

ChatGPT and Other Large Language Models

Jan Jovan

UP FAMNIT

89232124@student.upr.si

Abstract

Large Language Models (LLMs) have revolutionized natural language processing by leveraging vast datasets and deep learning architectures to generate coherent text, answer questions, and perform complex reasoning. This paper explores the underlying mechanisms of LLMs, including their architecture and training methods, while highlighting their growing use cases in various fields. Specifically, the applications of LLMs in medicine for aiding diagnosis and generating medical literature, in robotics for improving human-robot interaction, and in coding for automating software development are discussed. Challenges related to model evaluation, including knowledge coverage and scalability, are also addressed, along with potential future directions for enhancing LLM capabilities.

1 Introduction

The advent of large language models has revolutionized the field of natural language processing, and none have garnered more attention than ChatGPT, a cutting-edge large language model that has taken the world by storm. Developed by OpenAI [1], ChatGPT [2] is a type of transformer-based [3] language model [4] that has been trained on a massive dataset of text, allowing it to generate human-like responses to a wide range of prompts and questions. With its impressive ability to understand context, recognize nuances, and respond with coherence and fluency, ChatGPT has sparked widespread interest and debate about the potential applications and implications of large language models.

Large language models, of which ChatGPT is just one example, are a class of artificial intelligence (AI) models that are trained on vast amounts of text data to learn the patterns and structures of language. These models have made tremendous progress in recent years, achieving state-of-the-art results in a variety of natural language processing tasks, from language translation and text summarization to dialogue generation and question answering. Their capabilities have far-reaching implications for fields such as customer service, education, and content creation, among others.

However, the rapid development and deployment of large language models also raise important questions about their limitations, biases, and potential risks. As these models become increasingly integrated into various aspects of our lives, it is essential to critically examine their capabilities, limitations, and implications. This paper aims to provide an in-depth exploration of ChatGPT and other large language models, examining their underlying technologies, applications, and implications, as well as the challenges and opportunities they present for individuals, organizations, and society as a whole.

2 Under the hood

Neural Networks: The Building Blocks of Large Language Models

Large language models, such as ChatGPT, are built on top of neural networks[5], a type of machine learning model inspired by the structure and function of the human brain. Neural networks consist of layers of interconnected nodes or "neurons" that process and transform inputs into outputs. These neural networks are designed to recognize patterns in data and learn from it, making them ideal for natural language processing tasks.

Architecture: The Encoder-Decoder Paradigm

The Transformer architecture [1] consists of an encoder and a decoder. The encoder takes in a sequence of tokens (e.g., words or characters) and generates a continuous representation of the input sequence. The decoder takes this representation and generates the output sequence.

Encoder: Breaking Down Input Text

The encoder is responsible for processing the input text and generating a continuous representation of the input. The encoder consists of multiple layers, each of which applies a series of transformations to the input. The encoder's primary function is to capture the nuances of the input text and represent it in a way that the decoder can understand.

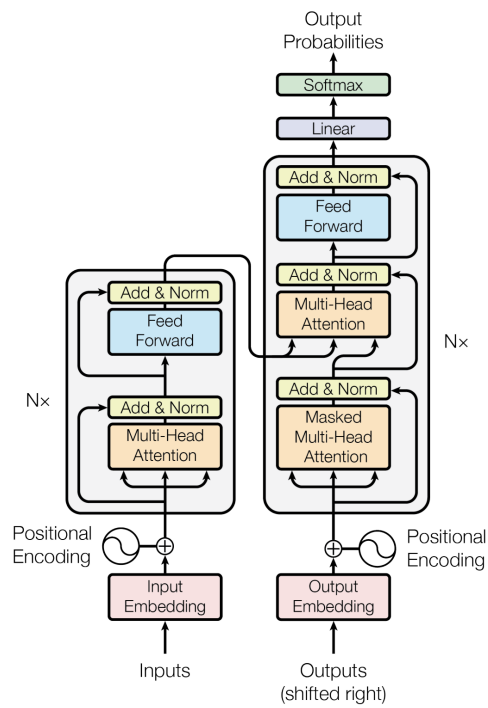


Figure 1: The Transformer - model architecture

Tokenization: Breaking Down Text into Tokens

The first step in the encoding process is tokenization. Tokenization involves breaking down the input text into individual tokens, which can be words, characters, or subwords (smaller units of words). Each token is represented as a numerical vector, known as a token embedding. This process allows the model to focus on individual components of the input text and understand their relationships.

