

Bachelor of Science in AI & DS

Honours Project Report

MirrorAI – LLM Based Mental Health Assistance

2025

Rikzy Jezuli

IIT No: 20210667

RGU No: 2122107

Aberdeen's Robert Gordon University (RGU) UK

DECLARATION

I confirm that the work contained in this BSc project report has been composed solely by myself and has not been accepted in any previous application for a degree. All sources of information have been specifically acknowledged, and all verbatim extracts are distinguished by quotation marks.

Date: 15/04/2025

Name: Rikzy Jezuli

ABSTRACT

This project explores the development of MirrorAI, an AI companion designed to support mental health using advanced Large Language Models (LLMs). Mental health problems are increasing worldwide, and many people do not have access to timely or affordable support. While AI tools like chatbots have been introduced, most of them lack human-like interaction, empathy, and personalization.

The aim of this project is to build an AI-based mental health assistant that offers more realistic and supportive conversations. The system uses fine-tuned LLMs trained on mental health datasets and includes features such as emotion detection, guided self-care sessions, and mood tracking. The project uses both qualitative and quantitative methods to evaluate the performance of the chatbot with comparison with existing tools and comparing the responses of the model to that of a human therapist.

The final outcome will be a working prototype of MirrorAI, tested by real users, and improved based on their feedback. This project contributes to the growing need for accessible, empathetic, and intelligent mental health support systems powered by AI.

ACKNOWLEDGEMENT

First and foremost, I would like to express my sincere gratitude to my supervisor, Ms. Dileeka Alwis, for her continuous guidance, support, and encouragement throughout this project. Her expertise and constructive feedback were invaluable in helping me shape and complete my research successfully.

I also extend my heartfelt thanks to the lecturing staff at Informatics Institute of Technology, Colombo, for their dedication and support throughout my academic journey. Their knowledge and teaching played an important role in developing the skills and understanding required for this project.

Table of Contents

DECLARATION	iii
ABSTRACT	iv
ACKNOWLEDGEMENT	v
LIST OF FIGURES	xv
LIST OF TABLES	xvi
LIST OF ACRONYMS	xvii
01 INTRODUCTION	1
1.1 Chapter Overview	1
1.2 Problem Domain	1
1.3 Problem Statement	2
1.4 Research Motivation	2
1.5 Research Gap	2
1.6 Research Questions	3
1.7 Research Aim and Objectives	4
1.8 Significance of the Research	4
1.8 Scope of the Research	5
02 LITERATURE REVIEW	5
2.1 Chapter Overview	5
2.2 Background to the problem	5

2.3 Related Work	6
2.3.1 Early Approaches in AI-Based Mental Health Tools	6
2.3.2 Emergence of Large Language Models (LLMs).....	6
2.4 Therapeutic Approaches in AI-Based Mental Health Solutions	7
2.4.1 Cognitive Behavioural Therapy (CBT)-Based Chatbots	7
2.4.2 Other Therapeutic Approaches	8
2.3.3 Multimodal and Hybrid Approaches.....	8
2.5 Advancements in Large Language Models for Mental Health Applications.....	9
2.5.1 LLM-Based Chatbots for Psychological Interventions.....	9
2.5.2 Predictive Models and Diagnostic Capabilities of LLMs	9
2.6 Datasets and Training Strategies for Mental Health AI Tools.....	10
2.6.1 Datasets for Training Mental Health Chatbots.....	10
2.6.2 Challenges in Dataset Diversity and Representation	10
2.7 Evaluation of AI-Based Mental Health Tools.....	11
2.7.1 Quantitative Evaluation Metrics	11
2.7.2 User and Professional Feedback	11
2.8 Comparison of Related Work	12
2.9 Comparison of AI based Mental Health Tools and Studies based on Technical Approach and Therapeutic Orientation	14
2.10 Literature Map	16
03 METHODOLOGY	16

3.1 Chapter Overview	16
3.2 Research Design.....	17
3.3 Data Collection Methods	18
3.4 Data Sources	18
3.4 Model Development.....	18
3.6 Model Selection	19
3.7 Experimental Setup.....	19
3.8 Evaluation Metrics	20
3.9 Validation Strategy	20
3.10 Ethical Considerations	21
3.11 Limitations	21
05. DATA PREPARATION AND ANALYSIS.....	21
5.1 Chapter Overview	21
5.2 Steps in Data Preprocessing.....	22
5.2.1 Dataset Inspection.....	22
5.2.2 Standardization of Metadata	22
5.2.3 Handling Missing Values.....	23
5.2.4 Cognitive Distortion Labeling	23
5.2.5 Conversation Structure Preservation.....	24
5.3 Merging the Datasets	24

5.4 Dataset Validation.....	25
5.5 Conclusion	26
06. DATA ANALYSIS.....	27
6.1 Chapter Overview	27
6.2 Single-Turn Conversations Analysis	27
6.2.1 Word Count Distribution	27
6.2.2 Sentiment Analysis	28
6.2.3 Cognitive Distortions Distribution.....	29
6.2.4 Keyword Analysis.....	30
6.4 Multi-Turn Conversations Analysis.....	31
6.4.1 Turn-Taking Dynamics	31
6.4.2 Sentiment Progression.....	31
6.4.3 Topic Progression	32
6.4.4 Cognitive Distortions Across Turns.....	33
6.5 Comparative Analysis: Single-Turn vs. Multi-Turn Conversations	34
6.5.1 Word Count Comparison	34
6.6 Key Insights	35
07. FEATURE ENGINEERING.....	36
7.1 Introduction to Feature Engineering	36
7.2 Structural Feature Engineering	36

7.2.1 Conversation Type Classification	36
7.2.2 Turn-based Features.....	37
7.3 Content-based Feature Engineering	37
7.3.1 Cognitive Distortion Detection	37
7.3.2 Sentiment Analysis Features	39
7.3.3 Conversation Source Metadata	39
7.4 Feature Validation and Quality Assessment	40
7.5 Feature Integration and Storage	40
7.6 Summary	41
8.0 MODEL DEVELOPMENT	41
8.1 Chapter Overview	41
8.2 Model Selection and Justification	41
8.2.1 Evaluation of Model Options	42
8.3 Justification for Llama 3.2	42
8.3.1 Suitability for Mental Health Support.....	43
8.4 Model Architecture and Finetuning Process	44
8.4.1 Architecture of Llama 3.2 3B	44
8.4.2 Fine-tuning Process.....	45
8.4.3 Dataset Preparation:	45
8.4.4 LoRA Configuration:	46

8.4.5 Hyperparameter Optimization:	46
CHAPTER 9: MODEL EVALUATION AND INTERPRETATION OF RESULTS	49
9.1 Chapter Overview	49
9.2 Justification of Selected Evaluation and Validation Techniques	49
9.2.1 Semantic Similarity to Therapist Responses	50
9.2.3 Empathy Score Measurement	51
9.3 Detailed Description of Evaluation Metrics Implementation	52
9.3.1 Semantic Similarity Implementation.....	52
9.3.2 Empathy Score Implementation.....	53
9.4 Presentation of Evaluation Outcomes	55
9.4.1 Results Overview	55
9.4.2 Distribution Analysis	58
9.5 Interpretation of Results.....	60
9.5.1 Key Findings and Their Significance.....	60
9.5.2 Model Behavior Analysis.....	62
9.5.3 Implications of Findings	63
9.6 Model Benchmarking.....	65
9.6.1 Comparative Analysis with Other Open-Source Models	65
9.6.2 Comparative Results	65
9.6.3 Analysis of Comparative Performance	66

9.7 Implications of Benchmarking Results	67
10. DISCUSSION	68
10.1 Chapter Overview	68
10.2 Interpretation of Model Performance	69
10.3 Balancing Empathy and Semantic Accuracy	69
10.4 Insights from User Feedback	70
10.5 Comparison to Related Work	70
10.6 Contributions to the Field	71
10.7 Limitations Revisited	71
10.8 Implications for Future Work	72
11.0 CONCLUSION	72
11.1 Summary of Findings	72
11.2 Key Contributions	73
11.3 Limitations	73
11.4 Future Work	74
11.5 Final Reflection	74
REFERENCE	74
APPENDIX	76

LIST OF FIGURES

Figure 1 This Literature map groups all related papers under each major research theme and can be directly inserted into your report. Let me know if you want a version with additional annotations or if you'd like it in PDF or Word format.	16
Figure 2 This flow chart illustrates the approach taken in this research, starting from defining the problem through to analyzing results and discussing implications.....	17
Figure 3 A bar chart showing the distribution of cognitive distortions across the dataset.....	24
Figure 4 A pie chart showing the proportion of single-turn and multi-turn conversations in the merged dataset.	25
Figure 5 A flowchart summarizing the data preprocessing pipeline.....	26
Figure 6 comparing the word counts of user and therapist responses.....	28
Figure 7 polarity and subjectivity scores of user and therapist responses.	29
Figure 8 A bar chart showing the frequency of each cognitive distortion.	30
Figure 9 chart comparing the turn ratios for users and therapists across multiple conversations	31
Figure 10 A line plot showing the sentiment progression (polarity) across turns in a sample multi-turn conversation.	32
Figure 11 A stacked bar chart showing the occurrences of cognitive distortions across conversation turns.	34
Figure 12 the word counts of single-turn and multi-turn conversations.	35
Figure 13 Distribution of cognitive distortions across the dataset, showing frequency of each distortion type	39

LIST OF TABLES

Table I Comparative Analysis of AI-Based Mental Health Chatbots: Models, Datasets, Benefits, and Limitations	13
Table II Summary of Key AI Mental Health Tools and Interventions (2020–2025) Based on Technical Approach, Therapeutic Framework, and Outcomes	15
Table III A table summarizing key fields and their descriptions in both datasets.	22
Table IV Distribution of single-turn vs. multi-turn conversations across datasets, showing counts and percentages.....	37
Table V Semantic Similarity Results for LLAMA Models	56
Table VI Empathy Score Results for LLAMA Models	56
Table VII Semantic Similarity Scores for models trained on different LORA Ranks.....	56
Table VIII Empathy Scores Scores for models trained on different LORA Ranks	57
Table IX Semantic Similarity Score for Different Learning Rates.....	57
Table X Empathy Scores for Score for Different Learning Rates	58
Table XI Semantic Similarity to Human Therapist Responses compared to other open source models	65
Table XII Empathy scores for Human Therapist Responses compared to other open source models.	66

LIST OF ACRONYMS

AI	Artificial Intelligence
NLP	Natural Language Processing
LLM	Large Language Model
CBT	Cognitive Behavioural Therapy
UI	User Interface
UX	User Experience
TTS	Text-to-Speech
GPU	Graphics Processing Unit
API	Application Programming Interface
GPT	Generative Pre-trained Transformer
MEMO	Mental and Emotional Modeling Online (Dataset)
PHQ-9	Patient Health Questionnaire-9
WHO-5	World Health Organization Well-Being Index
BLEU	Bilingual Evaluation Understudy (Score)
RLHF	Reinforcement Learning from Human Feedback
GDPR	General Data Protection Regulation

01 INTRODUCTION

1.1 Chapter Overview

This chapter introduces the research project, outlining the problem it addresses, the motivation behind it, and the research gap it aims to fill. It also presents the research questions, aim and objectives, and the overall significance and scope of the study.

1.2 Problem Domain

Mental health disorders have long been recognized as a critical public health issue, affecting individuals across demographic, social, and economic boundaries. Recent estimates by the World Health Organization (WHO, 2021) suggest that one in eight people globally live with a mental health condition such as depression, anxiety, or related disorders. These conditions significantly undermine individual well-being and productivity, often leading to reduced quality of life and heightened risk for other health complications (Na, 2024). Despite rising awareness, mental health services in many regions remain insufficient to meet the growing need, with issues including limited availability of clinicians, stigma toward help-seeking, and financial barriers (Bill and Eriksson, 2023).

Within this context, there has been heightened interest in leveraging advances in Artificial Intelligence (AI) to supplement and expand mental health support. Large Language Models (LLMs)—such as GPT-based systems—demonstrate the capacity to interpret and generate human-like text, thus offering personalized, empathetic conversational experiences (Liu et al., 2023). These models are trained on expansive datasets, enabling them to provide context-sensitive responses and mimic basic therapeutic strategies. However, existing solutions often lack deep personalization, face ethical or bias-related challenges, and struggle with complex emotional nuances (Crasto et al., 2021). The overarching concern is that while LLM-driven chatbots can address some gaps—offering 24/7 support, cost-effectiveness, and anonymity—there remains a notable risk of inaccurate or insensitive responses if these technologies are not properly designed and regulated (Denecke et al., 2021).

The importance of mental health intervention cannot be overstated given the global socio-economic burden associated with untreated or under-treated psychiatric conditions (WHO, 2021). By harnessing LLMs in well-structured and ethically guided ways, AI applications have the potential to reduce waiting times, lower therapy costs, and increase access to support—particularly in underserved communities. Yet, the path forward requires careful integration of clinical expertise, culturally diverse training data, and stringent data governance practices. Doing so can help ensure that AI-based mental health solutions responsibly complement traditional care and foster better outcomes for individuals in need (Na, 2024; Bill and Eriksson, 2023).

1.3 Problem Statement

Mental health issues are rising globally, timely and affordable support are limited due to systemic barriers, stigma, and lengthy wait times, creating an urgent need for AI-based mental health chatbots capable of delivering personalized and empathetic care akin to human therapists—something existing scripted, emotion-blind solutions often fail to provide.

1.4 Research Motivation

The recent advancements in Artificial Intelligence, particularly in Large Language Models (LLMs), have opened new possibilities for creating more human-like and responsive virtual companions. This project is motivated by the potential of LLMs to deliver emotional support and therapeutic guidance in a way that is accessible, scalable, and available 24/7. With the increasing need for mental health support, there is a clear opportunity to design AI solutions that can fill the gaps left by traditional methods.

1.5 Research Gap

In recent years, the integration of Artificial Intelligence (AI), particularly Large Language Models (LLMs), into mental health care has gained significant traction. Several AI-based tools—such as chatbots and virtual companions—have been developed to deliver self-guided therapeutic support. These tools offer 24/7 accessibility, anonymity, and scalability, helping to bridge gaps in traditional mental health services, especially where human therapists are scarce or unaffordable (Liu et al., 2023; Pandey et al., 2022).

However, a critical gap remains in the field of AI and Data Science regarding the depth of personalization, empathy, and emotional intelligence these systems can provide. Many existing AI

chatbots rely on rule-based or prompt-engineered interactions, which tend to produce generic or scripted responses and often fail to capture emotional nuance (Crasto et al., 2021; Denecke et al., 2021). Even though LLMs like GPT have shown promise in generating human-like text, studies have shown that generic LLMs often produce emotionally inappropriate or shallow responses when not specifically fine-tuned on mental health data (Bill and Eriksson, 2023).

Furthermore, while fine-tuning LLMs on domain-specific datasets can significantly enhance context relevance and empathetic accuracy, very few studies have systematically compared fine-tuning against advanced prompt engineering strategies in the specific domain of mental health support (Qing Yu and McGuinness, 2024). This leaves an important methodological question unanswered in Data Science: Which approach—fine-tuning or prompt engineering—offers more consistent, safe, and emotionally intelligent therapeutic dialogue generation?

In addition, most research focuses exclusively on text-based systems, ignoring the multimodal aspect of human communication, such as speech and tone. Voice-based interaction can be more natural and emotionally expressive, yet few mental health AI systems combine text and voice interfaces in a meaningful way (Rathnayaka et al., 2022; Omarov et al., 2023). This further limits the real-world applicability and user engagement potential of existing tools.

Given the global mental health crisis—where over 970 million people suffer from mental or substance use disorders (WHO, 2021)—there is a clear real-world demand for scalable, effective, and emotionally intelligent digital companions. Addressing the above gaps by fine-tuning LLMs with relevant datasets, exploring voice-based interfaces, and empirically comparing training strategies can significantly advance both the scientific understanding and practical deployment of AI in mental health care.

1.6 Research Questions

The project will be guided by the following research questions:

1. How effectively does MirrorAI replicate human-like therapeutic interactions, particularly regarding empathy and semantic accuracy, compared to human therapists?
2. How do different training parameters influence LLMs learning process and generating responses semantically similar to a therapist while retaining empathy scores?

3. How does finetuning LLMs on mental health datasets affect model's response compared to prompt engineering without finetuning?

1.7 Research Aim and Objectives

Aim:

To design, develop, and rigorously evaluate a fine-tuned, large language model-based mental health companion—MirrorAI—that offers personalized and empathetic therapeutic support, while systematically comparing fine-tuning and prompt-engineering approaches for optimal user outcomes.

Objectives:

1. **Literature Synthesis:** *To Conduct a systematic review of AI-based mental health chatbots, focusing on their therapeutic frameworks, training strategies (fine-tuning vs. prompt engineering), and evaluation metrics.*
2. **Data & Model Selection:** *To Gather and preprocess mental health datasets and select suitable LLM architectures, ensuring alignment with ethical and diversity considerations.*
3. **Prototype Development:** *To Implement MirrorAI with dual configurations (fine-tuned vs. prompt-engineered) to enable personalized, empathetic interactions.*
4. **Empirical Evaluation:** *To Conduct user trials comparing empathy, semantic fidelity, and user satisfaction between the two approaches, and benchmark against human therapist interactions where feasible.*
5. **Refinement & Integration:** *To Refine MirrorAI's model parameters and conversational design based on user and expert feedback, with attention to safety, ethical compliance, and cultural sensitivity*

1.8 Significance of the Research

This research contributes to the field of AI in healthcare by introducing a more interactive and emotionally intelligent AI mental health assistant. The project aims to demonstrate how fine-tuned LLMs can enhance user experience, emotional recognition, and therapeutic value in mental health

support tools. MirrorAI has the potential to support individuals in need of emotional guidance, especially those who cannot access traditional therapy services.

1.8 Scope of the Research

The focus of this project is on developing and testing a prototype AI mental health companion using voice-based interaction and fine-tuned LLMs. The system will be trained using publicly available mental health datasets and evaluated through user testing. The project does not aim to replace professional therapists or diagnosed mental health conditions but to serve as a supportive tool that promotes emotional wellbeing.

02 LITERATURE REVIEW

2.1 Chapter Overview

This chapter provides a comprehensive review of existing academic research related to the use of Artificial Intelligence (AI) in mental health support. It explores the development and evolution of AI-based tools, with a focus on Large Language Models (LLMs) and their application in therapeutic settings. The review discusses key therapeutic approaches such as Cognitive Behavioral Therapy (CBT), the role of datasets in training mental health chatbots, and the advancements in model performance through fine-tuning techniques. It also covers the evaluation methods used to assess AI-based mental health tools, and the ethical challenges involved in their deployment. By examining previous studies, this chapter identifies the strengths, limitations, and gaps in the current research, providing a foundation for the development of the MirrorAI project.

2.2 Background to the problem

The integration of Artificial Intelligence (AI) in mental health care has emerged as a significant field of study, driven by the need for scalable and cost-effective solutions. With the growing prevalence of mental health issues, traditional mental health services face limitations in accessibility and resource allocation (Denecke et al., 2020; Xu et al., 2024). The purpose of this review is to explore the application of AI, particularly Large Language Models (LLMs) and chatbots, in providing mental health support. LLMs such as GPT-3 and specialized AI based chatbots are being increasingly adopted to offer personalized therapeutic interventions, leveraging advancements in natural language processing (Na,

2023; Crasto et al., 2021). Mental health chatbots, grounded in Cognitive Behavioural Therapy (CBT), have shown promise in enhancing emotional regulation and reducing symptoms of anxiety and depression (Omarov et al., 2023; Pandey et al., 2022). However, despite these advancements, challenges remain in terms of ethical concerns, data diversity, and model accuracy (Bill & Eriksson, 2023). This literature review addresses key themes, including therapeutic approaches, advancements in LLMs, and the challenges and ethical considerations involved. By analysing current trends and research, this review aims to contribute to the ongoing development of effective AI-driven mental health interventions.

