh, S. Ganguli, and V. Pandey, “Trinity: A no-code ai platform for complex spatial datasets,” in *Proceedings of the 4th ACM SIGSPATIAL International Workshop on AI for Geographic Knowledge Discovery*, pp. 33–42, 2021.
- [709] B. Allen, S. Agarwal, J. Kalpathy-Cramer, and K. Dreyer, “Democratizing ai,” *Journal of the American College of Radiology*, vol. 16, no. 7, pp. 961–963, 2019.
- [710] L. Sundberg and J. Holmström, “Democratizing artificial intelligence: How no-code ai can leverage machine learning operations,” *Business Horizons*, 2023.
- [711] D. Kreuzberger, N. Kühl, and S. Hirschl, “Machine learning operations (mlops): Overview, definition, and architecture,” *IEEE Access*, 2023.