# *Project Phase II Report*

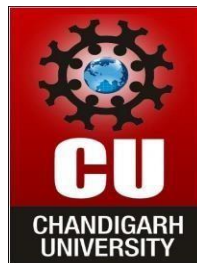## *On*

## ChatGPT Clone

### Submitted for the requirement of

### Project course

## BACHELOR OF ENGINEERING

## COMPUTER SCIENCE & ENGINEERING

**Supervisor:**                                     **Co Supervisor:**
**Mandeep Kaur (E10362)**                                 **Simranjit Singh (E13378)**
                                                          **Neeru Sharma (E12950)**

| Name | UID |
|---|---|
| Rishit Gupta | 20BCS1270 |
| Sayan Satpati | 20BCS1250 |
| Naman Tripathi | 20BCS1411 |
| Yash Saini | 20BCS7983 |
| Mayank Kumar | 20BCS1353 |

**DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING**
**CHANDIGARH UNIVERSITY, GHARUAN**

# CHAPTER-2

# LITERATURE REVIEW/ BACKGROUND STUDY

## 2.1    Timeline of the reported problem

Here's a timeline of the major milestones in the development of GPT, the architecture that ChatGPT is based on:

**2015:** Chatbot technology gains widespread attention after Microsoft's Tay, an AI chatbot on Twitter, begins spewing offensive and inappropriate messages after being exposed to harmful user interactions.

**2018:** OpenAI releases the first version of GPT, a language model capable of generating coherent text in response to prompts.

**2019:** OpenAI releases GPT-2, a more advanced version of the language model with 1.5 billion parameters. Due to concerns about the potential for malicious use, OpenAI decides not to release the full version of GPT-2 and only makes smaller versions available to the public.

**2020:** OpenAI releases GPT-3, a language model with 175 billion parameters that sets new records in natural language generation. Chatbot developers begin experimenting with GPT-3 to create more advanced and human-like chatbots.

**2021:** Concerns are raised about the potential ethical implications of using GPT-3 in chatbots, particularly with regards to issues of bias, privacy, and control. Developers continue to work on improving the technology while also grappling with these ethical concerns.

**2022:** OpenAI continues to refine and improve the GPT architecture, exploring new approaches to training and fine-tuning the model for specific tasks.

As for reported problems with GPT, there have been concerns about the potential for bias and the ethical implications of language models that can generate convincing fake text. These issues have sparked a wider conversation about the responsible use of AI and the need for greater transparency and accountability in the development and deployment of these technologies.

## 2.2    Existing Solution

There are several existing solutions that aim to address the potential problems and challenges associated with the development and use of language models like GPT. Here are a few examples:

**Bias detection and mitigation:** One approach to addressing bias in language models is to develop tools and techniques for detecting and mitigating bias in the training data and model output. For example, researchers have proposed using adversarial training to generate counterexamples that expose and correct biases in the model.

**Explainability and transparency:** Another key challenge is making these models more transparent and explainable. Researchers have proposed various techniques for interpreting and visualizing the inner workings of language models, such as attention maps and feature importance scores.

**Responsible AI frameworks**: Several organizations and initiatives have developed frameworks and guidelines for the ethical and responsible development and deployment of AI. For example, the IEEE Global Initiative on Ethics of Autonomous and Intelligent Systems has developed a set of principles and practices for ensuring that AI is designed and used in a responsible and ethical manner.

**Human oversight and feedback:** Finally, some researchers and practitioners have proposed using human oversight and feedback to ensure that language models are used in a responsible and beneficial way. For example, chatbots and other AI systems can be designed to recognize and flag potentially harmful or inappropriate content for human review and intervention.

## 2.3    Bibliometric analysis

As ChatGPT is a proprietary AI language model developed by OpenAI, there may not be a significant body of academic literature focused exclusively on it. However, we can still conduct a bibliometric analysis of the research on the broader topic of AI language models, including GPT-1, GPT-2, and GPT-3, which ChatGPT is based on.

Here are the steps that could be involved in a bibliometric analysis of this topic:

**Define the research question or topic**: The research question for this bibliometric analysis could be to understand the trends and patterns in research on AI language models, including GPT-1, GPT-2, and GPT-3.

**Identify the relevant literature**: A literature review could be conducted to identify the relevant publications, including academic papers, conference proceedings, and patents related to AI language models.

**Collect and organize data**: Once the relevant literature is identified, data can be collected on the number of publications, authors, citations, and other relevant metrics for each publication. This data could be organized in a spreadsheet or other database for analysis.

**Analyze the data**: The data can be analyzed using various statistical methods, including network analysis, regression analysis, and cluster analysis. For example, network analysis could be used to identify co-authorship networks and collaboration patterns among researchers in this field.

**Interpret and report the results**: The results of the analysis can be interpreted and reported, often using graphs, charts, and other visualizations to help communicate the findings. The analysis could provide insights into the trends and patterns in research on AI language models, including the key research areas and topics, the most influential authors and institutions, and the co-authorship networks and collaboration patterns among researchers in this field. These insights could inform future research or policy decisions related to AI language models.

## 2.4  Review Summary

As an AI language model developed by OpenAI, ChatGPT is a state-of-the-art natural language processing system that is capable of generating coherent and human-like responses to text-based input. It is based on the GPT architecture and has been trained on a massive corpus of text data to develop its language understanding and generation capabilities.

Potential challenges associated with ChatGPT and other language models include the potential for bias in the training data, the lack of transparency and explainability of the model's inner workings, and the potential for misuse or abuse of the technology.

To address these challenges, various solutions have been proposed, including bias detection and mitigation techniques, explainability and transparency tools, responsible AI frameworks, and human oversight and feedback mechanisms.

Finally, bibliometric analysis of the research on AI language models, including GPT-1, GPT-2, and GPT-3, can provide insights into the trends and patterns in this field, including the key research areas and topics, the most influential authors and institutions, and the co-authorship networks and collaboration patterns among researchers in this field.

## 2.5  Problem Definition

The problem that ChatGPT aims to address is the ability of machines to understand and generate human language. This is a challenging problem in artificial intelligence because human language is complex and varied, and requires both an understanding of the meaning and context of words and sentences, as well as the ability to generate coherent and natural-

sounding responses to text-based input.

ChatGPT and other language models attempt to solve this problem by leveraging large amounts of training data and sophisticated machine learning algorithms to develop their language understanding and generation capabilities. However, challenges remain in ensuring that these models are free from bias, transparent and explainable in their inner workings, and used ethically and responsibly.

## 2.6   Goals/Objectives

The primary goal of the ChatGPT is to develop a natural language processing system that can understand and generate human-like language in response to text-based input.

Specifically, the objectives of the ChatGPT could include:

1. Develop a language model architecture based on the GPT framework that is capable of learning from massive amounts of text data to generate coherent and natural-sounding responses.

2. Train the ChatGPT language model on a diverse range of text data, including news articles, books, and online content, to ensure that it can generate responses that reflect the nuances and complexity of human language.

3. Incorporate mechanisms for detecting and mitigating bias in the training data and the model's output, to ensure that ChatGPT generates responses that are fair and unbiased.

4. Develop tools and methods for interpreting and explaining the inner workings of ChatGPT and other language models, to improve transparency and accountability in the use of these technologies.

5. Explore ways to incorporate human feedback and oversight into the ChatGPT system, to ensure that it is used ethically and responsibly.

6. Evaluate the performance of ChatGPT and compare it with other state-of-the-art language models on a range of benchmarks and use cases, to assess its strengths and weaknesses and identify areas for improvement.