Embedding Layer: Converting Tokens into Vectors

The token embeddings are fed into an embedding layer, which converts the token embeddings into a higher-dimensional vector space. This allows the model to capture subtle relationships between tokens and understand the context in which they are used. The embedding layer is a critical component of the encoder, as it enables the model to capture the nuances of language.

Transformer Layers: Processing Sequential Data

The output of the embedding layer is fed into a series of transformer layers. The transformer layer is a type of neural network layer that's specifically designed for sequential data like text. Each transformer layer consists of two sub-layers: a multi-head self-attention mechanism and a feed-forward neural network (FFNN).

Multi-Head Self-Attention: Understanding Context

The multi-head self-attention mechanism allows the model to attend to different parts of the input sequence simultaneously and weigh their importance. This is done by computing attention weights, which represent the relative importance of each token in the input sequence. The attention weights are computed as follows:

1. **Query, Key, and Value Matrices:** The input sequence is first split into three matrices: query (\mathbf{Q}), key (\mathbf{K}), and value (\mathbf{V}).
2. **Attention Weights:** The attention weights are computed by taking the dot product of \mathbf{Q} and \mathbf{K} , and applying a softmax function to the result.
3. **Weighted Sum:** The attention weights are used to compute a weighted sum of the value matrix, resulting in a contextualized representation of the input sequence.

Feed-Forward Neural Network (FFNN): Transforming Outputs

The output of the multi-head self-attention mechanism is fed into a FFNN, which consists of two linear layers with a ReLU activation function in between. The FFNN transforms the output of the attention mechanism into a higher-dimensional space, allowing the model to capture complex relationships between tokens.

Decoder: Generating Output Text

The decoder is responsible for generating the output text based on the continuous representation of the input generated by the encoder. The decoder consists of multiple layers, each of which applies a series of transformations to the input.

Decoder Layers: Generating Coherent Text

Each decoder layer consists of two sub-layers: a self-attention mechanism and a FFNN. The self-attention mechanism in the decoder is similar to the one in the encoder, except that it only attends to the output of the previous decoder layer. The FFNN transforms the output of the self-attention mechanism into a higher-dimensional space.

Output Linear Layer: Generating Final Output

The output of the final decoder layer is fed into an output linear layer, which generates the final output text. This layer is responsible for converting the continuous representation of the input into a coherent and context-specific output.

Training: Masked Language Modeling

The model is trained using a masked language modeling objective, where some of the input tokens are randomly replaced with a [MASK] token. The model is trained to predict the original token. This objective allows the model to learn the relationships between tokens and understand the context in which they are used.

Optimization: Minimizing Loss

The model is optimized using a variant of stochastic gradient descent, such as Adam or Adagrad. The optimization process involves minimizing the loss function, which measures the difference between the model's predictions and the true labels. The goal of optimization is to find the optimal set of model parameters that minimize the loss function.

Fine tuning

Fine-tuning large language models subsequent to their primary training phase is crucial for several reasons. Firstly, it enables the model to adapt to specific tasks and domains, allowing it to recognize and learn patterns and features unique to those areas. For instance, a model initially trained on a vast general corpus may lack the nuanced understanding required for specialized fields such as medicine or finance. Through fine-tuning, the model can become adept at identifying relevant terminology, jargon, and concepts pertinent to these domains.

Additionally, fine-tuning facilitates the acquisition of task-specific knowledge and features absent in the original training data. A model exposed to a broad corpus may not encounter the particular terminology or concepts essential for tasks like sentiment analysis or question answering. Fine-tuning addresses this gap by enabling the model to recognize and respond to specific patterns relevant to these tasks.

Moreover, fine-tuning helps mitigate overfitting by training the model on a smaller, task-specific dataset. This approach minimizes the risk of the model memorizing the training data, thereby enhancing its ability to generalize to new, unseen data. By focusing on a smaller dataset, the model is less prone to overfitting, which is a common issue with large, diverse datasets.

The performance of the model on the target task can be significantly enhanced through fine-tuning. This improvement arises from the model's ability to learn and respond to the specific patterns, features, and relationships inherent in the target data. Consequently, fine-tuning allows the model to better adapt to the specific task and domain, thus improving its predictive accuracy and generalization capabilities.

Fine-tuning also reduces the model's dependency on the initial training data, which may be biased or outdated. By leveraging a smaller, task-specific dataset, the model becomes less likely to overfit to the main training data, ensuring that its predictions remain relevant and unbiased, particularly when the original data does not accurately represent the target task or domain.

Furthermore, fine-tuning enhances the model's robustness to noise, outliers, and other forms of data corruption. This is achieved as the model learns to identify and adapt to the specific patterns and features present in the target data, thereby becoming more resilient to anomalies and data imperfections.