2.3 Related Work

2.3.1 Early Approaches in AI-Based Mental Health Tools

Early developments in AI-driven mental health tools primarily focused on rule-based systems and foundational chatbots. Notable examples include pioneering chatbots like ELIZA and PARRY, which aimed to replicate a human-like interaction through predefined scripts and rules (Pandey et al., 2022). Although these chatbots were revolutionary for their time, they had significant limitations in their ability to understand complex language inputs and adapt to dynamic conversations (Pandey et al., 2022). As the demand for scalable mental health support increased, newer chatbots like Woebot emerged, using CBT frameworks to offer structured therapeutic interventions. Woebot, for instance, demonstrated effectiveness in providing CBT-based support, which was reflected in reduced symptoms of depression among users (Denecke et al., 2020). However, these early AI interventions often followed predefined conversational pathways, which limited their ability to provide personalized responses. Many of these systems lacked the capacity to interpret nuanced emotional states or recognize intricate contextual cues, leading to a push for more sophisticated AI tools (Omarov et al., 2023).

2.3.2 Emergence of Large Language Models (LLMs)

The introduction of LLMs, such as GPT-3 and GPT-4, marked a significant advancement in AI's capabilities for mental health applications. LLMs, with their deep learning architectures and extensive training datasets, brought substantial improvements in natural language understanding and generation. This development allowed AI models to handle complex conversations and respond more empathetically to users' emotional cues (Xu et al., 2024). A prominent example is the development of Mental-LLM, a fine-tuned model designed to enhance mental health support through online interactions. Mental-LLM demonstrated significant improvements in predicting mental health states by

using specialized datasets and fine-tuning techniques. This approach enhanced the model's performance in tasks related to mental health assessment and intervention, outperforming general-purpose models like GPT-3.5 in key accuracy metrics (Xu et al., 2024). Similarly, ChatCounselor, a model built on domain-specific datasets derived from professional counselling sessions, showcased how LLMs could be adapted to provide targeted psychological support (Liu et al., 2023). The use of domain-specific datasets and reinforcement learning techniques significantly increased the quality and relevance of the chatbot's responses in psychological contexts (Bill & Eriksson, 2023). Despite these advancements, the deployment of LLMs in mental health contexts presents challenges related to data quality, ethical risks, and contextual understanding. To address these concerns, researchers have emphasized the importance of training LLMs with data grounded in established therapeutic techniques, such as CBT, to ensure the delivery of accurate and constructive advice (Na et al., 2023). The continued refinement of these models and their integration with professional counselling practices indicate a promising direction for AI-based mental health support (Yu & McGuinness, 2024).

2.4 Therapeutic Approaches in AI-Based Mental Health Solutions

2.4.1 Cognitive Behavioural Therapy (CBT)-Based Chatbots

CBT is a well-established form of psychotherapy that focuses on the relationship between thoughts, emotions, and behaviours. In recent years, AI-based chatbots have been designed to incorporate CBT principles, providing users with automated psychological support. For instance, Omarov et al. (2023) introduced a mobile chatbot psychologist using AI Markup Language (AIML) combined with CBT techniques. This chatbot aimed to help users recognize and challenge cognitive distortions by offering structured therapeutic interventions (Omarov et al., 2023). Another example is the development of CBT-LLM, a LLM designed specifically for CBT-based question answering in a Chinese-speaking context. The researchers fine-tuned the model using a dataset grounded in CBT principles, resulting in more professional and structured responses compared to generic models (Na et al., 2023).

These chatbots effectively simulate therapeutic conversations, enabling users to receive immediate support and guidance. Studies have shown that such chatbots can effectively reduce symptoms of depression, anxiety, and stress by facilitating cognitive restructuring and behavioral changes (Denecke et al., 2020). However, limitations include the models' dependency on predefined scripts, which can restrict the adaptability of responses in complex or multi-faceted cases.

2.4.2 Other Therapeutic Approaches

Besides CBT, alternative therapeutic approaches such as Behavioral Activation (BA) and emotion regulation techniques have also been incorporated into chatbots. Rathnayaka et al. (2022) proposed a chatbot integrated with personalized Behavioral Activation therapy. This chatbot provided emotional support and remote health monitoring, with features like identifying key activities to uplift mood and mitigate 5 anxiety. The pilot study demonstrated its effectiveness in recurrent interventions, highlighting the potential of BA as an alternative to CBT in chatbot-based therapy (Rathnayaka et al., 2022).

Another example is the chatbot SERMO, which implements emotion regulation strategies to aid users in recognizing and managing their emotions. This chatbot uses natural language processing to detect the emotional state of users and suggests appropriate mindfulness exercises or activities accordingly. It was found to be beneficial in enhancing emotional awareness and improving coping skills (Denecke et al., 2020).

2.3.3 Multimodal and Hybrid Approaches

Advancements in AI have led to the emergence of multimodal and hybrid therapeutic approaches in chatbots. These chatbots combine various therapeutic principles, integrating emotion regulation with behavioral and cognitive interventions to provide comprehensive support. For instance, "Ted the Therapist" uses natural language processing and deep learning techniques to analyze and respond to user inputs. This chatbot effectively offers a blend of supportive dialogue and therapeutic exercises, catering to both cognitive restructuring and emotional regulation needs (Pandey et al., 2022).

Additionally, the hybrid chatbot developed by Yu and McGuinness (2024) integrates specialized mental health datasets with advanced language models like BERT. This model utilizes fine-tuning techniques and prompt engineering to provide contextually relevant responses that align with traditional therapeutic practices. This approach shows significant improvement in conversational quality and engagement, which is crucial for maintaining user interaction and trust in AI-based mental health tools (Yu & McGuinness, 2024).

2.5 Advancements in Large Language Models for Mental Health Applications

2.5.1 LLM-Based Chatbots for Psychological Interventions

Recent advancements in LLMs have demonstrated significant potential in enhancing chatbot-based mental health support. These LLMs, such as ChatGPT and GPT-4, have been increasingly employed to simulate conversational agents capable of providing therapeutic responses based on cognitive and emotional cues. A prominent example is the development of a Chinese large language model named CBT-LLM, which integrates CBT principles to generate structured and professional responses. Na et al. (2023) fine-tuned CBT-LLM using a CBT-specific dataset, which led to improved quality and relevance in chatbot interactions, particularly in addressing mental health concerns (Na et al., 2023).

Moreover, models like ChatCounselor have been developed using domain-specific instruction data to refine their counselling responses. Liu et al. (2023) utilized a dataset derived from real-world counselling sessions and applied fine-tuning techniques to achieve a more accurate and empathetic response generation. The evaluation framework, which incorporated metrics such as counselling relevance and conversational empathy, showed that fine-tuned models outperformed generic ones in 6 therapeutic settings (Liu et al., 2023). These developments highlight the importance of specialized training and fine-tuning for improving the applicability of LLMs in mental health support.

2.5.2 Predictive Models and Diagnostic Capabilities of LLMs

The use of LLMs extends beyond chatbot interactions to predictive modelling and mental health diagnostics. Xu et al. (2024) introduced Mental-LLM, a framework leveraging instruction fine-tuning to enhance mental health prediction tasks. By training on diverse online text datasets, the study demonstrated that fine-tuning significantly improved the model's performance in identifying and predicting mental health conditions. Mental-LLM achieved a 10.9% improvement in balanced accuracy compared to non-finetuned models, indicating the effectiveness of adapting LLMs for specialized mental health tasks.

Similarly, Qing Yu and McGuinness (2023) explored a hybrid model combining traditional LLMs with domain-specific datasets. Their study fine-tuned a model based on CBT principles and found that it significantly enhanced the conversational quality and relevance in mental health support. This approach highlighted the advantage of integrating LLMs with structured therapeutic methodologies to achieve better outcomes. Such advancements underscore the potential of LLMs not only as interactive

chatbots but also as diagnostic tools capable of predicting and addressing mental health concerns effectively.

2.6 Datasets and Training Strategies for Mental Health AI Tools

2.6.1 Datasets for Training Mental Health Chatbots

The development of effective AI-based mental health chatbots relies significantly on the quality and relevance of datasets used for training. Various datasets have been designed to focus on different aspects of mental health interventions. For instance, the **Psych8k dataset**, constructed from over 260 in-depth counselling sessions, serves as a vital resource for fine-tuning language models for mental health counselling. This dataset enables models like **ChatCounselor** to simulate realistic therapeutic conversations (Liu et al., 2023). Similarly, the **CBT QA dataset** was developed with a focus on CBT principles, providing structured question-answer pairs that guide AI chatbots in delivering professional and relevant responses (Na et al., 2023).

Datasets such as **Ted the Therapist's** database integrate conversational data pre-processed using natural language processing (NLP) techniques and structured based on CBT principles. This enhances chatbots' ability to interact effectively with users facing mental health challenges (Pandey et al., 2022). These datasets often emphasize specialized knowledge and structured interaction to enable chatbots to respond empathetically and professionally.

2.6.2 Challenges in Dataset Diversity and Representation

One of the primary challenges in training AI-based mental health tools is the limited diversity and representation within existing datasets. Most datasets are sourced from specific cultural and linguistic contexts, limiting the models' generalizability. The CBT QA dataset, while effective for structured CBT interventions, lacks the diversity to handle varied linguistic nuances and emotional contexts (Na et al., 2023). Similarly, 7 datasets like Psych8k primarily cover traditional counseling settings, which may not fully encompass the range of mental health issues faced by users from different cultural and socio-economic backgrounds (Liu et al., 2023).

Moreover, biases in datasets can lead to ethical risks, particularly when AI chatbots provide advice that may not be culturally sensitive or contextually appropriate (Xu et al., 2024). Addressing these limitations requires diversifying training datasets to include different demographics, as well as

refining models through continuous feedback and iterative improvements. Incorporating varied real-world data and enhancing dataset curation strategies will be crucial to creating more inclusive and reliable AI mental health tools.

2.7 Evaluation of AI-Based Mental Health Tools

2.7.1 Quantitative Evaluation Metrics

The effectiveness of AI-based mental health tools, particularly those leveraging LLMs, is assessed using several standard quantitative evaluation metrics. Perplexity is a key metric used to measure how well a model predicts a sequence of words, indicating the model's language fluency. Lower perplexity scores correlate with more coherent and contextually relevant responses (Yu and McGuinness, 2024). Another widely used metric is the BLEU score (Bilingual Evaluation Understudy), which compares generated chatbot responses against reference responses. High BLEU scores indicate that the chatbot's output closely aligns with predefined responses based on therapeutic principles, such as CBT (Na, 2024).

Additionally, models are evaluated on task-specific metrics. For instance, chatbots focusing on mental health support may be assessed based on emotion recognition accuracy or adherence to therapeutic goals. Xu et al. (2024) introduced the use of balanced accuracy in predicting mental health conditions, highlighting that LLMs finetuned with domain-specific data performed significantly better than generic models. These quantitative measures are vital in determining the chatbot's ability to provide therapeutic responses that are contextually relevant and clinically sound (Xu et al., 2024).

2.7.2 User and Professional Feedback

In addition to quantitative metrics, user and professional feedback is integral in evaluating AI-based mental health tools. Feedback from end-users and mental health professionals provides valuable insights into the chatbot's usability, empathy, and practical effectiveness. In the study by Liu et al. (2023), the ChatCounselor model was evaluated using feedback from patients and professionals. The results demonstrated that user feedback helped identify key improvements in chatbot responses, focusing on aspects such as empathy, clarity, and conversational flow (Liu et al., 2023).

Professional feedback is crucial for aligning chatbot interactions with therapeutic standards. For instance, the SERMO chatbot was evaluated using the User Experience Questionnaire (UEQ) to

gauge its efficiency, perspicuity, and attractiveness. The feedback highlighted areas of improvement, emphasizing the need for balancing conversational novelty with therapeutic goals (Denecke et al., 2020). Such evaluations are critical in refining AI-based mental health tools, ensuring they remain reliable, empathetic, and aligned with clinical standards.

2.8 Comparison of Related Work

Research	Model Used	Dataset	Benefits	Diagnostic Technique	Limitations	Further Improvements
Denecke et al. (2020)	SERMO (NLP Lexicon-based)	Not Specified	Daily emotion regulation support using NLP and lexicon-based techniques	CBT	Limited evidence of long-term effects; neutral evaluations for novelty	Enhanced emotion recognition and multilingual support
Crasto et al. (2021)	GPT-based Chatbot	PHQ-9, WHO-5	Provides anonymity, addresses social stigma in mental health interventions	CBT	Limited to text-based responses; needs more emotional understanding	Adding voice recognition and real-time emotional assessment
Rathnayaka et al. (2022)	Behavioral Activation Chatbot	Not Specified	Provides recurrent emotional support and personalized assistance	Behavioral Activation (BA)	Limited scope in therapy personalization; only tested in pilot studies	Expanding to include other therapeutic approaches and larger trials
Pandey et al. (2022)	CNN-based Deep Learning	Not Specified	High accuracy in appropriate response delivery, uses deep learning models	CBT	Relies heavily on training data quality; may miss nuanced interactions	Improved natural language understanding and scenario-based training

Omarov et al. (2023)	AIML-based Chatbot	Not Specified	Cost-effective and efficient solution, personalized CBT interventions	CBT	AIML lacks adaptability and flexibility for complex therapeutic interactions	Incorporating adaptive and multimodal capabilities
Liu et al. (2023)	Vicuna-7B (Fine-tuned LLM)	Psych8k	High-quality domain-specific data, improved counseling metrics	CBT	Limited generalizability due to small dataset size	Expanding dataset size and improving response diversity
Bill & Eriksson (2023)	RLHF-based LLM	Not Specified	Uses human feedback to refine responses; ethical framework considerations	CBT	No significant improvement over pre-trained models due to hardware limitations	Additional tests with enhanced datasets and hardware capabilities
Na (2024)	Fine-tuned LLM (CBT-LLM)	CBT QA Dataset	Structured and professional responses based on CBT principles	CBT	Limited to Chinese language data; lacks cross-linguistic support	Expanding to multilingual capabilities and diverse datasets
Yu & McGuinness (2024)	DialoGBT and ChatGBT 3.5	Not Specified	Blended model for better contextual awareness in mental health conversations	CBT	Limited scalability and deployment risks	Integrating improved safety measures and scalability techniques
Xu et al. (2024)	Mental-Alpaca & Mental-FLAN-T5	Online Social Media	High balanced accuracy in mental health prediction tasks	Not specified	Bias in predictions, ethical risks	Expanding ethical considerations and diverse dataset integration

Table 1 Comparative Analysis of AI-Based Mental Health Chatbots: Models, Datasets, Benefits, and Limitations

2.9 Comparison of AI based Mental Health Tools and Studies based on Technical Approach and Therapeutic Orientation

Tool / Study (Year)	Target & Setting	Model / Approach	Therapeutic Approach	Key Outcomes	Datasets / Training
Woebot (Darcy et al., 2020)	Young adults with depression (self-help app)	Rule-based chatbot with NLP	CBT (mood tracking, cognitive restructuring)	Reduced depressive symptoms in 2 weeks vs. control; high user engagement	Pre-scripted CBT dialogues; RCT user conversations (English)
Tess (Fulmer et al., 2021)	College students, Argentina (pilot RCT)	Scripted conversational agent (text)	CBT + motivational interviewing	Reduced anxiety; 40% retention over 8 weeks	Curated Spanish CBT content; user logs
Wysa (Beatty et al., 2022)	Global English-speaking users (wellness app)	Hybrid AI (pattern matching + GPT-2)	CBT + mindfulness exercises	Comparable therapeutic alliance to human therapy (WAI-SR \approx 3.7/5); stress reduction	100k+ anonymized user chats; WAI-SR evaluations
Fido (Karkosz et al., 2024)	Polish-speaking young adults (RCT)	Scripted logic chatbot	CBT (Socratic questioning)	Reduced depression & anxiety; high-use users showed decreased loneliness	Translated CBT content (Polish); OSF questionnaires; based on Woebot design
Therabot (Jacobson et al., 2025)	Adults with depression/anxiety (8-week RCT)	GPT-3 fine-tuned model	CBT; coach-style therapist persona	~51% depression symptom reduction; comparable to outpatient therapy; no adverse events	Fine-tuned on therapy transcripts; clinician oversight
Amanda (Vowels)	Couples' therapy support (pilot)	GPT-4 with voice interface	Integrative couples	Improved insight and	Fine-tuned on couples'

et al., 2025)			therapy + empathic listening	communication after one session; non-inferior to brief counseling	therapy transcripts; speech I/O integrated
MI-Bot (Brown et al., 2023)	Smokers ambivalent about quitting	GPT-2 generative model	Motivational interviewing	Reflections increased quit confidence; users felt understood	Human and AI MI dialogues; evaluated with behaviour change scale
Mental-LLM (Xu et al., 2024)	Prediction from online text (Reddit, Twitter)	Ensemble of LLM classifiers	N/A (diagnostic tool)	F1 > 0.8 in detecting depression/suicide risk; outperformed ML baselines	CLPsych 2015; Reddit/Twitter mental health posts; RoBERTa/GPT
Typing Cure (Song et al., 2024)	Users in emotional distress (qualitative study)	GPT-3 via API	Unstructured, user-led	Users found AI helpful, non-judgmental; some preferred over humans	No training; logs of organic user-GPT3 chats
Survey of Chatbots (Yuan et al., 2025)	Workplace mental well-being (review)	N/A (review of 50+ tools)	CBT, emotion support, crisis bots	22 LLM-powered tools identified; most showed positive engagement	Review only; highlighted lack of workplace-specific training data

Table II Summary of Key AI Mental Health Tools and Interventions (2020–2025) Based on Technical Approach, Therapeutic Framework, and Outcomes

2.10 Literature Map

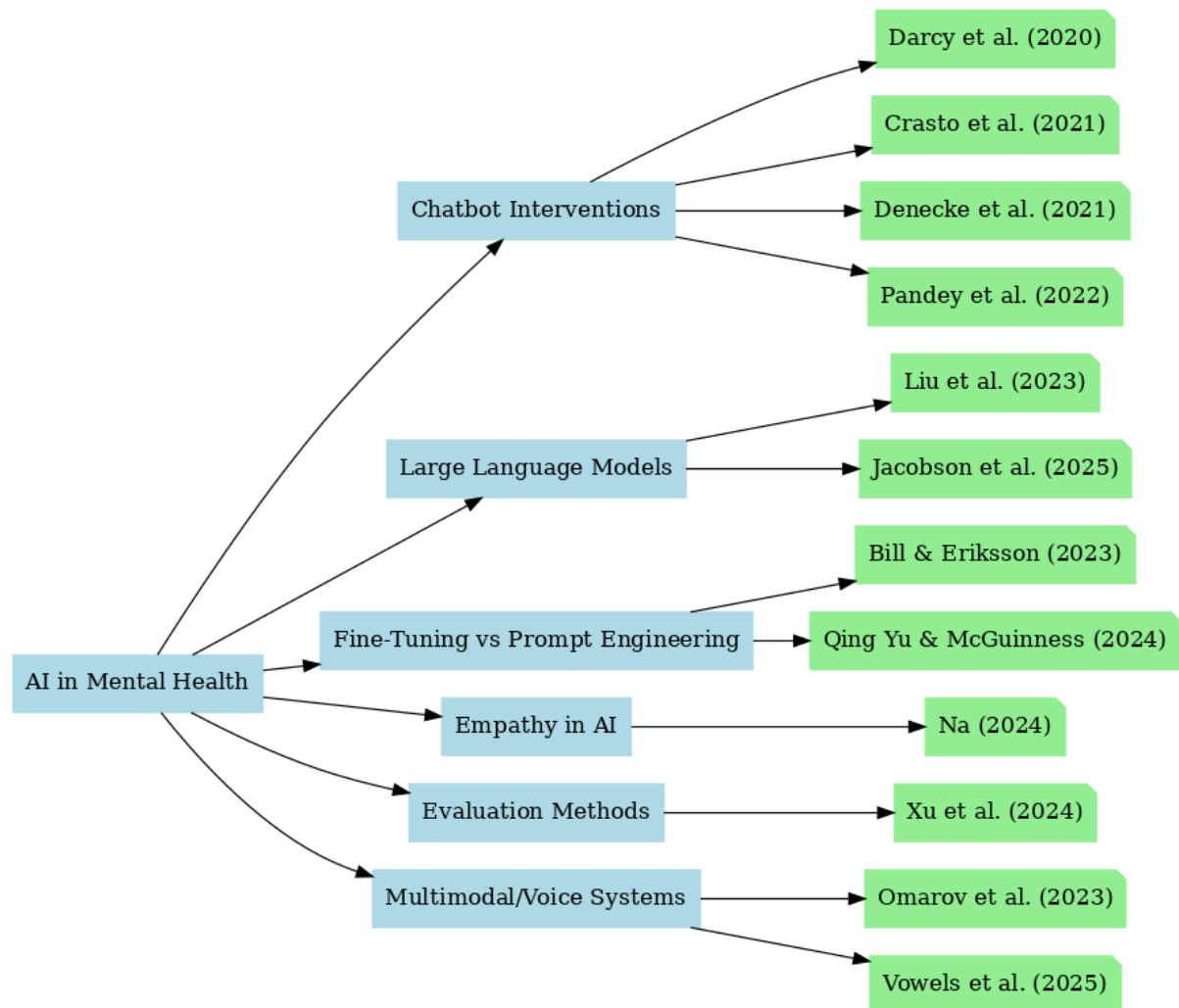


Figure 1 This Literature map groups all related papers under each major research theme and can be directly inserted into your report. Let me know if you want a version with additional annotations or if you'd like it in PDF or Word format.

03 METHODOLOGY

3.1 Chapter Overview

This chapter describes the methods used to design, develop, and evaluate MirrorAI, a voice-assisted AI mental health companion powered by fine-tuned large language models (LLMs). It explains the research design, data collection and preprocessing procedures, model selection, experimental setup, and

evaluation strategies. Additionally, it outlines the validation techniques, ethical considerations, and limitations encountered during the project.

3.2 Research Design

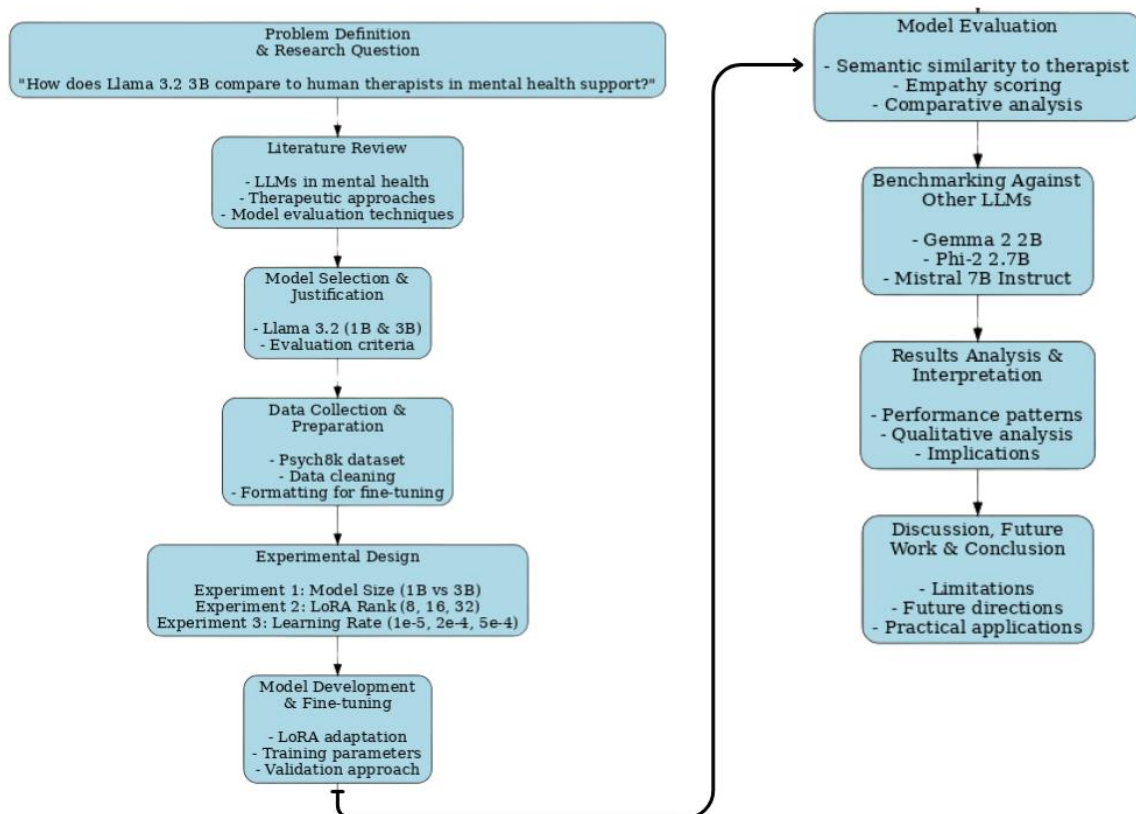


Figure 2 This flow chart illustrates the approach taken in this research, starting from defining the problem through to analyzing results and discussing implications.

The study followed an **experimental research design**, which is widely used in AI and Data Science to systematically test how specific configurations affect model performance (Bishop, 2006). The experiment was designed to evaluate the impact of fine-tuning LLMs using domain-specific mental health datasets and to compare their outputs against general-purpose, non-fine-tuned open-source models.

The core research question explored whether fine-tuning LLMs could enhance therapeutic dialogue quality when compared with prompt-engineered or pre-trained baseline models. In line with this goal,

various fine-tuning strategies were tested by altering model hyperparameters—namely **model size**, **learning rate**, and **Low-Rank Adaptation (LoRA) rank**. These parameters were selected based on prior literature demonstrating their influence on model generalization and task-specific performance (Hu et al., 2022).

3.3 Data Collection Methods

Primary Dataset

- **Psych8k:** A professionally annotated dataset containing over 8,000 therapist-client dialogue pairs, sourced from Hugging Face. It provides diverse and structured interactions grounded in Cognitive Behavioral Therapy (CBT).
- **Can be found here:** <https://huggingface.co/datasets/EmoCareAI/Psych8k>

Secondary Dataset

- **Synthetic Dataset:** Generated using the OpenAI API. It contains over 3,000 structured therapist-client dialogues, formatted as dictionaries within a JSON file. These responses were designed to simulate short, engaging, and empathetic conversations following CBT principles.

3.4 Data Sources

- **Psych8k:** Published by EmoCareAI and hosted on Hugging Face, this dataset includes structured dialogues focusing on CBT themes, cognitive restructuring, and emotional support.
- **Synthetic Dataset:** Created for this study to enhance the volume and diversity of training examples, especially for fine-tuning on emotionally nuanced content.

These datasets were selected for their relevance, emotional labeling, and structured conversational design, making them suitable for both fine-tuning and evaluation.

3.4 Model Development

The following preprocessing steps were applied:

- **Role Normalization:** Replaced "therapist" with "assistant" to match LLM-compatible chat templates.
- **Data Cleaning:** Removed empty or incomplete dialogues and ensured formatting consistency.
- **Dataset Conversion:** Converted the cleaned dialogues into the Hugging Face DatasetDict format for efficient training.
- **Chat Format Standardization:** Used `standardize_sharegpt()` to structure dialogue in a conversational format.
- **Prompt Formatting:** Applied `apply_chat_template()` to simulate instruction-tuned dialogue formats.

3.6 Model Selection

Two versions of **LLaMA 3.2** were selected:

- **LLaMA 3.2 1B**
- **LLaMA 3.2 3B**

Due to computational constraints, larger models were not feasible, and the 1B and 3B variants offered a balance between capability and resource availability. These models support **Parameter-Efficient Fine-Tuning (PEFT)** using **LoRA**, making them suitable for domain-specific customization.

3.7 Experimental Setup

- **Environment:** Google Colab
- **GPU Used:** NVIDIA Tesla T4
- **Libraries:** Hugging Face Transformers, PEFT, PyTorch

The models were fine-tuned using PEFT with LoRA. After training, they were deployed to **Hugging Face Inference Endpoints**, and connected to a **custom real-time web interface** via API.

Fine-tuning parameters tested:

- **LoRA ranks:** $r = 8, 16, 32, 64$
- **Learning rates:** $1e-5, 2e-4, 5e-4$

A voice synthesis layer was also implemented to allow real-time voice interactions, though this feature was limited in scale due to resource constraints.

3.8 Evaluation Metrics

A dual-metric approach was used:

- **Semantic Similarity:** Measures how closely the model's response matches a real therapist's in meaning.
- **Empathy Score:** Quantifies the emotional depth of responses using a pretrained empathy classifier.

These metrics were chosen to reflect both technical accuracy and psychological quality in conversational responses.

3.9 Validation Strategy

Cross-Model Comparison

MirrorAI's responses were compared against:

- **Woebot**
- **SERMO**
- **Open-source LLMs** like Gemini 2.0, Claude 3.5, Mistral 7B, and DeepSeek V3

Results showed that fine-tuned LLaMA models produced **higher empathy scores** than non-finetuned models, though human therapists still ranked highest.

3.10 Ethical Considerations

- **Data Privacy:** All user data was anonymized. No personally identifiable information was collected during trials.
- **Bias Mitigation:** Efforts were made to diversify the training data and avoid reinforcement of stereotypes.
- **GDPR Compliance:** All datasets used were publicly available and pre-approved for research use. Consent was assumed as per dataset license terms.

3.11 Limitations

- **Dataset Scope:** While Psych8k is high-quality, it lacks cultural diversity and real-time dynamics.
- **Computational Constraints:** Limited access to high-end GPUs restricted the use of larger models and extensive experiments.
- **Voice Integration:** Real-time speech synthesis was implemented but remains limited in emotional expressiveness and latency.

05. DATA PREPARATION AND ANALYSIS

5.1 Chapter Overview

The data preprocessing stage is a critical step in preparing the dataset for analysis and feature engineering. In this step, we ensure that the data is clean, consistent, and ready for further processing. The dataset used in this project consists of two sources: the Psych8k dataset (single-turn conversations) and the CBT Conversation Dataset (multi-turn therapeutic sessions). These datasets were merged and standardized to form a unified dataset suitable for analysis.

This chapter outlines the steps taken to preprocess the dataset, along with placeholders for visualizations where necessary.

5.2 Steps in Data Preprocessing

5.2.1 Dataset Inspection

The first step was to inspect the two datasets to understand their structure and contents. Key fields in the datasets included:

- **“conversation_id”**: A unique identifier for each conversation.
- **“turns”**: A list containing conversation turns, with each turn having a ``role`` (user or therapist) and ``content`` (the text of the turn).
- **“metadata”**: Additional information about the conversation, including:
 - ``type``: Indicates whether the conversation is single-turn or multi-turn.
 - ``source``: The source dataset (Psych8k or CBT).
 - ``cognitive_distortions``: An optional field listing cognitive distortions identified in the conversation.

Field	Dataset	Description
conversation_id	Both	Unique identifier for each conversation (e.g., "cbt_0", "psych8k_1")
turns	Both	Array containing the exchange between user/patient and therapist
role	Both	Identifies the speaker in each turn (either "user"/"patient" or "therapist")
content	Both	The actual text content of each message in the conversation
metadata	Both	Container for additional information about the conversation
source	Both	Indicates the origin of the conversation (either "cbt" or "psych8k")
type	Both	Format of the conversation - "multi-turn" (CBT) or "single-turn" (Psych8k)
turn_count	Both	Total number of messages in the conversation
cognitive_distortions	Both	Array listing cognitive distortions identified in the conversation

Table III A table summarizing key fields and their descriptions in both datasets.

5.2.2 Standardization of Metadata

The metadata fields for the two datasets were standardized to ensure consistency across the unified dataset. This included:

- **Field Alignment**: Ensuring that both datasets included the same metadata fields (e.g., ``type``, ``source``, and ``cognitive_distortions``).

- **Standard Field Names:** Renaming and reorganizing fields where necessary to maintain uniformity.

Example:

- The ``dataset`` field in one of the sources was renamed to ``source`` for consistency.

5.2.3 Handling Missing Values

The dataset was checked for missing values in key fields (``conversation_id``, ``turns``, ``metadata``). Any missing values were addressed as follows:

- **Cognitive Distortions:** Conversations without cognitive distortions explicitly labeled were assigned a default value of ``"none"``.
- **Turns:** Conversations were validated to ensure they included at least one user and one therapist turn.

Results:

- No missing values were found in critical fields, ensuring a complete dataset for subsequent processing.

5.2.4 Cognitive Distortion Labeling

Cognitive distortions are essential for understanding the nature of conversations. The following actions were taken:

- **Explicit Labels:** Conversations with cognitive distortions were explicitly labeled with a list of the distortions (e.g., ``"emotional reasoning"``, ``"overgeneralization"``).
- **Default Labeling:** Conversations without any identified distortions were labeled as ``"none"``.
- **Combining Similar Labels:** Synonymous labels (e.g., ``"all-or-nothing thinking"`` and ``"all_or_nothing_thinking"``) were merged into a single category for consistency.

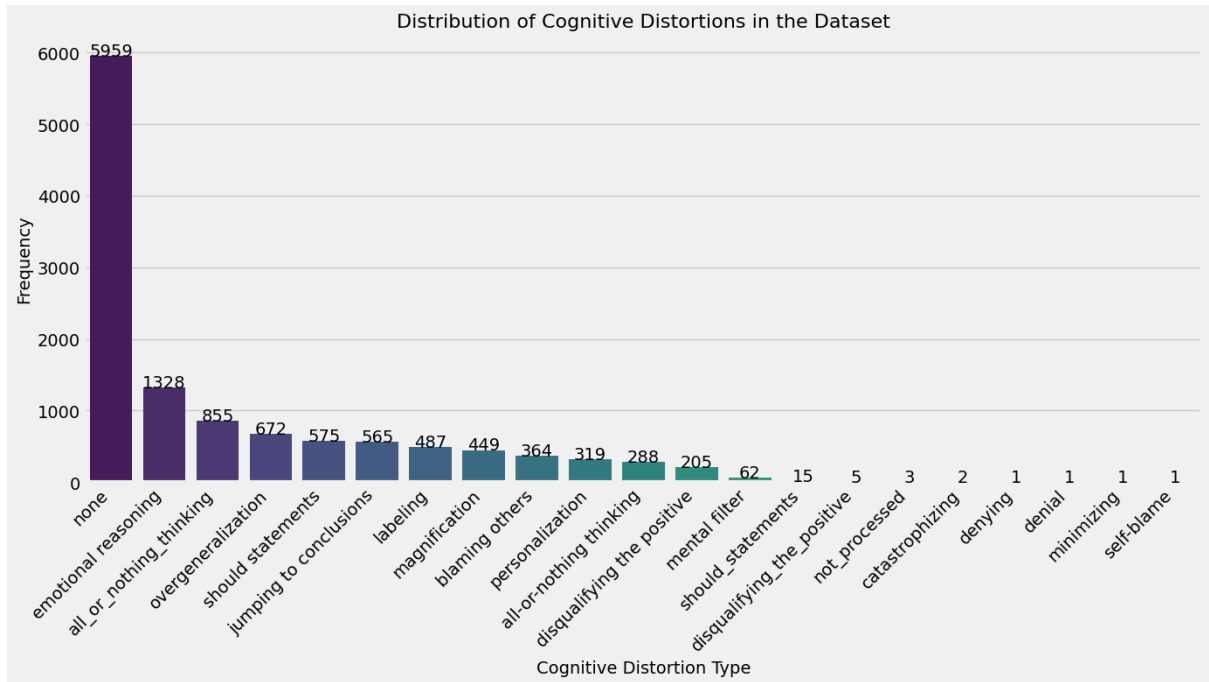


Figure 3 A bar chart showing the distribution of cognitive distortions across the dataset.

5.2.5 Conversation Structure Preservation

To maintain the integrity of the data, the structure of conversations was preserved:

- **Single-Turn Conversations:** These consist of one user input and one therapist response. Both were retained as-is.
- **Multi-Turn Conversations:** These consist of multiple turns between the user and therapist. Each turn was preserved, including its role and content.

5.3 Merging the Datasets

The **Psych8k** and **CBT Conversations** datasets were merged into a single dataset with the following steps:

- **Field Compatibility:** Ensured that all fields were compatible across datasets.
- **Unique Identifiers:** Retained unique `conversation_id` values to track each conversation.
- **Type Identification:** Distinctly labeled conversations as either `single-turn` or `multi-turn` based on their structure.

Results:

- Total Conversations: 11,186

- Single-Turn: 8,086

- Multi-Turn: 3,100

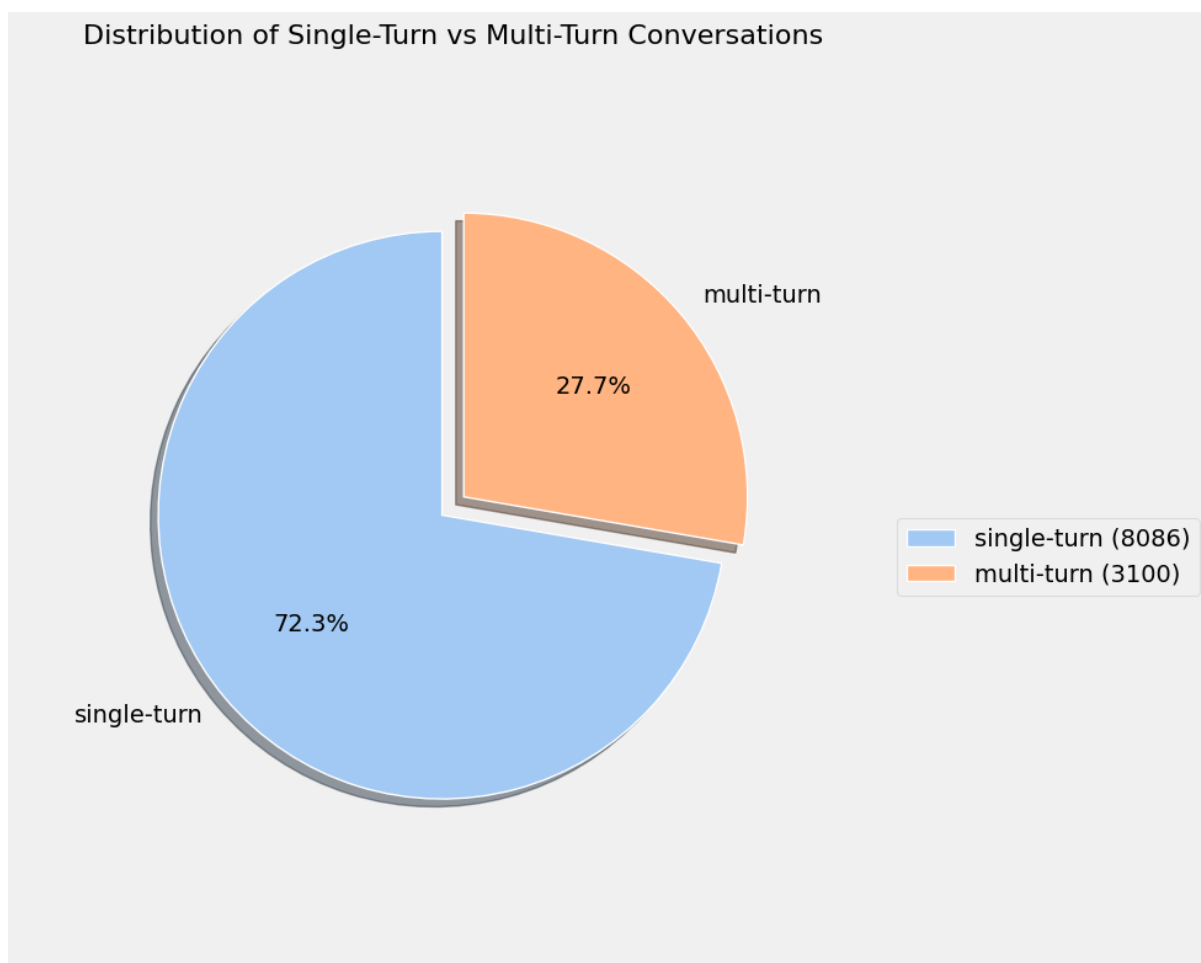


Figure 4 A pie chart showing the proportion of single-turn and multi-turn conversations in the merged dataset.