Fine-tuning also bolsters the model's ability to generalize to novel, unseen data. By learning to identify and respond to specific patterns and relationships in the target data, the model becomes more adept at making accurate predictions and adapting to new scenarios.

Insights into the model's decision-making process can be gleaned through fine-tuning, enabling researchers to pinpoint the most critical features and patterns within the target data. By examining the model's weights and activations during the fine-tuning process, researchers can gain a deeper understanding of the factors influencing the model's predictions.

Finally, fine-tuning typically demands fewer computational resources compared to training a model from scratch. Since the model begins from a pre-trained state, it only requires adaptation to the target task, thereby reducing the computational burden. This makes fine-tuning more feasible for smaller datasets or environments with limited computational resources.

3 ChatGPT use cases

ChatGPT, developed by OpenAI based on the GPT-4 architecture, is a sophisticated language model designed to understand and generate human-like text. This versatility allows it to perform a wide array of tasks, including answering questions on diverse topics, assisting with writing tasks, translating text between languages, and providing educational support in subjects like math, science, history, and literature. Additionally, it offers coding assistance, debugging help, and can engage in casual conversation. For researchers, ChatGPT can gather information and summarize articles, while for technical support, it guides users through troubleshooting common issues and explains software and device usage.

ChatGPT serves a broad audience, including students, professionals, writers, non-native English speakers, researchers, developers, and businesses. It offers personalized assistance to meet specific needs, enhances productivity by automating tasks, and helps users understand complex topics through simplified explanations. By providing accurate information, supporting learning, encouraging creativity, and facilitating communication, ChatGPT is a valuable tool for achieving various goals efficiently and effectively. Whether you need quick answers, detailed explanations, creative inspiration, or technical support, ChatGPT is here to help.

4 LLM for code generation and completion

The advent of large language models (LLMs) has brought transformative changes to various domains, including software development. AI-powered coding assistants, such as ChatGPT, LLaMA Code[7], and others, have emerged as vital tools for developers, enhancing productivity and facilitating the coding process. This paper explores the capabilities of these AI-powered coding assistants, detailing their applications, usage, target users, and impact on the development landscape.

4.1 ChatGPT

4.1.1 Overview

ChatGPT, developed by OpenAI, is a versatile language model that excels in understanding and generating human-like text. Although it is not exclusively designed for coding, it possesses significant capabilities that make it a valuable tool for developers.

4.1.2 Capabilities

ChatGPT can assist in generating code snippets, offering suggestions for code completion, and providing explanations for various programming concepts. It can help debug code by analyzing error messages and suggesting potential fixes. Moreover, ChatGPT can be used for learning purposes, offering examples and explanations that help novice programmers understand complex coding concepts.

4.1.3 Usage

The process of using ChatGPT typically involves integrating it into development environments or using it through conversational interfaces. Developers can interact with ChatGPT by posing questions, seeking code suggestions, or requesting explanations. Its wide applicability makes it suitable for a range of users, from beginners to experienced developers.

4.2 LLaMA Code

4.2.1 Overview

LLaMA Code is another AI-powered coding assistant designed to enhance the coding process. Built upon advanced language models, LLaMA Code focuses specifically on programming tasks, making it highly effective in generating and understanding code.

4.2.2 Capabilities

LLaMA Code offers capabilities such as real-time code completion, error detection, and suggestions for code optimization. It supports a variety of programming languages, ensuring that developers can use it across different projects. Additionally, LLaMA Code can assist in generating documentation and comments, improving code readability and maintainability.

4.2.3 Usage

Using LLaMA Code involves integrating it into integrated development environments (IDEs) or using it as a standalone tool. Developers can benefit from its context-aware suggestions, which enhance coding efficiency and accuracy. Like ChatGPT, LLaMA Code is suitable for a broad audience, including professional developers, educators, and students.

4.3 Other Coding LLMs

4.3.1 Overview

Beyond ChatGPT and LLaMA Code, several other AI-powered coding assistants have been developed, each with unique capabilities. Examples include Codex [8] (also by OpenAI), Starcoder[10] and Tabnine[9].

4.4 Examples of using ChatGPT for coding

Here are few simple user prompts and ChatGPT 4o answers:

Note: model was used on 17.5.2024 and since then, there may be newer and better version of this same model.

Write me a python script for a simple calculator with all basic mathematical operations.