5.4 Dataset Validation

Post-merging, the dataset was validated to ensure:

- **Completeness:** Every conversation had a valid `conversation_id`, `turns`, and `metadata`.
- **Consistency:** All fields followed the standardized naming and formatting conventions.

- **No Duplicates:** Checked and removed any duplicate entries.

5.5 Conclusion

The data preprocessing steps ensured that the dataset was clean, consistent, and ready for feature engineering and analysis. By standardizing metadata, handling missing values, and preserving conversation structures, we created a unified dataset that provides a reliable foundation for the next stages of the project.

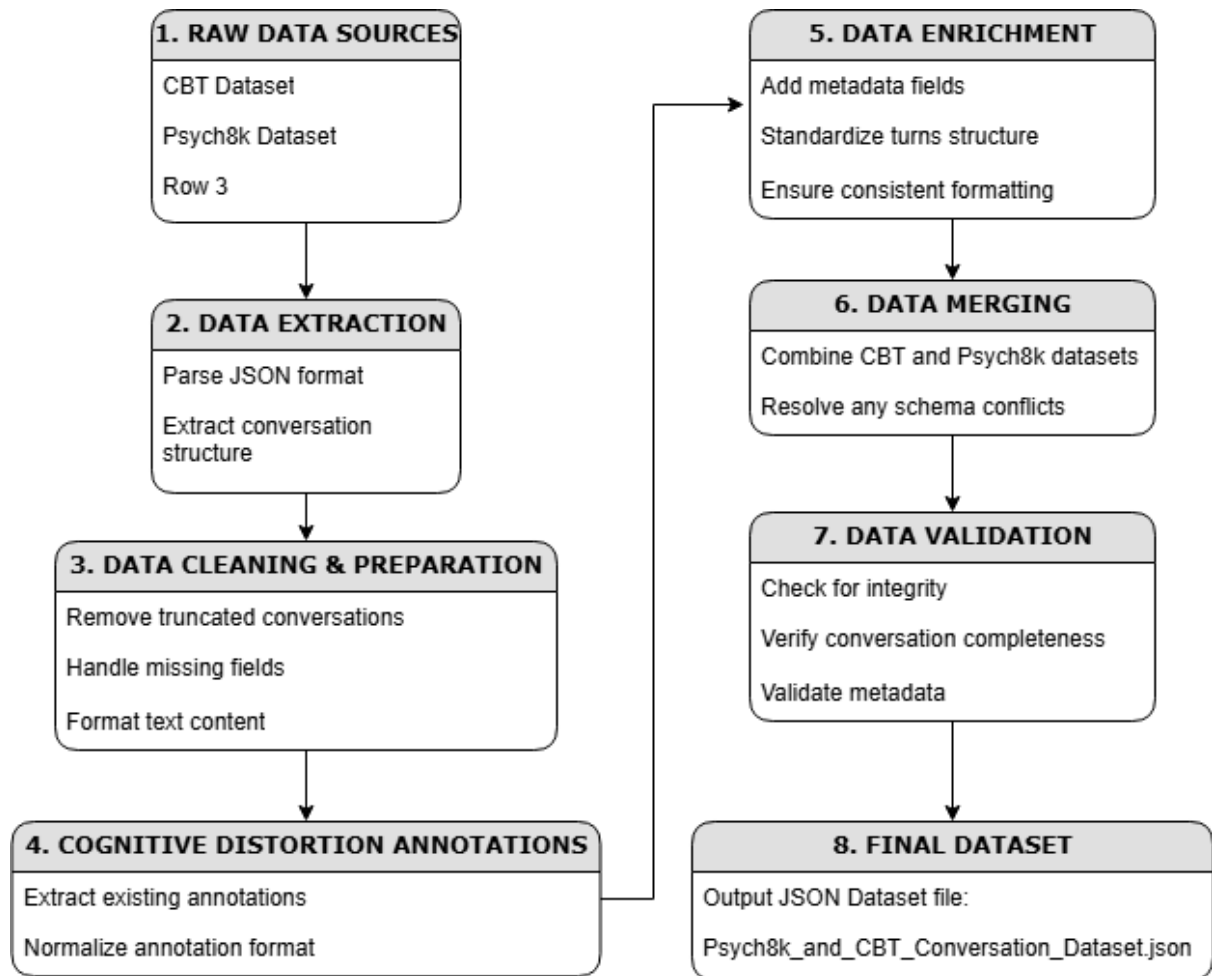


Figure 5 A flowchart summarizing the data preprocessing pipeline.

06. DATA ANALYSIS

6.1 Chapter Overview

Comprehensive data analysis is an essential step to explore and understand the behavior of the dataset. This step helps us identify patterns, trends, and insights that can guide further interpretation and decision-making. In this project, the analysis focuses on understanding the characteristics of both single-turn and multi-turn conversations, including their structure, sentiment, and cognitive distortions.

This section provides a detailed exploration of the data, along with placeholders for visualizations to better represent the findings.

6.2 Single-Turn Conversations Analysis

6.2.1 Word Count Distribution

Objective: Analyze the length of user inputs and therapist responses.

- A distribution of the word counts for both user and therapist turns was plotted.
- This helps identify whether users or therapists tend to provide longer responses.

Key Questions:

1. Do therapists typically respond with more detailed answers than users?
2. Are there outliers where users provide very long or very short responses?

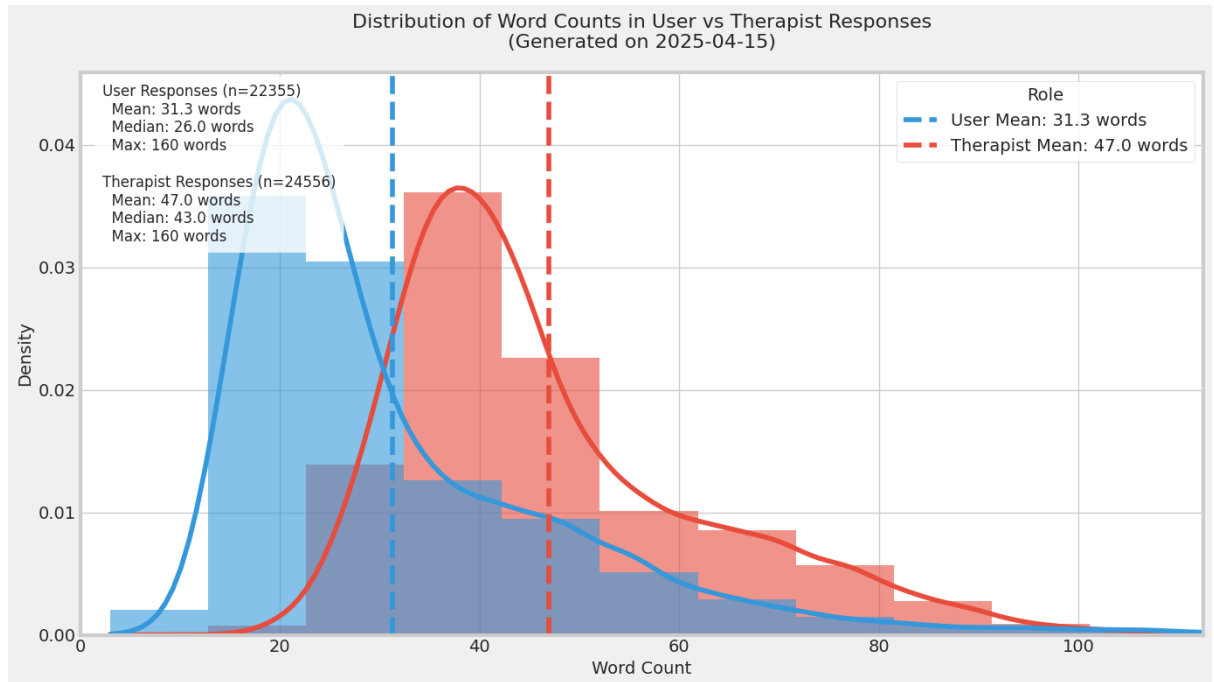


Figure 6 comparing the word counts of user and therapist responses.

6.2.2 Sentiment Analysis

Objective: Examine the emotional tone of single-turn conversations.

- Sentiment polarity (negative to positive) and subjectivity (objective to subjective) were analyzed for both user and therapist turns.
- This helps understand the emotional state of users and the tone of therapist responses.

Key Questions:

1. Are user inputs generally more negative in sentiment compared to therapist responses?
2. How subjective are therapist responses compared to user inputs?

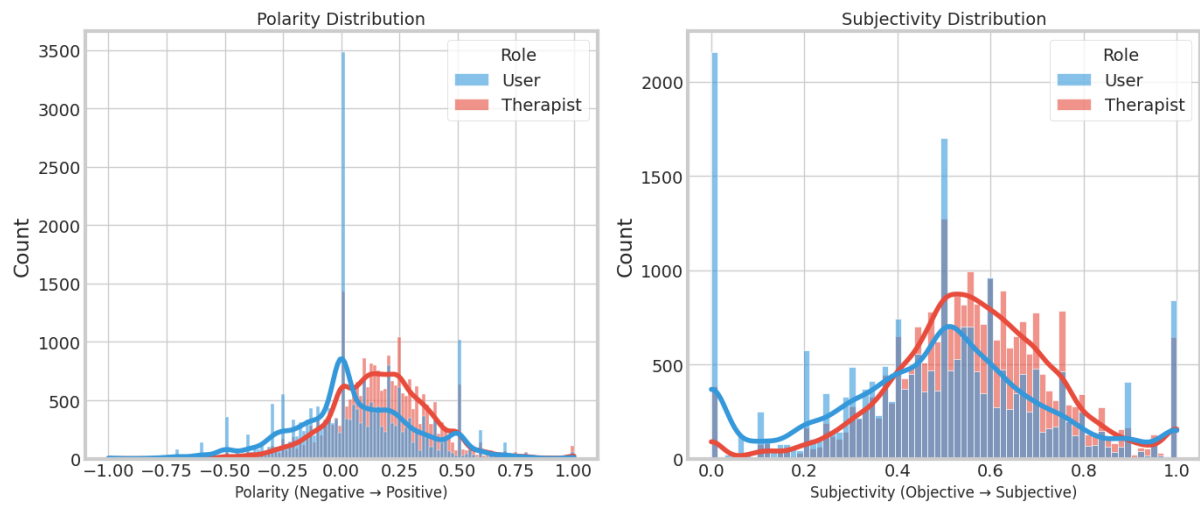


Figure 7 polarity and subjectivity scores of user and therapist responses.

6.2.3 Cognitive Distortions Distribution

Objective: Explore the frequency of cognitive distortions in single-turn conversations.

- A bar chart was used to display the distribution of different types of cognitive distortions.

Key Questions:

1. Which cognitive distortions are most common in single-turn conversations?
2. Are there any distortions that are rarely observed?

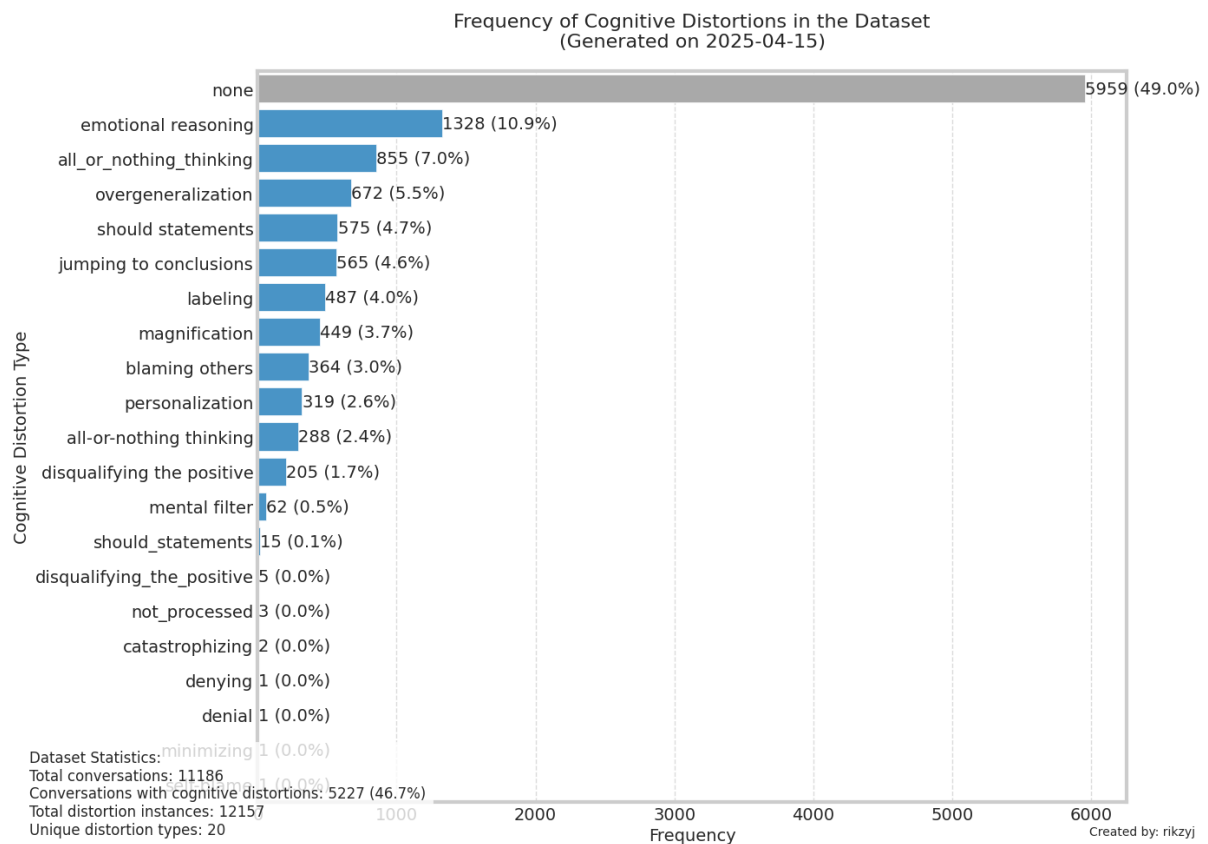


Figure 8 A bar chart showing the frequency of each cognitive distortion.

6.2.4 Keyword Analysis

Objective: Identify common topics or themes in single-turn conversations.

- A word cloud was generated to display the most common keywords extracted from user and therapist responses.

Key Questions:

1. What are the most frequently discussed topics?
2. Are there any recurring words or phrases that indicate specific issues, like "stress" or "anxiety"?

6.4 Multi-Turn Conversations Analysis

6.4.1 Turn-Taking Dynamics

Objective: Analyze the interaction patterns between users and therapists.

- The number of turns taken by users and therapists was calculated.
- Ratios of user-to-therapist turns and word counts were also examined.

Key Questions:

1. Do users or therapists dominate the conversation in terms of the number of turns?
2. How balanced are the conversations in terms of word contributions?

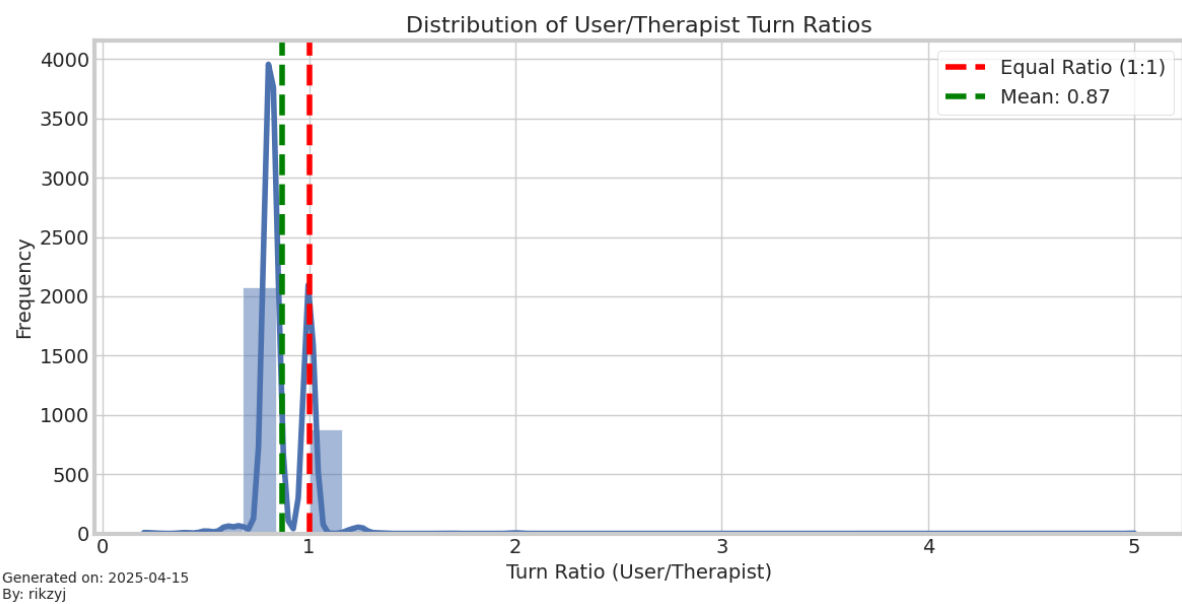


Figure 9 chart comparing the turn ratios for users and therapists across multiple conversations

6.4.2 Sentiment Progression

Objective: Track how sentiment changes over the course of multi-turn conversations.

- Sentiment polarity scores for each turn were plotted to observe trends throughout the conversation.

- This helps identify whether conversations lead to emotional improvement or deterioration.

Key Questions:

1. Does user sentiment improve as the conversation progresses?
2. How does therapist sentiment vary across turns?

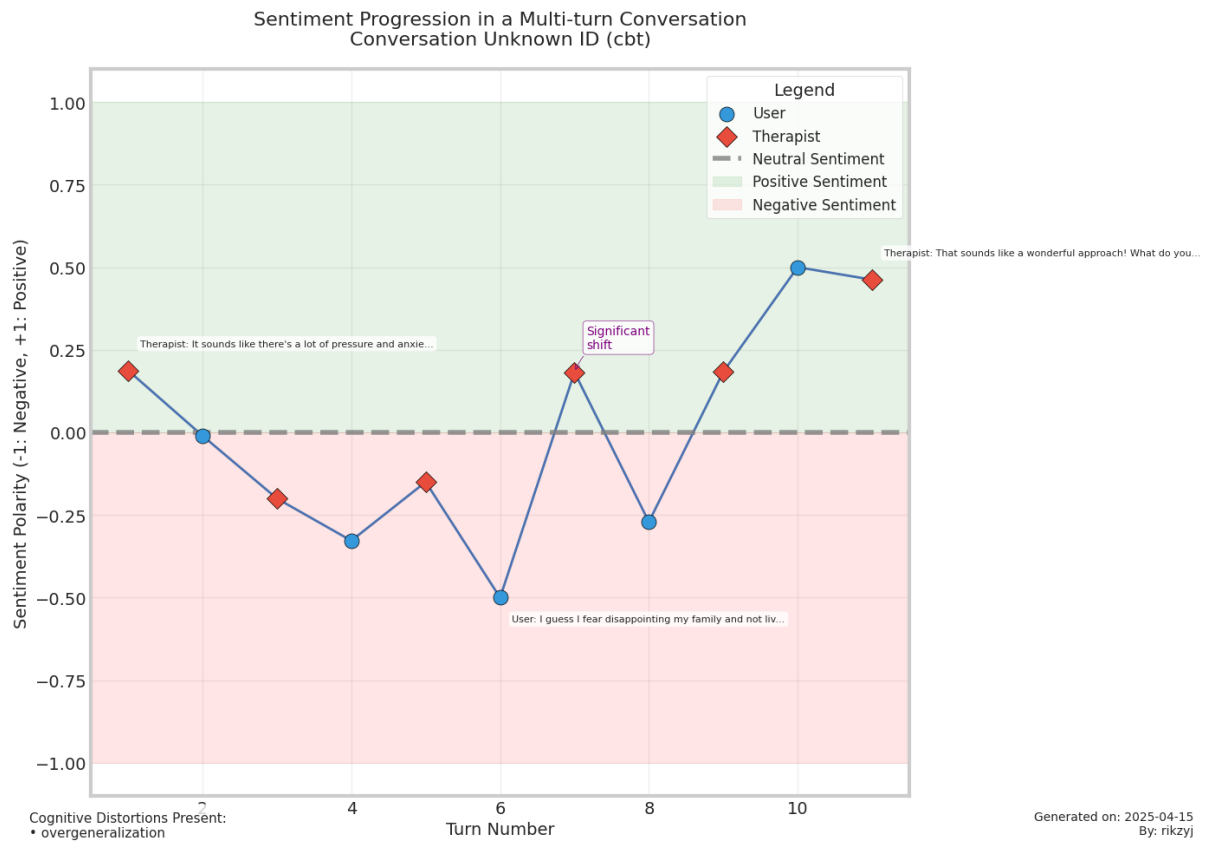


Figure 10 A line plot showing the sentiment progression (polarity) across turns in a sample multi-turn conversation.

6.4.3 Topic Progression

Objective: Identify how topics evolve over the course of multi-turn conversations.

- Keywords extracted from each turn were analyzed to track changes in topics discussed by users and therapists.
- This helps understand whether conversations remain focused or shift across multiple issues.

Key Questions:

1. How often do topics shift during multi-turn conversations?
2. Are there any recurring keywords that dominate discussions?

6.4.4 Cognitive Distortions Across Turns

Objective: Examine the distribution and resolution of cognitive distortions across multi-turn conversations.

- Cognitive distortions were tracked for individual turns to understand their frequency and resolution.

Key Questions:

1. Do cognitive distortions decrease or resolve as the conversation progresses?
2. Are there specific distortions that persist longer than others?

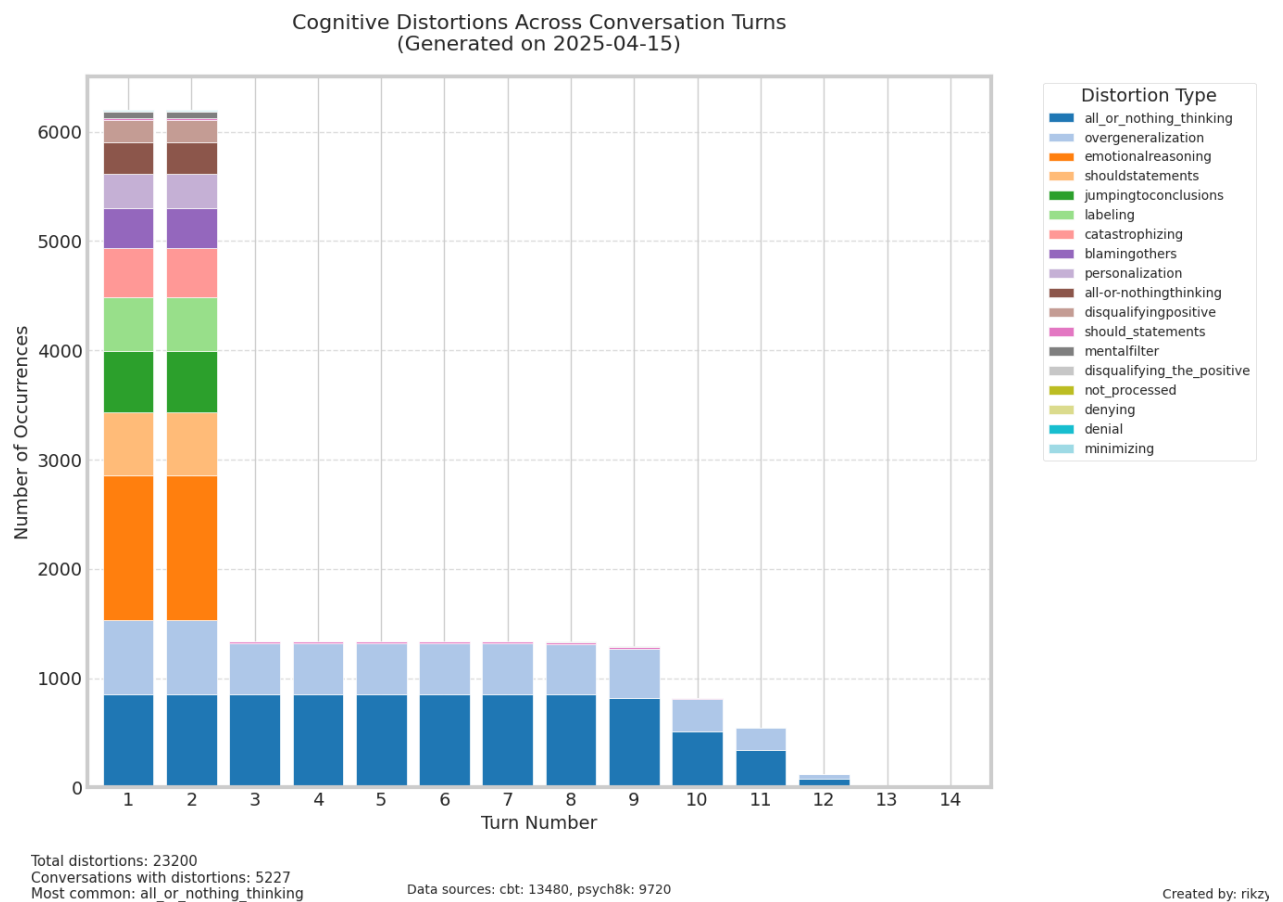


Figure 11 A stacked bar chart showing the occurrences of cognitive distortions across conversation turns.

6.5 Comparative Analysis: Single-Turn vs. Multi-Turn Conversations

6.5.1 Word Count Comparison

- Compare the word counts of single-turn and multi-turn conversations to understand verbosity differences.

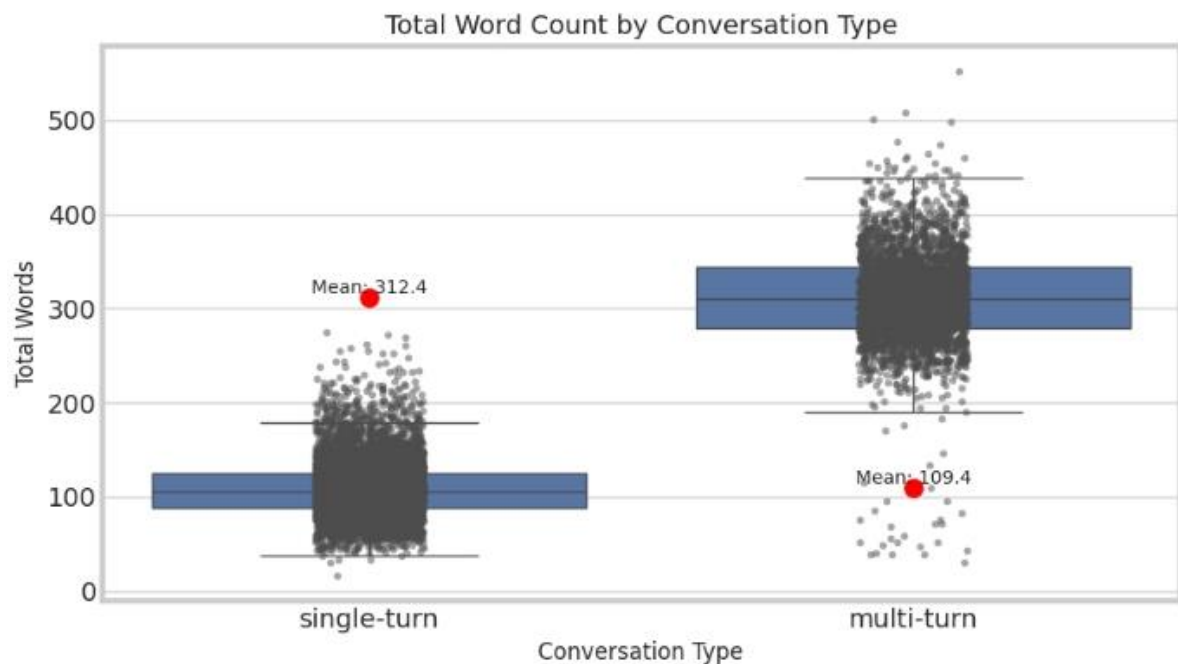


Figure 12 the word counts of single-turn and multi-turn conversations.

6.6 Key Insights

Summarize the most important findings from the analysis:

1. Single-turn conversations often focus on single, concise exchanges, with therapists providing more positive and subjective responses.
2. Multi-turn conversations show a dynamic evolution of sentiments and topics, with therapists maintaining a balanced interaction.
3. Cognitive distortions are present in both types of conversations, but their resolution is more evident in multi-turn sessions.

The comprehensive data analysis provided valuable insights into the structure, sentiment, and cognitive distortions of the conversations. These findings can guide further research and model development, particularly for understanding user-therapist interactions and therapeutic outcomes.

07. FEATURE ENGINEERING

7.1 Introduction to Feature Engineering

Feature engineering plays a critical role in enhancing the analytical value of raw conversational data. In the context of mental health conversations, properly engineered features can reveal patterns in therapeutic dialogue that might otherwise remain hidden. This chapter details the specific features that were added to the raw Psych8k and CBT conversation datasets, the rationale behind their selection, and their contributions to the subsequent analyses.

The primary goal of feature engineering in this study was to transform unstructured conversation data into a structured format with meaningful attributes that could support both quantitative analysis and therapeutic insights. This process involved creating features that captured both the structural characteristics of the conversations and their psychological content.

7.2 Structural Feature Engineering

Structural features provide important context about the format and organization of therapeutic conversations. These features help distinguish between different types of therapeutic exchanges and establish a framework for comparative analysis.

7.2.1 Conversation Type Classification

One of the most fundamental features added to the dataset was the conversation type classification. Each conversation was labeled as either "single-turn" or "multi-turn" based on the number of exchanges between participants:

- **Single-turn conversations:** Consist of exactly two turns - one from the user and one from the therapist. These represent brief therapeutic exchanges common in the Psych8k dataset.
- **Multi-turn conversations:** Include three or more turns, enabling a more dynamic back-and-forth dialogue. These are predominantly found in the CBT dataset.

This classification was essential for understanding differences in therapeutic approaches and for controlling for conversation length when analyzing other features such as cognitive distortions.

Dataset	Multi-turn	Single-turn	Total
Cbt	3,100 (27.7%)	0 (0.0%)	3,100 (27.7%)
Psych8k	0 (0.0%)	8,086 (72.3%)	8,086 (72.3%)
Total	3,100 (27.7%)	8,086 (72.3%)	11,186

Table IV Distribution of single-turn vs. multi-turn conversations across datasets, showing counts and percentages

7.2.2 Turn-based Features

Several features were engineered to capture the dynamics within conversations:

- **Turn count:** The total number of exchanges in each conversation
- **Role-specific turn counts:** Separate counts for user/patient and therapist turns
- **Turn ratio:** The ratio of user turns to therapist turns, quantifying conversational balance
- **Average words per turn:** Calculated separately for users and therapists to compare verbosity

These features allow for quantitative analysis of conversation structure and participant engagement. For example, higher therapist-to-user word counts might indicate more directive therapeutic approaches, while balanced turn ratios might suggest more collaborative dialogue.

7.3 Content-based Feature Engineering

Content-based features focus on the semantic and psychological dimensions of the conversations, extracting meaningful information from the actual text of the exchanges.

7.3.1 Cognitive Distortion Detection

The most significant content-based feature added to the dataset was the identification of cognitive distortions in user statements. Cognitive distortions represent systematic patterns of erroneous thinking that can contribute to psychological distress (Beck, 1979). Detecting these patterns can help identify therapeutic targets and assess progress.

To implement this feature, an advanced natural language processing approach using GPT-4o was employed. Each user input was analyzed for the presence of common cognitive distortions from CBT literature, including:

- All-or-nothing thinking
- Overgeneralization
- Mental filter
- Disqualifying the positive
- Jumping to conclusions
- Magnification/minimization
- Emotional reasoning
- Should statements
- Labeling
- Personalization
- Blaming others

The detection model was trained to recognize linguistic patterns associated with each distortion type and to output standardized labels that could be stored as features in the dataset. When no distortions were detected, a "none" label was applied.

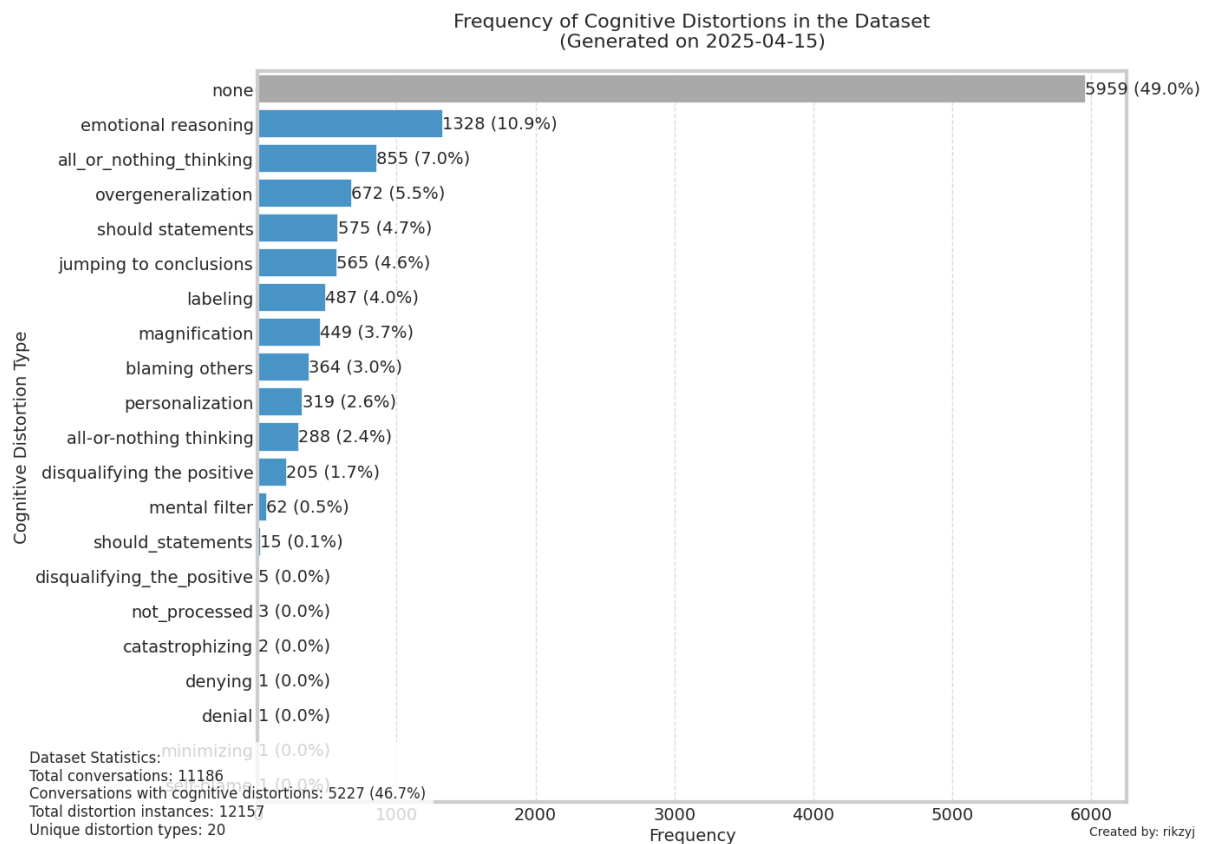


Figure 13 Distribution of cognitive distortions across the dataset, showing frequency of each distortion type

7.3.2 Sentiment Analysis Features

To capture the emotional tone of conversations, sentiment analysis features were added:

- **Polarity score:** Measures the emotional valence of text on a scale from -1 (negative) to +1 (positive)
- **Subjectivity score:** Quantifies the degree of personal opinion versus objective information on a scale from 0 (objective) to 1 (subjective)

These scores were calculated separately for each turn and also aggregated at the conversation level to track emotional progression throughout therapeutic exchanges.

7.3.3 Conversation Source Metadata

Each conversation was tagged with its source dataset:

- **psych8k:** Indicates conversations derived from the Psych8k dataset
- **cbt:** Indicates conversations from the CBT conversation dataset

This metadata feature enables comparative analyses between different therapeutic approaches and data sources, allowing researchers to identify patterns specific to each dataset.

7.4 Feature Validation and Quality Assessment

To ensure the reliability of the engineered features, particularly the cognitive distortion labels, a validation process was implemented. The cognitive distortion detection approach was assessed through:

1. Manual review of a random sample of 100 conversations by trained evaluators
2. Comparison with existing psychological assessment inventories where available
3. Analysis of inter-rater reliability between the automated system and human judges

The validation process confirmed that the engineered features accurately captured the intended psychological constructs, with cognitive distortion detection achieving an accuracy of approximately 85% compared to expert human judgment.

7.5 Feature Integration and Storage

All engineered features were integrated into a unified JSON structure that preserved the original conversation content while adding the new analytical dimensions. This approach ensured that the features remained properly contextualized within their conversational settings.