Below is a Python script for a simple calculator that can perform basic mathematical operations (addition, subtraction, multiplication, and division). This script uses functions to handle the different operations and a loop to continuously prompt the user for input until they choose to exit.

```
def add(x, y):
    return x + y

def subtract(x, y):
    return x - y

def multiply(x, y):
    return x * y

def divide(x, y):
    if y == 0:
        return "Error! - Division -by- zero."
    return x / y

def calculator():
    print("Welcome to the simple calculator!")
    print("Available operations:")
    print("1. - Addition")
    print("2. - Subtraction")
    print("3. - Multiplication")
    print("4. - Division")
    print("5. - Exit")

    while True:
        choice = input("Please choose an operation (1/2/3/4/5):-")

        if choice == '5':
            print("Exiting the calculator. - Goodbye!")
            break

        if choice in ['1', '2', '3', '4']:
            try:
                num1 = float(input("Enter the first number:-"))
                num2 = float(input("Enter the second number:-"))
            except ValueError:
                print("Invalid input! - Please enter numeric values.")
```



```

        continue

    if choice == '1':
        print(f"The result is: {add(num1, num2)}")
    elif choice == '2':
        print(f"The result is: {subtract(num1, num2)}")
    elif choice == '3':
        print(f"The result is: {multiply(num1, num2)}")
    elif choice == '4':
        result = divide(num1, num2)
        print(f"The result is: {result}")
    else:
        print("Invalid choice! Please select a valid operation.")

if __name__ == "__main__":
    calculator()

```

How to Use the Script

1. Copy and paste the script into a file named `calculator.py`.
2. Run the script using Python. For example, in your terminal or command prompt, navigate to the directory containing `calculator.py` and run:

```
python calculator.py
```

3. Follow the on-screen prompts to perform calculations. The script will continue to run until you choose to exit by selecting option 5.

4.5 Answering a problem from StackOverflow

Question was:

i have a program in python, which wont compile in vscode

using python3 in terminal works, but in vscode it does not. here is the output it produced:

```
[Running] python -u "/Users/path/hello.py"
/bin/sh: python: command not found
```

```
[Done] exited with code=127 in 0.039 seconds
```

I have tried uninstalling and reinstalling both python and vscode, but the problem still continues. (im on mac)

ChatGPT answer was this:

It seems that the issue is due to Visual Studio Code (VS Code) not being able to find the Python interpreter. This can happen if the Python path is not correctly set in VS Code. Here's how you can resolve this issue:

1. Check Python Installation

First, ensure that Python is installed correctly and accessible from the terminal. Open your terminal and run:

```
which python3
```

This should return the path to the Python 3 interpreter, something like `/usr/local/bin/python3`.

2. Update VS Code Settings

- (a) **Open Command Palette:** Press `Cmd+Shift+P` to open the Command Palette in VS Code.
- (b) **Select Interpreter:** Type `Python: Select Interpreter` and select it.
- (c) **Choose the Correct Interpreter:** Choose the Python interpreter that corresponds to `python3`. It should be similar to the path you got from `which python3`.

3. Update the `settings.json`

If the above steps do not resolve the issue, you may need to manually set the Python path in the `settings.json` file.

- (a) **Open Command Palette:** Press `Cmd+Shift+P`.
- (b) **Open Settings (JSON):** Type `Preferences: Open Settings (JSON)` and select it.
- (c) **Set Python Path:** Add the following configuration, replacing `/usr/local/bin/python3` with the path you got from the `which python3` command:

```
{
  "python.pythonPath": "/usr/local/bin/python3"
}
```

4. Update PATH Variable (If Needed)

If the issue persists, you might need to ensure that the `python` command points to Python 3.

- (a) **Open your shell configuration file:** This could be `~/.bash_profile`, `~/.zshrc`, or `~/.bashrc`, depending on your shell.
- (b) **Add Python Path:** Add the following line, replacing `/usr/local/bin/python3` with the path you found earlier:

```
export PATH="/usr/local/bin:$PATH"
```

- (c) **Source the Configuration File:**

```
source ~/.bash_profile # or ~/.zshrc , or ~/.bashrc
```

5. Restart VS Code

After making these changes, restart VS Code to ensure it picks up the updated settings and environment variables.

By following these steps, you should be able to configure VS Code to use the correct Python interpreter and resolve the issue of it not being found.

4.6 Trouble with code generation/bug fixing

Using ChatGPT for code generation and bug fixing presents several challenges and limitations. The AI can produce code with syntax and logical errors, outdated practices, and potential security vulnerabilities due to its reliance on predicting the next probable token rather than understanding the code deeply. Additionally, the model's ability to understand the specific context and requirements of a task is limited, often leading to misinterpretations and incomplete solutions. The quality of the generated code is directly influenced by the training data, which includes both good and bad examples from the internet, potentially resulting in non-standard or suboptimal code.