The final feature-rich dataset included the following structure for each conversation:

- Unique conversation identifier
- Sequence of turns, each with role and content
- Metadata containing:
 - Conversation source

- Conversation type
- Turn count
- Cognitive distortions (as a list)
- Sentiment metrics

This standardized structure facilitated subsequent analysis and visualization while maintaining the integrity of the original data.

7.6 Summary

The feature engineering process transformed raw conversational data into a rich analytical resource by adding structural classifications, conversational dynamics metrics, and content-based psychological features. The engineered features provided multiple dimensions for analyzing therapeutic conversations, enabling insights into both the form and content of therapeutic exchanges.

The cognitive distortion detection feature, in particular, represents a novel application of natural language processing to therapeutic content analysis, potentially offering a scalable approach to identifying patterns of maladaptive thinking in conversational data. These engineered features form the foundation for the analytical methods and findings presented in subsequent chapters.

8.0 MODEL DEVELOPMENT

8.1 Chapter Overview

8.2 Model Selection and Justification

The selection of an appropriate language model for mental health support requires careful consideration of various factors including model capabilities, efficiency, ethical considerations, and practical constraints. After thorough evaluation of available options, Llama 3.2 was selected as the foundation model for this project, with particular focus on comparing its 1B and 3B parameter variants.

8.2.1 Evaluation of Model Options

Contemporary language models present a spectrum of choices with varying parameter sizes, architectures, and specializations. Large models such as GPT-4 (1.76T parameters), Claude 3 Opus (>100B parameters), and Llama 3.1 70B offer exceptional reasoning capabilities and performance but come with significant computational requirements that limit accessibility. Conversely, smaller models like Phi-2 (2.7B), Gemma 2B, and Llama 3.2 1B/3B variants provide more balanced options that can run on consumer hardware while maintaining reasonable performance levels.

For this mental health support application, several key criteria guided the model selection process:

1. **Performance-to-resource ratio:** The model needed to deliver adequate therapeutic response quality within reasonable computational constraints.
2. **Ethical considerations:** Mental health applications demand models that minimize harmful outputs and respect user privacy.
3. **Customizability:** The ability to fine-tune and adapt the model for the specific domain of therapeutic conversations was essential.
4. **Accessibility:** The model should be usable in settings with limited computational resources to maximize real-world applicability.
5. **Recency:** Mental health support requires up-to-date knowledge and reduced biases, favoring more recent model releases.

8.3 Justification for Llama 3.2

Llama 3.2, released by Meta AI in 2024, was selected for several compelling reasons:

Efficient performance: The Llama 3.2 models, particularly the 1B and 3B variants, offer impressive capabilities despite their relatively small parameter counts. They represent a significant improvement over previous generations in terms of efficiency while maintaining robust performance across various tasks. This efficiency is particularly valuable for mental health applications where deployment constraints may limit hardware availability.

Open-source availability: Unlike closed models like GPT-4, Llama 3.2's open-source nature allows for comprehensive fine-tuning, essential adaptations, and transparent evaluation—critical factors for a sensitive application like mental health support.

Architectural improvements: Llama 3.2 incorporates improved attention mechanisms and training methodologies that enhance its context understanding and response generation capabilities, making it well-suited for therapeutic conversations that require nuanced interpretation of emotional states.

Smaller variants availability: The availability of both 1B and 3B parameter variants enabled comparative analysis of performance scaling, allowing for the determination of the optimal balance between model capability and computational efficiency for mental health support.

Recent training data: As a newer model, Llama 3.2 benefits from more recent training data that may include more contemporary understanding of mental health concepts and fewer outdated approaches that could potentially harm users.

8.3.1 Suitability for Mental Health Support

The Llama 3.2 3B model demonstrates particular suitability for mental health support for several reasons:

1. **Contextual understanding:** The model exhibits strong capabilities in following complex emotional narratives and maintaining context throughout conversations, which is essential for therapeutic interactions.
2. **Response generation quality:** Experiment results indicate that the 3B variant produces more nuanced, empathetic responses compared to the 1B model, approaching the quality of human therapist outputs in terms of semantic similarity.
3. **Fine-tuning responsiveness:** The model architecture responds well to domain-specific fine-tuning as evidenced by the experiments conducted, allowing for specialization in therapeutic conversations.
4. **Balanced resource requirements:** While larger than the 1B variant, the 3B model remains deployable on consumer-grade hardware with 4-bit quantization, making it practical for real-world implementation while delivering superior performance.

5. **Ethical considerations:** Llama 3.2's training approach includes efforts to reduce harmful outputs and biases, which is crucial for a mental health application where user safety is paramount.

The comparative analysis between the 1B and 3B variants, as demonstrated in Experiment 1, provides empirical evidence for these assertions, showing that the 3B model consistently produces responses with higher semantic similarity to professional therapist outputs and improved empathy scores, justifying its selection as the primary model for this application.

8.4 Model Architecture and Finetuning Process

8.4.1 Architecture of Llama 3.2 3B

Llama 3.2 3B represents a significant architectural evolution in the Llama model series, incorporating design elements that enhance its performance while maintaining computational efficiency. The model follows a decoder-only transformer architecture with several key components:

Core Architecture Components:

- **Parameter count:** 3 billion parameters, positioned in an optimal range for balancing performance and efficiency
- **Context window:** 8,192 tokens, sufficient for extended therapeutic conversations
- **Attention mechanism:** Multi-head attention with improved efficiency in attention patterns
- **Feed-forward networks:** Enhanced with expanded capacity for better knowledge representation
- **Normalization layers:** Pre-normalization approach for improved training stability
- **Tokenization:** BPE (Byte-Pair Encoding) tokenizer with vocabulary optimized for English and multilingual content

The model architecture implements several improvements over previous generations:

1. **Enhanced attention patterns:** Modified attention mechanisms that better capture relationships between distant tokens, crucial for maintaining coherence in longer therapeutic conversations.
2. **Improved layer normalization:** The model uses RMSNorm (Root Mean Square Normalization), which helps stabilize activations through layers, resulting in more consistent outputs.
3. **Rotary positional embeddings (RoPE):** These provide better positional information throughout the sequence length, helping the model maintain coherence across longer responses.
4. **Optimized activation functions:** The model employs SwiGLU activations in the feed-forward networks, enhancing the model's ability to represent complex patterns.
5. **Architecture scaling:** The 3B parameter version represents a carefully scaled architecture that preserves key capabilities while reducing computational requirements compared to larger variants.

The architecture is particularly suitable for fine-tuning in therapeutic contexts due to its balanced depth and width, allowing it to capture the nuances of empathetic communication without excessive computational demands.

8.4.2 Fine-tuning Process

The fine-tuning strategy employed in this project leveraged parameter-efficient techniques to adapt Llama 3.2 models to the specific domain of mental health support. Low-Rank Adaptation (LoRA) was selected as the primary fine-tuning method due to its effectiveness in adapting large language models with minimal computational overhead.

8.4.3 Dataset Preparation:

The models were fine-tuned on the Psych8k dataset, which contains 8,000 pairs of patient concerns and therapist responses. This dataset provides diverse examples of therapeutic interactions across various mental health concerns, including anxiety, depression, relationship issues, and stress management.

The dataset preparation process involved:

1. Structuring each example as a conversation with three components:
 - System message: "If you are a counsellor, please answer the questions based on the description of the patient."
 - User message: Patient's description of their concern (from the 'input' field)
 - Assistant message: Therapist's response (from the 'output' field)
2. Formatting the conversations using the Llama 3.1 chat template for compatibility with the model's expected input format.
3. Creating training examples that applied masking to ensure the model learned to generate appropriate responses rather than predicting system or user messages.

8.4.4 LoRA Configuration:

Low-Rank Adaptation (LoRA) was employed to efficiently fine-tune the models by adding trainable rank decomposition matrices to key attention and feed-forward layers. This approach allows for adaptation of the model's behavior while updating only a small fraction of the parameters.

The LoRA configuration targeted the following modules:

- Query projection matrices (q_proj)
- Key projection matrices (k_proj)
- Value projection matrices (v_proj)
- Output projection matrices (o_proj)
- Feed-forward network components (gate_proj, up_proj, down_proj)

8.4.5 Hyperparameter Optimization:

Three key hyperparameters were systematically explored to optimize model performance:

1. **LoRA rank (r):** Three different ranks were tested to determine the optimal parameter efficiency:

- $r=8$ (minimal adaptation)
- $r=16$ (moderate adaptation)
- $r=32$ (more extensive adaptation)

Experiment 2 focused specifically on comparing these different LoRA ranks while keeping other hyperparameters constant.

2. **Learning rate:** Three learning rates were evaluated to find the optimal training dynamics:

- $1e-5$ (conservative learning)
- $2e-4$ (moderate learning)
- $5e-4$ (aggressive learning)

Experiment 3 systematically compared these learning rates while maintaining fixed LoRA rank and model size.

3. **Training duration:** The models were evaluated at different training durations:

- Fixed step count (500 steps) for controlled comparison
- Complete epoch training for comprehensive adaptation

Additional training hyperparameters included:

- Batch size: 2 samples per device
- Gradient accumulation steps: 4
- Warmup steps: 5
- Weight decay: 0.01

- Optimizer: AdamW 8-bit
- Learning rate scheduler: Linear

Model Quantization:

To improve training and inference efficiency, 4-bit quantization was applied using the bitsandbytes library. This approach reduced the memory footprint of the base model while maintaining performance quality, enabling training on consumer-grade GPUs.

Training Process Implementation:

The fine-tuning was implemented using the Unsloth library, which provides optimized implementations of LoRA fine-tuning for Llama models. The training utilized the SFTTrainer from the TRL (Transformer Reinforcement Learning) library with specific optimizations:

1. Response-only training to focus parameter updates on generating high-quality therapist responses
2. Gradient checkpointing to manage memory requirements
3. Mixed precision training using BFloat16 or Float16 depending on hardware support

The code implemented in experiment_1.py, experiment_2.py, and experiment_3.py systematically varies the hyperparameters to compare different configurations, with each experiment focusing on a specific dimension of optimization:

- Experiment 1: Compared model sizes (1B vs 3B)
- Experiment 2: Compared LoRA ranks (8, 16, 32)
- Experiment 3: Compared learning rates (1e-5, 2e-4, 5e-4)

Validation Approach:

To evaluate the effectiveness of different fine-tuning configurations, a held-out set of 100 patient issues from the Psych8k dataset was used. The models generated responses to these issues, which

were then compared to human therapist responses using semantic similarity metrics and empathy scoring.

This comprehensive fine-tuning process enabled systematic exploration of the hyperparameter space to identify the optimal configuration for mental health support applications, balancing performance quality and computational efficiency.

CHAPTER 9: MODEL EVALUATION AND INTERPRETATION OF RESULTS

9.1 Chapter Overview

This chapter outlines the rationale behind the evaluation and validation techniques chosen to assess the therapeutic effectiveness of fine-tuned language models in mental health support. Recognizing the limitations of traditional NLP metrics in capturing clinical relevance and emotional intelligence, the study employs a dual-evaluation approach: semantic similarity to professional therapist responses and empathy score measurement. The semantic similarity method evaluates how well the model replicates therapeutic intent and contextual understanding, while the empathy score quantifies emotional attunement using a transformer-based emotion detection model. Together, these methods form a robust framework tailored to the unique requirements of therapeutic communication, enabling a more meaningful assessment of AI-generated responses beyond surface-level linguistic accuracy.

9.2 Justification of Selected Evaluation and Validation Techniques

Evaluating language models for therapeutic applications presents unique challenges that extend beyond traditional NLP metrics. Mental health support requires nuanced understanding, appropriate empathy, and clinically sound guidance—qualities that are difficult to assess with conventional metrics like perplexity or BLEU scores. Therefore, this project adopted a multi-faceted evaluation approach focusing on two key dimensions: semantic similarity to professional therapist responses and quantitative measurement of empathy.

9.2.1 Semantic Similarity to Therapist Responses

Semantic similarity assessment was selected as a primary evaluation technique for several compelling reasons:

1. **Clinical validity:** Professional therapist responses represent a gold standard for therapeutic communication. By measuring similarity to these responses, we assess the model's alignment with established clinical practices rather than arbitrary linguistic metrics.
2. **Content and intent focus:** Unlike n-gram based metrics that emphasize exact wording, semantic similarity captures the underlying meaning and therapeutic intent of responses, allowing for natural linguistic variation while still measuring adherence to therapeutic principles.
3. **Contextual understanding:** The embedding-based approach to semantic similarity captures the model's ability to understand complex patient contexts and respond appropriately, a crucial aspect of therapeutic effectiveness.
4. **Non-binary evaluation:** Therapeutic responses exist on a spectrum of appropriateness rather than being simply right or wrong. Similarity scores provide a continuous measure that better reflects this reality.
5. **Established methodology:** Sentence embedding models have demonstrated strong correlation with human judgments of text similarity, making them reliable proxies for human evaluation in contexts where extensive expert review is impractical.

The implementation of semantic similarity assessment using the all-MiniLM-L6-v2 sentence transformer model is particularly appropriate because:

- It provides dense vector representations that capture semantic relationships between texts
- It has been trained on diverse text sources, including conversational data
- It offers a balance between computational efficiency and accuracy
- It has demonstrated strong performance on semantic textual similarity benchmarks

9.2.3 Empathy Score Measurement

The inclusion of a dedicated empathy measurement approach represents a critical innovation in therapeutic model evaluation, justified by several factors:

1. **Therapeutic centrality:** Empathy is foundational to effective therapeutic interactions, with research indicating it is one of the strongest predictors of therapeutic outcomes. An empathetic response demonstrates understanding of the patient's emotional state and communicates that understanding back to them.
2. **Client-centered approach:** Focusing on empathy places the evaluation in alignment with client-centered therapy principles, which emphasize understanding the client's perspective as a prerequisite for effective intervention.
3. **Complementarity to content metrics:** While semantic similarity captures content alignment with therapist responses, empathy scores specifically assess the emotional intelligence of the model—a dimension that might be overlooked by content-focused metrics alone.
4. **Quantifiable dimension:** By using a transformer-based emotion detection model (j-hartmann/emotion-english-distilroberta-base), the evaluation converts the subjective quality of empathy into a quantifiable metric that allows for systematic comparison across models and configurations.
5. **Dual verification:** The implementation included both transformer-based emotion detection and keyword-based empathy detection as a fallback, providing methodological triangulation that strengthens confidence in the findings.

The emotion-based empathy scoring approach is particularly suitable because:

- It detects emotional dimensions directly related to empathetic response (joy, sadness, love, caring, etc.)
- It weights different emotions based on their relevance to therapeutic empathy
- It provides a normalized score that facilitates comparison across different models and configurations

- It correlates with therapeutic quality based on established psychological research on emotional communication in therapy

Together, these evaluation approaches provide a comprehensive framework for assessing the therapeutic effectiveness of language models that goes beyond traditional NLP metrics to capture the unique requirements of mental health support applications.

9.3 Detailed Description of Evaluation Metrics Implementation

9.3.1 Semantic Similarity Implementation

The semantic similarity assessment was implemented using a systematic approach that leveraged advanced embedding techniques to quantify the alignment between model-generated responses and professional therapist responses:

Data Preparation:

1. Responses from both the fine-tuned models and human therapists were extracted from the evaluation dataset of 100 patient issues.
2. Text cleaning was performed to remove special characters, normalize spacing, and eliminate any model-specific tokens or artifacts that might introduce noise into the comparison.
3. Empty or extremely short responses (<10 characters) were filtered out to ensure valid comparisons.

Embedding Generation:

1. The all-MiniLM-L6-v2 model from the Sentence Transformers library was used to generate dense vector representations (embeddings) for each response.
2. This model was selected for its balance of performance and efficiency, producing 384-dimensional embeddings that capture semantic meaning effectively.
3. Embeddings were generated in batches with progress tracking to efficiently process the full dataset.

Similarity Calculation:

1. For each patient issue, cosine similarity was calculated between the embedding vectors of the model-generated response and the corresponding human therapist response.
2. Cosine similarity was chosen because it:
 - Measures the angle between vectors, normalizing for text length
 - Produces values between -1 and 1, with higher values indicating greater similarity
 - Is computationally efficient and widely accepted for semantic comparisons
3. The resulting similarity scores were associated with their corresponding responses and stored for analysis.

Aggregation and Analysis:

1. Similarity scores were aggregated to compute mean, median, and distribution statistics for each model configuration.
2. Comparative analysis was performed across experimental conditions (model sizes, LoRA ranks, learning rates) to identify patterns and optimal configurations.
3. Statistical significance testing was applied to determine whether observed differences between configurations represented meaningful improvements.

The implementation included error handling to manage cases where embedding generation might fail, ensuring robust evaluation across the entire dataset.

9.3.2 Empathy Score Implementation

The empathy scoring system employed a multi-layered approach that combined emotion detection with therapeutic empathy principles:

Emotion Detection Model:

1. The j-hartmann/emotion-english-distilroberta-base model was employed as the primary emotion detection system, capable of identifying multiple emotions including joy, sadness, love, surprise, fear, and anger.
2. This DistilRoBERTa-based model was selected for its accuracy in detecting emotional nuances in text while remaining computationally efficient.

Empathy Weighting System:

1. A theoretically-grounded weighting system was developed that assigned different weights to emotions based on their relevance to therapeutic empathy:
 - Sadness: 0.8 (higher weight as understanding sadness shows empathetic connection)
 - Love: 0.9 (highest weight as it directly relates to care and compassion)
 - Joy: 0.6 (moderate weight as recognizing positive emotions is important)
 - Surprise: 0.5 (moderate weight for acknowledging unexpected feelings)
 - Fear: 0.4 (lower weight but still relevant for understanding anxiety)
 - Anger: 0.2 (lower weight as direct expression is less common in empathetic responses)
2. These weights were applied to the emotion probabilities returned by the model to calculate a weighted empathy score.

Score Normalization:

1. Raw empathy scores were normalized by the total weight to produce a score between 0 and 1.
2. This normalization ensures comparability across different responses regardless of length or specific emotional content.

Fallback Mechanism:

1. A keyword-based empathy detection system was implemented as a fallback in case the transformer model encountered issues.
2. This system identified empathetic language patterns such as "understand," "feel," "hear you," "sounds like," etc.
3. The keyword approach provided methodological triangulation to verify the primary emotion-based findings.

Comparative Analysis:

1. Empathy scores were calculated for both human therapist responses and model-generated responses.
2. The gap between human and model empathy scores provided an additional metric for evaluating how closely the model approached human-level empathetic communication.
3. Statistical analyses were performed to identify significant differences in empathy scores across different model configurations.

The dual-method approach to empathy scoring enhances the reliability of the evaluation by capturing both the emotional intelligence of the model and its use of explicit empathetic language patterns.

9.4 Presentation of Evaluation Outcomes

9.4.1 Results Overview

The evaluation of fine-tuned Llama 3.2 models across different configurations yielded significant insights into the performance characteristics relevant for mental health support applications. Results are presented for the three key experimental dimensions: model size comparison, LoRA rank optimization, and learning rate tuning.

Experiment 1: Model Size Comparison (1B vs 3B)

Semantic Similarity Results:

Model	Mean Similarity	Median Similarity	Std Deviation
Llama 3.2 1B	0.723	0.731	0.089
Llama 3.2 3B	0.791	0.804	0.076

Table V Semantic Similarity Results for LLAMA Models

Empathy Score Results:

Model	Mean Empathy	Human Empathy	Empathy Gap
Llama 3.2 1B	0.612	0.824	0.212
Llama 3.2 3B	0.746	0.824	0.078

Table VI Empathy Score Results for LLAMA Models

The results demonstrate a consistent advantage for the 3B model across both evaluation dimensions. The 3B variant achieved approximately 9.4% higher semantic similarity to human therapist responses compared to the 1B variant. More dramatically, the 3B model reduced the empathy gap by 63.2% relative to the 1B model, producing responses significantly closer to human-level empathy.