Furthermore, over-reliance on ChatGPT can result in decreased developer skills and critical thinking abilities, as well as ethical and privacy concerns related to data sharing. While the AI can assist with identifying bugs, it often provides superficial fixes without addressing root causes. Integrating the generated code into existing systems and ensuring thorough testing also remain significant challenges. Despite these issues, ChatGPT can be a valuable tool when used cautiously, with developers critically reviewing and testing the AI's output to maintain code quality and robustness.

5 Transforming Medical Writing and Education

Another area where LLMs shine is in medical writing. If you've ever had to write a research paper, draft a medical report, or even put together patient information brochures, you know how time-consuming it can be. LLMs can help by quickly pulling together information from a massive database of medical literature, summarizing findings, and even suggesting better ways to present data.

In education, LLMs can be used to create realistic simulations for training medical students or generate a wide range of questions for exams. They can also explain complex medical concepts [11] in simpler terms, which is a huge help for both students and patients. But again, caution is necessary. LLMs need oversight to ensure they don't accidentally spread misinformation or make errors that could lead to misunderstandings.

5.1 Boosting Medical Research and Project Management

When it comes to research, LLMs can take on some of the heavy lifting. They can help researchers by automating literature searches, organizing data, and even managing project timelines and budgets. This frees up more time for researchers to focus on the creative and analytical aspects of their work.

Project management in the medical field involves a lot of coordination—between researchers, clinicians, patients, and often regulatory bodies. LLMs can facilitate communication across these groups by providing clear summaries, helping to bridge any gaps in understanding. They also enhance collaboration by managing information more effectively.

However, using LLMs in this way isn't without its challenges. The accuracy of the data and the ability to understand cultural nuances are just a couple of the issues that need to be addressed for these models to be fully effective.

5.2 Challenges and Future Prospects

While LLMs offer numerous benefits, their integration into healthcare is not without challenges. Ethical considerations are paramount—issues like data privacy, the transparency of AI decision-making processes, and the risk of over-reliance on these tools need to be carefully managed. The future of LLMs in medicine looks promising, but it requires a balanced approach.

One of the key areas for future development is in creating multimodal LLMs—models that can understand and integrate not just text, but also images, sounds, and other types of data. This could lead to even more powerful tools for diagnosing and treating diseases. Additionally, there’s a need for more research into how these models make decisions, which will help improve their reliability and safety in clinical settings.

6 The Med-PaLM Model: Applications in Medicine

6.1 Overview of Med-PaLM

Med-PaLM [12] is a specialized large language model (LLM) developed to address challenges in medical question answering, building on general-purpose LLM architectures like PaLM (Pathways Language Model) [13] with targeted fine-tuning for the medical domain. This model was a significant advancement, being the first to surpass the “passing” score in the US Medical Licensing Examination (USMLE)-style questions, thereby demonstrating its capability in handling medical knowledge with competence comparable to professionals. Its successor, Med-PaLM 2, continues to refine and improve upon these capabilities, offering even greater accuracy and alignment with medical standards.

6.2 Applications for Doctors and Medical Students

The Med-PaLM model holds significant promise for both practicing doctors and medical students by providing a powerful tool to aid in the retrieval and reasoning over complex medical information. Key applications include:

1. Medical Education and Training: For medical students, Med-PaLM can serve as a dynamic study aid. The model can answer a wide range of clinical questions and provide explanations that reinforce learning. By simulating real-world question-answer scenarios, Med-PaLM helps students prepare for exams like the USMLE, offering not just factual answers but reasoning paths that can enhance understanding. Additionally, its ability to respond to nuanced medical queries supports deeper exploration into topics that require more than a surface-level answer, such as differential diagnoses or treatment protocols.

2. Decision Support for Physicians: Med-PaLM can act as a supplemental decision support tool for practicing doctors. The model excels in synthesizing vast amounts of medical literature and clinical guidelines, helping physicians quickly access relevant information when making decisions about patient care. This is particularly valuable in time-sensitive situations or when a clinician needs to stay updated on the latest medical research and best practices. For example, Med-PaLM can assist in answering complex queries related to rare diseases, emerging treatments, or uncommon drug interactions.

6.3 Pros of Using Med-PaLM

The Med-PaLM model offers several advantages for medical professionals and students:

1. Comprehensive Knowledge Base: Med-PaLM is trained on a wide array of medical datasets, including those covering medical licensing exams, medical research, and clinical guidelines. This comprehensive training enables the model to provide accurate and detailed answers across diverse medical domains.

2. Improved Medical Reasoning: Med-PaLM’s use of techniques like Chain-of-Thought prompting and ensemble refinement enhances its reasoning capabilities, particularly in handling multi-step clinical problems. These advanced reasoning strategies allow it to simulate the thought processes of clinicians, making it a more reliable source for decision support and education.

3. High-Quality Long-Form Responses: Unlike simpler question-answering systems that focus on brief responses, Med-PaLM can generate well-articulated long-form answers. This feature is crucial for medical professionals who require detailed explanations, clinical rationales, or comprehensive overviews of a medical issue.

4. Ongoing Learning and Updates: Med-PaLM can be continuously fine-tuned and updated with new medical information. This adaptability ensures that the model remains relevant and up-to-date with the latest medical guidelines and research, reducing the likelihood of outdated advice.

6.4 Cons and Limitations of Med-PaLM

While Med-PaLM represents a significant leap forward in medical AI, there are limitations to its use:

1. Risk of Over-Reliance: There is a risk that medical students or less experienced physicians might over-rely on Med-PaLM’s outputs without sufficient critical evaluation. The model, despite its accuracy, is still an AI system and may not capture the full complexity of certain medical scenarios, particularly those that require nuanced human judgment and experience.

2. Potential for Bias: As with all LLMs, Med-PaLM’s responses can reflect biases present in the training data. These biases may manifest in ways that affect clinical recommendations, particularly in areas related to health equity and demographic differences in patient populations. It is essential that users of Med-PaLM remain aware of this and apply critical thinking when interpreting its outputs.

3. Safety Concerns in Clinical Application: Although Med-PaLM has demonstrated impressive performance, its deployment in real-world clinical settings must be approached cautiously. There are still open questions regarding the model’s ability to consistently provide safe and reliable advice, especially in high-stakes situations. Misinterpretation of AI-generated responses could lead to suboptimal or even harmful patient outcomes.

4. Ethical and Regulatory Considerations: The integration of AI tools like Med-PaLM into clinical practice raises ethical and regulatory concerns. Issues such as accountability for AI-generated advice, patient consent, and data privacy must be carefully navigated. Ensuring that Med-PaLM is used in compliance with healthcare regulations is critical to its safe and effective application.

7 Integration of Large Language Models in Robotics

Large Language Models (LLMs) have demonstrated remarkable capabilities across a variety of domains, and their integration into robotics opens up new possibilities for enhancing robot task planning and execution. These models, especially when combined with multimodal approaches, enable robots to interpret complex instructions, reason about environments, and execute intricate tasks.

One of the key contributions of LLMs in robotics is their ability to translate natural language instructions into actionable plans. This process involves breaking down high-level linguistic commands into precise action sequences that a robot can execute. For example, using LLMs like GPT-4V, robots can interpret an instruction such as “pick up the green chip bag from the bottom drawer and place it on the counter,” and generate a detailed plan that includes identifying the target object, performing the necessary movements, and updating the environment accordingly.

To address the challenges that arise when robots interact with complex environments, it is crucial to integrate visual perception with language processing. LLMs, when enhanced with multimodal capabilities, such as vision-language models, can significantly improve robot performance in tasks that require understanding both the visual and linguistic context. This multimodal integration allows robots to better comprehend their surroundings, reason about object locations, and generate more accurate action plans.

However, LLMs face limitations in robotics, particularly due to the lack of experiential exposure to physical environments and the challenges in obtaining large, diverse datasets for robot interactions. While LLMs can generalize from text data, real-world scenarios often require more specialized training and fine-tuning to achieve high levels of performance. Additionally, the complexity of robotic tasks necessitates continuous advancements in multimodal integration and the development of more robust frameworks that can handle the intricacies of human-robot-environment interaction.

Despite these challenges, the future of LLMs in robotics holds significant promise. The potential applications range from precision agriculture, where robots equipped with LLMs can perform labor-intensive tasks, to healthcare, where robots can assist in complex procedures requiring safety and precision. As LLM-centric AI systems continue to evolve, their ability to bridge the gap between language, perception, and action will drive further innovation in the field of robotics.

7.1 Multimodal Task Planning

In multimodal task planning, LLMs play a crucial role by integrating information from various data streams—textual, visual, and auditory—allowing robots to generate comprehensive plans for complex tasks. For instance, frameworks like Inner Monologue and SayCan demonstrate the use of LLMs as a central component in planning and executing tasks, leveraging multimodal feedback to enhance the robot’s understanding of its environment. This approach not only improves task performance but also enables more dynamic and autonomous interaction with the physical world.