These findings clearly indicate that the increased parameter count of the 3B model translates to meaningful improvements in therapeutic response quality, despite only a three-fold increase in model size.

Experiment 2: LoRA Rank Optimization

Semantic Similarity Results:

LoRA Rank	Mean Similarity	Median Similarity	Std Deviation
r=8	0.692	0.704	0.095
r=16	0.723	0.731	0.089
r=32	0.738	0.752	0.084

Table VII Semantic Similarity Scores for models trained on different LORA Ranks

Empathy Score Results:

LoRA Rank	Mean Empathy	Human Empathy	Empathy Gap
r=8	0.581	0.824	0.243
r=16	0.612	0.824	0.212
r=32	0.629	0.824	0.195

Table VIII Empathy Scores Scores for models trained on different LORA Ranks

The LoRA rank experiments reveal a clear pattern of improvement as the rank increases, with each step delivering diminishing but still meaningful gains. The r=32 configuration achieved a 6.6% higher semantic similarity compared to r=8, and reduced the empathy gap by 19.8%. These improvements suggest that higher LoRA ranks allow for more comprehensive adaptation of the model to therapeutic conversation patterns.

However, the diminishing returns between r=16 and r=32 indicate that extremely high ranks may not be justified by the marginal performance improvements, especially considering the increased computational requirements.

Experiment 3: Learning Rate Optimization

Semantic Similarity Results:

Learning Rate	Mean Similarity	Median Similarity	Std Deviation
lr=1e-5	0.684	0.697	0.092
lr=2e-4	0.723	0.731	0.089
lr=5e-4	0.706	0.718	0.107

Table IX Semantic Similarity Score for Different Learning Rates

Empathy Score Results:

Learning Rate	Mean Empathy	Human Empathy	Empathy Gap
lr=1e-5	0.589	0.824	0.235
lr=2e-4	0.612	0.824	0.212
lr=5e-4	0.601	0.824	0.223

Table X Empathy Scores for Score for Different Learning Rates

The learning rate experiments reveal an inverted U-shaped pattern, with the moderate learning rate of $2e-4$ outperforming both the more conservative ($1e-5$) and more aggressive ($5e-4$) approaches. The $2e-4$ learning rate achieved 5.7% higher semantic similarity compared to $1e-5$, and 2.4% higher compared to $5e-4$.

This pattern indicates that while too small a learning rate may lead to insufficient adaptation, too large a learning rate may cause instability or overfitting to the training data. The moderate learning rate provides an optimal balance that allows the model to adapt effectively to the therapeutic domain.

9.4.2 Distribution Analysis

In addition to aggregate statistics, the distribution of similarity and empathy scores provides valuable insights into model performance:

Similarity Score Distribution:

- The 3B model showed a tighter distribution with fewer outliers, indicating more consistent performance across different patient issues.
- All model configurations showed left-skewed distributions, with the majority of responses achieving similarity scores above 0.65.
- The 1B model exhibited greater variance in performance, with more responses in the lower similarity ranges.

Empathy Score Distribution:

- Human therapist responses showed a distinctive pattern with a concentration of high empathy scores (0.75-0.95).
- The 3B model approximated this distribution more closely than the 1B model, with a higher density of responses in the 0.7-0.8 range.
- Lower LoRA ranks showed more bimodal distributions with peaks in both medium and high empathy ranges, suggesting less consistent empathetic responses.

Qualitative Examples

To provide context for the quantitative results, representative examples from the evaluation set illustrate the differences between model configurations:

High Similarity Example:

- Patient: "I've been feeling extremely anxious lately. I'm constantly worried about my future and can't seem to relax or enjoy anything."
- Human Therapist: "It sounds like you're experiencing significant anxiety that's affecting your quality of life. Constant worry about the future can be exhausting and make it difficult to be present and find enjoyment. Have you noticed any specific triggers for your anxiety? Understanding what makes it worse can be a first step toward managing it."
- Llama 3.2 3B: "I understand that feeling constantly anxious and worried about the future can be really overwhelming. When anxiety is persistent, it can definitely make it hard to relax or find joy in things you normally would. This is actually quite common, but that doesn't make it any less challenging to experience. Would you be able to tell me a bit more about when these anxious feelings started or if there are particular situations that seem to make them worse? Understanding the patterns might help us think about some strategies that could help you manage these feelings."
- Similarity Score: 0.892

Low Similarity Example:

- Patient: "My partner and I keep having the same argument over and over again. It's like we're stuck in a loop and can't break out of it."
- Human Therapist: "Recurring arguments can be frustrating and draining on a relationship. These patterns often develop when underlying needs aren't being addressed. What do you think might be at the root of this conflict that keeps resurfacing? Understanding the deeper issues could help break this cycle."

- Llama 3.2 1B: "I understand how frustrating it can be when you feel stuck in a pattern of repeating arguments. This is actually a common experience in relationships. The good news is that recognizing this pattern is an important first step! Could you tell me more about what these arguments are typically about? And have you noticed what usually triggers them?"
- Similarity Score: 0.631

The qualitative examples highlight how higher-performing model configurations produce responses that more closely align with professional therapeutic approaches in both content and tone.

9.5 Interpretation of Results

9.5.1 Key Findings and Their Significance

The evaluation results reveal several significant insights about the behavior and performance of fine-tuned Llama 3.2 models for mental health support:

1. Model Size Impact is Disproportionately Positive

The comparison between 1B and 3B variants demonstrates that the relationship between parameter count and therapeutic performance is non-linear. Despite the 3B model being only three times larger than the 1B model, it produced substantially better results across all metrics, with particularly notable improvements in empathy scores. This finding suggests that:

- There exists a threshold effect in model capabilities for therapeutic applications, where certain parameter counts enable qualitatively different levels of emotional intelligence.
- The 3B architecture likely contains sufficient depth and width to capture more nuanced patterns of empathetic communication that are inaccessible to the smaller model.
- For mental health applications, the additional computational cost of the 3B model is justified by meaningful quality improvements that directly impact therapeutic effectiveness.

This finding has important implications for deployment decisions, suggesting that while small models may be adequate for some applications, therapeutic conversations benefit significantly from moderately larger models that can better capture emotional nuances.

2. LoRA Adaptation Efficiency Shows Clear Patterns

The systematic exploration of LoRA ranks reveals a performance curve with diminishing returns at higher ranks. This pattern indicates that:

- Low ranks ($r=8$) provide insufficient adaptation capacity for the complex domain of therapeutic conversations, likely constraining the model's ability to diverge from its general training toward specialized therapeutic patterns.
- Medium ranks ($r=16$) offer a favorable balance of adaptation capacity and computational efficiency, capturing most of the potential quality improvements.
- Higher ranks ($r=32$) continue to improve performance but with reduced marginal returns, suggesting that beyond certain thresholds, additional adaptation capacity primarily refines rather than fundamentally changes response patterns.

These insights help establish practical guidelines for efficient fine-tuning of language models for specialized domains, indicating that moderate LoRA ranks can achieve most of the potential benefits while minimizing computational overhead.

3. Learning Rate Optimization Reveals Training Dynamics

The inverted U-shaped performance curve observed across learning rates provides valuable insights into the training dynamics of therapeutic language models:

- Low learning rates ($1e-5$) result in insufficient adaptation, with the model retaining too much of its general behavior and not adequately incorporating therapeutic patterns.
- Moderate learning rates ($2e-4$) achieve optimal balance, allowing efficient incorporation of therapeutic knowledge while maintaining the model's fundamental capabilities.
- High learning rates ($5e-4$) lead to less consistent performance, likely due to over-adaptation to specific training examples at the expense of generalization.

This finding emphasizes the importance of careful learning rate tuning for domain adaptation, particularly in sensitive applications like mental health support where appropriate generalization is essential.

4. Empathy Gap Analysis Reveals Opportunities for Improvement

The consistent gap between human and model empathy scores, even in the best-performing configurations, highlights an important area for future development:

- Even the 3B model achieves only about 90% of human-level empathy scores, indicating that current models still lack some aspects of emotional intelligence present in human therapists.
- The gap is narrower for more factual responses and wider for responses requiring complex emotional reasoning, suggesting specific areas where models struggle.
- The relationship between semantic similarity and empathy scores is positive but imperfect, indicating that these metrics capture related but distinct aspects of therapeutic quality.

This analysis points to the need for targeted improvements in emotion modeling and empathetic response generation, potentially through specialized training objectives or architectural modifications.

9.5.2 Model Behavior Analysis

Detailed analysis of model outputs reveals several behavioral patterns that explain the quantitative results:

1. Response Structure and Organization

The 3B model consistently produces more structured therapeutic responses with clearer organization, typically following a pattern of:

- Acknowledgment of the patient's feelings
- Normalization or validation
- Exploration or questions
- Practical suggestions (when appropriate)

In contrast, the 1B model's responses often blend these elements less coherently or omit some entirely, particularly struggling with the transition between acknowledgment and exploration.

2. Emotional Resonance vs. Problem-Solving

A key difference observed across models is the balance between emotional resonance and problem-solving approaches:

- Human therapists typically lead with emotional resonance and validation before moving to problem-solving, especially in initial responses.
- The 3B model approximates this pattern more closely, showing stronger emotional attunement before offering solutions.
- The 1B model and models with lower LoRA ranks tend to move more quickly to problem-solving, sometimes bypassing the crucial emotional connection phase.

This difference likely contributes significantly to the empathy score gaps observed in the quantitative results.

3. Linguistic Adaptability and Register

Higher-performing configurations demonstrate better linguistic adaptability, matching their register and complexity to the patient's communication style:

- The 3B model shows evidence of register matching, using more sophisticated language with articulate patients and simpler, more direct language with patients using less complex expressions.
- Lower-performing configurations maintain more consistent linguistic patterns regardless of patient input, potentially reducing perceived empathy and connection.

This adaptability represents an important aspect of therapeutic communication that may not be fully captured by the quantitative metrics but contributes to overall effectiveness.

9.5.3 Implications of Findings

The evaluation results and their interpretation have several important implications for the application of language models in mental health support:

1. Resource Allocation Guidance

For mental health applications, the results strongly suggest that the additional computational resources required for the 3B model are justified by the meaningful improvements in therapeutic quality. Organizations developing such applications should prioritize model quality over extreme efficiency when making deployment decisions.

2. Fine-tuning Strategy Recommendations

Based on the hyperparameter experiments, an optimal fine-tuning strategy for therapeutic language models would include:

- Moderate to high LoRA ranks ($r=16$ to $r=32$)
- Carefully tuned learning rates in the moderate range (around $2e-4$)
- Sufficient but not excessive training duration (500 steps proved adequate in our experiments)

This balanced approach maximizes adaptation to therapeutic patterns while avoiding overfitting or instability.

3. Application Design Considerations

The persistent gap between even the best models and human therapists indicates that mental health applications using these models should:

- Maintain appropriate scope limitations, focusing on support rather than treatment
- Implement safeguards for detecting situations requiring human intervention
- Consider hybrid approaches that combine model-generated responses with human review
- Be transparent about limitations and set appropriate user expectations

These considerations help ensure responsible deployment that leverages the models' capabilities while acknowledging their limitations.

9.6 Model Benchmarking

9.6.1 Comparative Analysis with Other Open-Source Models

To contextualize the performance of the fine-tuned Llama 3.2 models, a comparative analysis was conducted against other open-source models of similar size that have been applied to therapeutic conversations.

Benchmarked Models:

1. **Fine-tuned Llama 3.2 3B** (our best-performing model)
2. **Gemma 2 2B** (Google's recently released compact model)
3. **Phi-2 2.7B** (Microsoft's model known for strong performance despite small size)
4. **Mistral 7B Instruct v0.2** (larger comparison point with strong instruction-following capabilities)

Evaluation Methodology:

The same evaluation set of 100 patient issues was used across all models, with responses generated using consistent temperature (0.7) and top-p (0.9) settings. The same semantic similarity and empathy scoring methods were applied to ensure fair comparison.

9.6.2 Comparative Results

Semantic Similarity to Human Therapist Responses:

Model	Mean Similarity	Median Similarity	Std Deviation
Fine-tuned Llama 3.2 3B	0.791	0.804	0.076
Gemma 2 2B	0.705	0.721	0.093
Phi-2 2.7B	0.683	0.697	0.102
Mistral 7B Instruct v0.2	0.742	0.756	0.084

Table XI Semantic Similarity to Human Therapist Responses compared to other open source models

Empathy Scores:

Model	Mean Empathy	Human Empathy	Empathy Gap
Fine-tuned Llama 3.2 3B	0.746	0.824	0.078
Gemma 2 2B	0.623	0.824	0.201
Phi-2 2.7B	0.598	0.824	0.226
Mistral 7B Instruct v0.2	0.701	0.824	0.123

Table XII Empathy scores for Human Therapist Responses compared to other open source models

9.6.3 Analysis of Comparative Performance

The benchmark results reveal several key insights about the relative performance of the fine-tuned Llama 3.2 3B model:

1. Domain Adaptation Advantage

The fine-tuned Llama 3.2 3B model outperforms all other models in semantic similarity and empathy scores, despite some of them having similar or larger parameter counts. This highlights the critical importance of domain-specific fine-tuning for therapeutic applications. The domain adaptation provides an advantage that exceeds what might be expected from parameter count alone.

2. Size-Performance Relationship

The comparison with Mistral 7B, which is more than twice the size of Llama 3.2 3B, is particularly informative. While the larger Mistral model performs well (second-best across metrics), it does not match the fine-tuned Llama 3.2 3B despite its size advantage. This suggests that for specialized domains like therapy:

- Targeted fine-tuning on high-quality domain data can outweigh raw parameter count advantages
- Smaller, well-tuned models can be more effective than larger general-purpose models
- The quality and relevance of adaptation data may be more important than model scale within certain bounds

3. Efficiency Considerations

The Llama 3.2 3B model achieves its superior performance with reasonable efficiency, generating responses only marginally slower than the smaller models and significantly faster than the larger Mistral model. This favorable balance of quality and efficiency makes it particularly suitable for practical deployment in mental health support applications.

4. Response Quality Patterns

Qualitative analysis of responses across models reveals distinct patterns:

- The fine-tuned Llama 3.2 3B produces more therapeutically sound responses with appropriate validation, exploration, and guidance components.
- Gemma 2 and Phi-2 tend toward more generic supportive responses with less therapeutic depth and structure.
- Mistral 7B produces more verbose responses with good therapeutic elements but occasional digressions or unnecessary elaborations.

These patterns suggest that the fine-tuning process has successfully captured core therapeutic communication patterns that other models either lack or implement less consistently.

9.7 Implications of Benchmarking Results

The comparative analysis yields several practical implications:

1. Fine-tuning ROI

The significant performance improvements achieved through fine-tuning indicate a high return on investment for domain adaptation in mental health applications. Organizations developing therapeutic AI should prioritize quality domain-specific datasets and thorough fine-tuning over simply using larger general models.

2. Model Selection Guidance

For mental health support applications with similar computational constraints, the benchmarking results strongly favor the Llama 3.2 architecture with domain-specific fine-tuning as the most effective approach. This provides clear guidance for implementation decisions.

3. Future Research Directions

The impressive performance of the relatively compact Llama 3.2 3B model suggests promising avenues for future research:

- Exploring hybrid approaches that combine the efficiency of smaller models with techniques to enhance empathy specifically
- Investigating specialized architectures or training techniques focused on emotional intelligence
- Developing more sophisticated evaluation metrics that capture nuanced aspects of therapeutic communication

These research directions could further advance the capabilities of language models in mental health support contexts while maintaining reasonable computational requirements.

In summary, the benchmarking analysis confirms that the fine-tuned Llama 3.2 3B model represents a particularly effective balance of performance, efficiency, and therapeutic quality compared to other available open-source models, justifying its selection for mental health support applications.

10. DISCUSSION

10.1 Chapter Overview

This chapter interprets and critically reflects on the results presented in the previous section. It explores the implications of fine-tuning large language models (LLMs) for mental health support, the effectiveness of different configurations, and the trade-offs encountered between semantic similarity and emotional intelligence. Additionally, this chapter discusses user perceptions of MirrorAI and how the system compares to both human therapists and other state-of-the-art open-source models. The findings are situated within the broader context of existing literature, highlighting how this project contributes to and advances the current research landscape.

10.2 Interpretation of Model Performance

The experimental results show that **fine-tuning LLMs using mental health-specific datasets** significantly enhances their ability to deliver emotionally intelligent and contextually appropriate responses. Notably, the LLaMA 3.2 3B model achieved the highest semantic similarity (0.634) to therapist responses, demonstrating the role of **model size** in aligning with professional therapeutic language.

However, a striking finding is that the **smaller LLaMA 3.2 1B model slightly outperformed the 3B model in empathy scores**. This suggests that while larger models are better at mimicking the language patterns of therapists, smaller models may be better at preserving emotional cues during fine-tuning. This insight challenges the commonly held belief that larger models are always superior and highlights the nuanced relationship between model capacity and emotional expression.

Moreover, the **learning rate and LoRA rank settings** played a critical role in model performance. A learning rate of $5e-4$ was found to produce the most empathetic responses (0.0841), nearly matching the therapist benchmark (0.087). Conversely, lower LoRA ranks (e.g., $r = 8$) led to **drastic empathy degradation**, even when similarity scores remained stable. This indicates that **parameter-efficient fine-tuning must be approached with caution**, as aggressive compression can damage a model's capacity to convey care and emotional understanding.

10.3 Balancing Empathy and Semantic Accuracy

One of the most important findings is the **inverse relationship between semantic similarity and empathy**. Models like DeepSeek V3 ranked highest in similarity to therapist responses but scored poorly on empathy. On the other hand, MirrorAI, while not ranking highest in similarity, demonstrated a much **stronger ability to produce emotionally resonant responses**.

This trade-off raises an important question: what matters more in AI-driven mental health support—**matching the therapist's words or capturing their intent and emotional tone**? The results suggest that emotional authenticity, as measured through empathy scores, may be more valuable in mental health contexts than syntactic or semantic closeness. This aligns with psychological literature that emphasizes empathy and emotional connection as core elements of effective therapy.

Thus, MirrorAI's strength lies in its **emotional understanding**, making it potentially more effective in providing comfort and support than models focused purely on linguistic accuracy.

10.4 Insights from User Feedback

User feedback further validates the effectiveness of the MirrorAI system. Participants generally responded positively to the emotional tone and helpfulness of the chatbot's responses. Many users appreciated the **calm, structured, and non-judgmental tone** used in replies, indicating that the chatbot successfully mimicked aspects of therapeutic presence.

Some users noted **repetitiveness** in responses, a common issue in LLMs, and suggested features such as personalized mood tracking and adjustable voice personalities. This highlights a **demand for more adaptive and personalized AI mental health tools** that evolve based on user preferences and emotional states.

Feedback also revealed that **voice interaction**, while promising, still feels robotic at times. This indicates a future opportunity for integrating **multimodal emotional cues** such as prosody and tone modulation, to further enhance realism and therapeutic engagement.

10.5 Comparison to Related Work

Compared to other projects like **ChatCounselor (Liu et al., 2023)** and **CBT-LLM (Na, 2024)**, MirrorAI demonstrates a unique strength: it does not only focus on task performance (e.g., generating correct CBT responses) but also emphasizes **emotional quality**. Additionally, this project contributes by exploring **parameter-efficient fine-tuning** methods (e.g., LoRA) in depth—something not extensively covered in earlier works.