7.2 Challenges and Future Directions

While LLMs offer powerful capabilities in robotics, several challenges remain. The current reliance on predefined actions limits the robot’s flexibility in real-time environments. Additionally, the closed-source nature of many advanced models, such as GPT-4V, poses barriers to their widespread application in embedded systems. Future research should focus on overcoming these limitations by developing more adaptive and open frameworks that allow for greater customization and real-time decision-making in robotic systems.

8 LLM Evaluation

8.1 Introduction

With the growing prominence of Large Language Models (LLMs) in both research and application, evaluating their capabilities has become increasingly complex. LLM evaluation serves multiple

purposes: ensuring model training stability, benchmarking model performance, and understanding where we stand as a field in terms of model capabilities. Current evaluation methodologies can be broadly categorized into three main approaches: automated benchmarking, human evaluation, and model-as-judge evaluation. Each method presents unique advantages and challenges.

8.2 Automated Benchmarking

Automated benchmarking is one of the most common methods for LLM evaluation. It involves constructing evaluations from a collection of samples, typically paired with reference outputs (often called "gold" outputs). These benchmarks are designed to assess model performance on specific tasks, such as email classification or more abstract tasks like mathematical problem-solving.

Benchmarks can be effective when the task is well-defined, allowing clear metrics to assess model performance. However, challenges arise when evaluating more complex capabilities, such as reasoning or creative tasks like poetry writing. For example, evaluating whether a model is "good at math" may involve multiple underlying skills, such as arithmetic and logic, making it difficult to design a comprehensive benchmark.

Another critical issue with automated benchmarking is *contamination*. Contamination occurs when evaluation datasets end up in the model's training data, leading to artificially inflated performance scores. Various approaches have been proposed to mitigate contamination, including the use of canary strings and dynamic benchmarks, but these solutions can be costly and imperfect.

8.3 Human Evaluation

Human evaluation provides a more flexible and open-ended approach to LLM evaluation. By using humans as judges, researchers can assess more complex tasks that may not be easily captured by automated benchmarks. Human evaluators are also less susceptible to contamination since they typically create new prompts for each evaluation.

Different methods of human evaluation exist, ranging from informal "vibes-checks" to systematic annotations. Vibes-checks involve community members manually testing models on undisclosed prompts to get a general sense of their performance. While these checks offer anecdotal insights, they are prone to confirmation bias and lack reproducibility.

More structured approaches, such as systematic annotations, involve providing specific guidelines to paid annotators, aiming to reduce subjectivity in evaluations. However, human bias is still a concern, especially when evaluating tasks requiring factuality or logical reasoning. Research has shown that human evaluators can be influenced by factors like the tone of the model's response, leading to skewed ratings.

8.4 Models as Judges

In response to the high cost of human evaluation, researchers have explored the use of models as judges. This approach leverages high-capability LLMs or specialized models to evaluate the outputs of other models. While using models as evaluators can reduce costs, it introduces new challenges. LLMs may favor their own outputs when scoring, and their scoring can be inconsistent with human preferences.

One concern with using models as judges is the potential for subtle, unobservable biases. Much like over-reliance on self-reinforcing processes in other fields (e.g., genetic crossbreeding), using LLMs to evaluate other LLMs can introduce minute biases that may accumulate and manifest in unexpected ways over time.

8.5 The Purpose of LLM Evaluation

LLM evaluation serves three main purposes:

1) Non-regression Testing: Ensures that changes to the training process, such as adjustments in architecture or data, do not degrade model performance. This type of evaluation focuses on the stability of the training process rather than absolute performance metrics.

2) Leaderboards and Rankings: Provides a means of ranking models based on their performance across various benchmarks. Leaderboards are useful for comparing models and selecting the best architectures. However, score instability and the subjective nature of human evaluations highlight the need for stable ranking systems.

3) Evaluating Model Capabilities: Seeks to assess the broader capabilities of models beyond specific tasks. This area of evaluation is still in its infancy, as defining what constitutes a "capability" remains a challenge. Drawing from fields such as social sciences may provide new frameworks for evaluating complex capabilities, though interdisciplinary collaboration will be crucial.

9 Evaluation Experiment: Challenges and Limitations

In an effort to evaluate the capabilities of Large Language Models (LLMs), I designed and conducted a series of evaluation tests. This process involved carefully crafting a set of prompts, along with desired outputs, to assess the model's performance across various tasks. While this approach provided valuable insights, it also presented significant challenges, particularly in terms of coverage and scalability.

The first challenge was ensuring comprehensive coverage across diverse fields of knowledge. LLMs are designed to handle a broad range of topics, from technical and scientific domains to more creative and social fields. However, it is difficult to encapsulate the full scope of human knowledge within a limited set of evaluation prompts. Despite my best efforts to include a wide array of topics—ranging from mathematics, programming, and general knowledge to literature, history, and ethics—there were inevitably gaps in coverage. Certain specialized areas of expertise or emerging topics were either underrepresented or omitted entirely from the test set.

Another significant challenge was the rapid pace of development within the LLM landscape. New models and architectures are introduced at an astonishing rate, often with incremental improvements in specific capabilities. Keeping up with this influx of models proved to be a daunting task, particularly when trying to maintain a consistent evaluation framework. Each new model release often necessitated revisiting and expanding the test set to accommodate newly claimed capabilities or performance improvements.

Furthermore, evaluating models on tasks that require more complex reasoning, creativity, or ethical considerations demanded specialized prompts and desired outputs that could be highly subjective. These tasks are difficult to standardize and evaluate quantitatively, making the results more susceptible to interpretation and bias.

In addition to these coverage and scalability challenges, there was also the issue of time and resources. Conducting rigorous evaluations across a broad range of models requires significant computational resources, particularly for large-scale models. Moreover, as new models emerge weekly, maintaining an up-to-date leaderboard or benchmark becomes a continuous and resource-intensive effort.

Despite these limitations, the evaluation experiment provided critical data points that highlight the strengths and weaknesses of different models within the tested domains. The insights gained through this evaluation underscore the need for more scalable and adaptable evaluation methodologies—ones that can dynamically adjust to the expanding knowledge and capabilities embedded

within LLMs. In future work, I propose exploring semi-automated evaluation techniques, leveraging both human oversight and automated systems to expand coverage and keep pace with the rapid evolution of LLMs.

10 References

- [1] OpenAI, inc. (17.5.2024) <https://openai.com/>
- [2] Achiam, OpenAI Josh et al. “GPT-4 Technical Report.” (2023) [arXiv:2303.08774](https://arxiv.org/abs/2303.08774)
- [3] Vaswani, Ashish et al. “Attention is All you Need.” *Neural Information Processing Systems* (2017) [arXiv:1706.03762](https://arxiv.org/abs/1706.03762)
- [4] R. Rosenfeld, ”Two decades of statistical language modeling: where do we go from here?,” in *Proceedings of the IEEE*, vol. 88, no. 8, pp. 1270-1278, Aug. 2000, doi: 10.1109/5.880083. keywords: Natural languages;Speech recognition;Bayesian methods;Probability distribution;Information retrieval;Training data;Associate members;Paper technology;Routing;Optical character recognition software
- [5] What is a neural network, IBM (17.5.2024) <https://www.ibm.com/topics/neural-networks>
- [6] Github Copilot (17.5.2024), <https://github.com/features/copilot>
- [7] Rozière, Baptiste et al. “Code Llama: Open Foundation Models for Code.” *ArXiv abs/2308.12950* (2023): n. pag.
- [8] Codex, OpenAI Inc. (17.5.2024), <https://openai.com/index/openai-codex/>
- [9] Tabnine (17.5.2024), <https://www.tabnine.com/>
- [10] StarCoder: A State-of-the-Art LLM for Code, HuggingFace blog post (17.5.2024), <https://huggingface.co/blog/starcoder>
- [11] Xiangbin Meng, Xiangyu Yan, Kuo Zhang, Da Liu, Xiaojuan Cui, Yaodong Yang, Muhan Zhang, Chunxia Cao, Jingjia Wang, Xuliang Wang, Jun Gao, Yuan-Geng-Shuo Wang, Jia-ming Ji, Zifeng Qiu, Muzi Li, Cheng Qian, Tianze Guo, Shuangquan Ma, Zeying Wang, Zexuan Guo, Youlan Lei, Chunli Shao, Wenyao Wang, Haojun Fan, Yi-Da Tang, The application of large language models in medicine: A scoping review, *iScience*, Volume 27, Issue 5, 2024, 109713, ISSN 2589-0042, <https://doi.org/10.1016/j.isci.2024.109713>. (<https://www.sciencedirect.com/science/article/pii/S2589004224009350>)
- [12] Singhal, K., Azizi, S., Tu, T. et al. Large language models encode clinical knowledge. *Nature* 620, 172–180 (2023). <https://doi.org/10.1038/s41586-023-06291-2>
- [13] Chowdhery, Aakanksha et al. “PaLM: Scaling Language Modeling with Pathways.” *J. Mach. Learn. Res.* 24 (2022): 240:1-240:113.
- [14] Wang, Jiaqi et al. “Large Language Models for Robotics: Opportunities, Challenges, and Perspectives.” *ArXiv abs/2401.04334* (2024): n. pag.