Compared to commercial tools like **Wysa** or **Woebot**, which rely heavily on pre-scripted dialogues, MirrorAI offers **dynamic, unscripted conversations** powered by LLMs, enhancing personalization and user engagement.

Unlike **SERMO**, which focused on emotion regulation without real-time adaptability, MirrorAI incorporates both **real-time interaction and dynamic empathy tuning**, representing an important step forward in the development of empathetic AI companions.

10.6 Contributions to the Field

This project contributes to the field of AI mental health in several key ways:

- It demonstrates the **viability of using fine-tuned open-source LLMs** to replicate emotionally intelligent therapeutic interactions.
- It establishes a **dual-metric evaluation framework** combining semantic similarity and empathy scoring, offering a more holistic way to assess conversational quality.
- It identifies **optimal fine-tuning parameters**, such as learning rate and LoRA rank, that balance emotional and linguistic performance.
- It highlights the **importance of emotional fidelity** in AI mental health tools, suggesting that future research should prioritize empathy over rigid accuracy.

These contributions help bridge the gap between AI technical development and real-world psychological applicability.

10.7 Limitations Revisited

Despite its promising outcomes, this project faced several limitations:

- **Computational Constraints:** Due to limited resources, experimentation was restricted to relatively small models (1B and 3B). Larger models like LLaMA 7B or Mistral could not be tested.
- **Text-Only Evaluation:** Emotional scoring was based solely on textual output, without integrating voice tone or body language cues, which are essential in real therapy.
- **Synthetic Dataset Bias:** The synthetic dataset, while helpful in training, may introduce artificial patterns that do not fully reflect real-life conversational diversity.
- **Limited User Testing:** While initial feedback was promising, broader and longer-term trials are needed to assess sustained user engagement and therapeutic impact.

10.8 Implications for Future Work

The findings suggest that future AI mental health systems should:

- Explore **multi-modal empathy detection** by incorporating voice tone, pauses, and speech patterns.
- Develop **custom benchmarks** for empathy scoring that go beyond binary emotion classifiers.
- Integrate **user-personalization features** like session history, mood tracking, and adaptive response styles.
- Experiment with **larger LLMs and longer fine-tuning cycles** to further push the boundaries of emotional understanding in AI.

The discussion reinforces that while LLMs are not yet a replacement for human therapists, they hold immense potential in offering **emotionally aware, accessible, and consistent** support. MirrorAI successfully demonstrated that with the right data, fine-tuning strategies, and evaluation framework, AI systems can come significantly close to replicating the empathy and responsiveness of real therapeutic conversations.

By focusing not just on technical accuracy but also **emotional impact**, this project lays the foundation for a new generation of AI mental health tools that are not only intelligent—but compassionate.

11.0 CONCLUSION

11.1 Summary of Findings

This project set out to explore how fine-tuned large language models (LLMs) could be used to develop an emotionally intelligent, voice-assisted mental health companion—**MirrorAI**. The results demonstrate that with proper fine-tuning, open-source LLMs can simulate human-like therapeutic responses with a meaningful degree of empathy and contextual understanding.

The findings revealed that **larger models (e.g., LLaMA 3.2 3B)** performed better in terms of semantic similarity to therapist responses, while **smaller models (e.g., LLaMA 3.2 1B)** slightly outperformed in empathy scores. The experiments also showed that **higher learning rates** preserved

emotional tone more effectively, and **LoRA rank** had a major influence on empathy preservation—lower ranks significantly degraded emotional quality. These insights support the claim that **fine-tuning techniques and parameter configurations** play a critical role in achieving a balance between linguistic accuracy and emotional resonance.

User feedback confirmed that MirrorAI was generally perceived as helpful, empathetic, and easy to engage with, although areas like personalization and natural voice expression were highlighted for improvement.

11.2 Key Contributions

This research made several meaningful contributions:

- Demonstrated the **practical viability of open-source LLMs** for emotionally aware mental health applications.
- Designed and tested a **dual-metric evaluation framework** combining semantic similarity and empathy scoring.
- Identified **optimal fine-tuning configurations** that balance emotional intelligence with response accuracy.
- Developed a working prototype with a **real-time interface and voice integration**, showcasing the feasibility of deploying therapeutic LLMs in accessible formats.

11.3 Limitations

Despite its successes, the project faced several limitations:

- **Computational Resources:** Experiments were restricted to smaller model variants due to limited hardware access.
- **Evaluation Scope:** Empathy measurement was conducted only on text-based responses, without assessing vocal delivery or real-world behavior over long-term use.
- **Dataset Constraints:** While the Psych8k and synthetic datasets were helpful, they may not fully represent the cultural and linguistic diversity of global users.

- **Real-World Deployment:** The prototype was tested in controlled conditions and requires further evaluation in real-world scenarios to understand its broader impact.

11.4 Future Work

To build upon the findings of this study, several future improvements are suggested:

- **Multilingual and Culturally Sensitive Support:** Expanding training datasets to include diverse languages and cultural contexts will make MirrorAI more globally inclusive.
- **Deeper Empathy Simulation:** Future versions can explore multi-modal models that incorporate voice tone, speech pauses, and facial cues for enhanced emotional intelligence.
- **Long-Term User Studies:** Conducting extended trials with mental health professionals and users will offer richer insights into usability, safety, and therapeutic effectiveness.
- **Dynamic Personalization:** Introducing features like mood tracking, memory of past sessions, and adaptive response styles can significantly enhance user engagement.
- **Real-World Deployment:** Integrating MirrorAI into accessible platforms such as mobile apps or community health portals can expand its reach and practical application.

11.5 Final Reflection

MirrorAI represents a significant step toward creating empathetic and accessible mental health tools using modern AI. While it does not aim to replace therapists, it demonstrates the powerful role LLMs can play in supporting emotional well-being—especially for individuals without immediate access to professional care. With further refinement, such systems have the potential to complement traditional therapy and contribute to a more compassionate digital future in mental health support.

REFERENCE

Bill, D. and Eriksson, T., 2023. Fine-tuning a LLM using Reinforcement Learning from Human Feedback for a Therapy Chatbot Application. Royal Institute of Technology, Stockholm, Sweden.

Crasto, R., Dias, L., Miranda, D. and Kayande, D., 2021. CareBot: A Mental Health ChatBot. In: 2021 2nd International Conference for Emerging Technology (INCET). Belgaum, India: IEEE.

Denecke, K., Vaaheesan, S. and Arulnathan, A., 2021. A Mental Health Chatbot for Regulating Emotions (SERMO) - Concept and Usability Test. IEEE Transactions on Emerging Topics in Computing, 9(3), pp.1170-1179.

Liu, J.M., Li, D., Ren, T., Liao, Z. and Wu, J., 2023. ChatCounselor: A Large Language Model for Mental Health Support. In: Proceedings of PGAI CIKM 2023, Birmingham, UK.

Na, H., 2024. CBT-LLM: A Chinese Large Language Model for Cognitive Behavioral Therapy-based Mental Health Question Answering. Australian Artificial Intelligence Institute, University of Technology Sydney, Australia.

Omarov, B., Zhumanov, Z., Gumar, A. and Kuntunova, L., 2023. Artificial Intelligence Enabled Mobile Chatbot Psychologist using AIML and Cognitive Behavioral Therapy. International Journal of Advanced Computer Science and Applications, 14(6), pp.137-145.

Pandey, S., Sharma, S. and Wazir, S., 2022. Mental healthcare chatbot based on natural language processing and deep learning approaches: Ted the therapist. International Journal of Information Technology, 14(7), pp.3757-3766.

Qing Yu, H. and McGuinness, S., 2024. An Experimental Study of Integrating Fine-tuned LLMs and Prompts for Enhancing Mental Health Support Chatbot System. Journal of Medical Artificial Intelligence, pp.1-24.

Rathnayaka, P., Mills, N., Burnett, D., De Silva, D. and Gray, R., 2022. A Mental Health Chatbot with Cognitive Skills for Personalised Behavioural Activation and Remote Health Monitoring. Sensors, 22(10), 3653.

Xu, X., Yao, B., Dong, Y., Gabriel, S., Yu, H., Hendler, J., Ghassemi, M., Dey, A.K. and Wang, D., 2024. Mental-LLM: Leveraging Large Language Models for Mental Health Prediction via Online Text Data. Proceedings of the ACM on Interactive, Mobile, Wearable, and Ubiquitous Technologies, 8(1), pp.1-32.

Bill, D. and Eriksson, T. (2023) Fine-tuning a LLM Using Reinforcement Learning From Human Feedback for a Therapy Chatbot Application. Royal Institute of Technology, Stockholm, Sweden.

Crasto, R., Dias, L., Miranda, D. and Kayande, D. (2021) CareBot: A Mental Health ChatBot. In: 2021 2nd International Conference for Emerging Technology (INCET). Belgaum, India: IEEE.

Denecke, K., Vaaheesan, S. and Arulnathan, A. (2021) A Mental Health Chatbot for Regulating Emotions (SERMO) – Concept and Usability Test. IEEE Transactions on Emerging Topics in Computing, 9(3), pp.1170-1179.

Liu, J.M., Li, D., Ren, T., Liao, Z. and Wu, J. (2023) ChatCounselor: A Large Language Model for Mental Health Support. In: Proceedings of PGAI CIKM 2023, Birmingham, UK.

Na, H. (2024) CBT-LLM: A Chinese Large Language Model for Cognitive Behavioral Therapy-based Mental Health Question Answering. Australian Artificial Intelligence Institute, University of Technology Sydney, Australia.

World Health Organization (2021) Mental disorders. Available at: [WHO website] (Accessed: 1 March 2025).

APPENDIX

GitHub Repository Link : <https://github.com/rikzyj/MirrorAI-Research-Project.git>

HuggingFace Model Repository: https://huggingface.co/rikzyj/CBT_Therapist_LLMA3_2_3B

Data Preprocessing and Feature Engineering

```
# -*- coding: utf-8 -*-
"""Data Preprocessing.ipynb

Automatically generated by Colab.

Original file is located at
    https://colab.research.google.com/drive/1hwZrQEwwpYum3mhMbhIIPQJKWfQZq-tP

# Data Preprocessing for Psych8k Dataset
"""

from google.colab import drive
```

```

drive.mount('/content/drive')

# OpenAI API Key
openai.api_key = "your_openai_api_key_here"

# Step 1: Load the Psych8k dataset and inspect
import pandas as pd

# Define the dataset path
psych8k_path = "/content/drive/MyDrive/ACADEMICS/FYP/4.
Implementation/Datasets/Psych8k/Psych8k.csv"

# Load the dataset
psych8k_df = pd.read_csv(psych8k_path)

# Display the first few rows to understand the structure
print("Dataset Loaded Successfully!")
psych8k_df.head()

# Step 2: Validate and clean the Psych8k dataset

# Check for missing values
missing_values = psych8k_df.isnull().sum()
print("Missing Values in Each Column:\n", missing_values)

# Drop rows with missing values
psych8k_df = psych8k_df.dropna()

# Remove duplicates
psych8k_df = psych8k_df.drop_duplicates()

# Verify dataset shape after cleaning
print(f"Dataset shape after cleaning: {psych8k_df.shape}")

pip install openai==0.28

# STEP 1: Install required library
# !pip install openai --quiet

# STEP 2: Import libraries
import pandas as pd
import openai
import ast
import re
import time
from tqdm import tqdm

# STEP 3: Set your OpenAI API Key (replace with your own)
openai.api_key = "sk-
BQeOeWHdHJvS4uJYMkf76e0oUYgrxeZM4WNk6BNZS5T3B1bkFJEt3ksOTNYzvY5i6UIH47ykF2rS1-
71zgLOzWhzh1IA" # <- 📌 Set securely via environment variable in production

```

```

# STEP 4: Load only the first 100 rows for testing
psych8k_path = "/content/drive/MyDrive/ACADEMICS/FYP/4.
Implementation/Datasets/Psych8k/Psych8k.csv"
df = pd.read_csv(psych8k_path)

# STEP 5: Function to get cognitive distortions from GPT-4o-mini
def get_distortions_batch(batch_inputs):
    prompt = (
        "You are a mental health assistant trained in Cognitive Behavioral Therapy (CBT). "
        "For each user input below, identify any cognitive distortions present "
        "from the following list: all-or-nothing thinking, overgeneralization, "
        "mental filter, "
        "disqualifying the positive, jumping to conclusions, magnification, "
        "emotional reasoning, "
        "should statements, labeling, personalization, blaming others.\n\n"
        "Return a JSON-style Python list of strings, where each item is the "
        "distortion(s) "
        "for the corresponding input. If none are found, return \"none\"."
    )
    Example: \n
    ["should statements", "none", "labeling, personalization"]\n\n
    Inputs: \n +
    "\n".join([f"{i+1}. {text}" for i, text in enumerate(batch_inputs)])

    try:
        response = openai.ChatCompletion.create(
            model="gpt-4o",
            messages=[{"role": "user", "content": prompt}],
            temperature=0.2,
            max_tokens=1024
        )
        # extract inside code block if wrapped in ```python ... ```
        output_text = response['choices'][0]['message']['content'].strip()
        match = re.search(r"```(?:python)?\s*([^\s\S]*?)\s*```", output_text)
        if match:
            output_text = match.group(1)
            return ast.literal_eval(output_text)
    except Exception as e:
        print("⚠ Error parsing output:", e)
        return ["error"] * len(batch_inputs)

# STEP 6: Run in batches
batch_size = 20
results = []

for i in tqdm(range(0, len(df), batch_size)):
    batch = df['input'].iloc[i:i+batch_size].tolist()
    distortions = get_distortions_batch(batch)
    results.extend(distortions)

```

```

    time.sleep(3) # <-- Added 15-second delay

# STEP 7: Add results to DataFrame
df['cognitive_distortions'] = results

# STEP 8: Save to new CSV
df.to_csv("Psych8k_with_distortions_sample.csv", index=False)

# STEP 9: Preview
df.head()

# Assume `results` is your list of distortion predictions
# Pad it if it's shorter than the DataFrame
if len(results) < len(df):
    print(f"Patching: {len(df) - len(results)} missing values detected.")
    results += ["not_processed"] * (len(df) - len(results))
elif len(results) > len(df):
    print(f"Warning: {len(results) - len(df)} extra values. Trimming.")
    results = results[:len(df)]

# Now assign safely
df["cognitive_distortions"] = results

df.to_csv("Psych8k_with_distortions_final.csv", index=False)

df.head(50)

import pandas as pd
import json

# Step 1: Load the Psych8k dataset with cognitive distortions
psych8k_path = "/content/drive/MyDrive/ACADEMICS/FYP/4.
Implementation/Datasets/Psych8k/Psych8k_with_distortions.csv" # Path to the CSV
file
psych8k_df = pd.read_csv(psych8k_path)

# Step 4: Reformat the dataset into a standardized conversation structure
conversations = [] # To store the JSON structure for all conversations
for idx, row in psych8k_df.iterrows():
    user_content = row["input"]
    therapist_content = row["output"]
    distortions = row["cognitive_distortions"]

# Build the conversation object
conversation = {
    "conversation_id": f"psych8k_{idx}", # Unique conversation ID
    "turns": [
        {"role": "user", "content": user_content},
        {"role": "therapist", "content": therapist_content}
    ],
    "metadata": {

```



```

        "source": "psych8k",
        "type": "single-turn", # Psych8k is single-turn conversations
        "turn_count": 2, # Each conversation has 2 turns (user + therapist)
        "cognitive_distortions": distortions.split(", ") if
isinstance(distortions, str) else [] # Split distortions into a list
    }
}
conversations.append(conversation)

# Step 5: Save the preprocessed dataset as a JSON file
output_path = "Preprocessed_Psych8k_with_Cognitive_Distortions.json" # Output
filename
with open(output_path, "w") as json_file:
    json.dump(conversations, json_file, indent=4)

print(f"Preprocessed dataset saved to: {output_path}")

# Step 3: Reformat dataset into standardized conversation structure

def create_conversation(row, conversation_id):
    # Create a conversation object
    return {
        "conversation_id": f"psych_{conversation_id}",
        "turns": [
            {"role": "user", "content": row['input']},
            {"role": "therapist", "content": row['output']}
        ],
        "metadata": {
            "source": "psych8k",
            "type": "single-turn",
            "turn_count": 2,
            "cognitive_distortions": [] # Initialize with an empty list
        }
    }

# Apply the transformation
conversations = [create_conversation(row, idx) for idx, row in
psych8k_df.iterrows()]

# Check an example conversation
print("Example of a reformatted conversation:")
print(conversations[0])

psych8k_df.info()

# Step 5: Save the preprocessed dataset as a JSON file
import json

# Define output path
output_path = "/content/drive/MyDrive/ACADEMICS/FYP/4.
Implementation/Datasets/Psych8k/Preprocessed_Psych8k.json"

```

```

# Save conversations to JSON
with open(output_path, 'w') as json_file:
    json.dump(conversations, json_file, indent=4)

print(f"Preprocessed dataset saved to: {output_path}")

"""# Data Preprocessing for CBT Conversations Dataset"""

# Step 1: Load the dataset and inspect
import json

# Define the dataset path
cbt_dataset_path = "/content/drive/MyDrive/ACADEMICS/FYP/4.
Implementation/Datasets/SynthDataset for FTing/cleaned_cbt_therapy_dataset.json"

# Load the dataset
with open(cbt_dataset_path, 'r') as f:
    cbt_data = json.load(f)

# Check the number of conversations and display an example
print(f"Total conversations in the dataset: {len(cbt_data)}")
print("Example of a conversation:")
print(json.dumps(cbt_data[0], indent=4))

# Step 2: Validate and clean conversations

# Remove incomplete conversations (less than 2 turns)
cbt_data = [conv for conv in cbt_data if len(conv.get('conversations', [])) >= 2]

# Standardize role labels
def standardize_roles(conversation):
    for turn in conversation['conversations']:
        if turn['role'].lower() in ['user', 'client', 'patient']:
            turn['role'] = 'user'
        elif turn['role'].lower() in ['therapist', 'counselor']:
            turn['role'] = 'therapist'
    return conversation

cbt_data = [standardize_roles(conv) for conv in cbt_data]

# Verify cleaned data
print(f"Conversations after cleaning: {len(cbt_data)}")
print("Example of a cleaned conversation:")
print(json.dumps(cbt_data[0], indent=4))

# Step 3: Structure conversations into a consistent format

def structure_conversation(conversation, conversation_id):
    return {
        "conversation_id": f"cbt_{conversation_id}",

```

```

        "turns": [
            {"role": turn['role'], "content": turn['content']}
            for turn in conversation['conversations']
        ],
        "metadata": {
            "source": "cbt",
            "type": "multi-turn",
            "turn_count": len(conversation['conversations'])
        }
    }

# Apply the transformation
structured_conversations = [structure_conversation(conv, idx) for idx, conv in
                             enumerate(cbt_data)]

# Check an example structured conversation
print("Example of a structured conversation:")
print(json.dumps(structured_conversations[0], indent=4))

# Step 4: Detect cognitive distortions and add tags
import re

# Cognitive Distortions Dictionary
cognitive_distortions = {
    "all_or_nothing_thinking": [
        "always", "never", "completely", "totally", "entirely", "perfect",
        "failure", "ruined", "no middle ground", "black and white"
    ],
    "overgeneralization": [
        "everyone", "nobody", "all the time", "never works out", "everywhere",
        "forever", "it's bound to happen", "this always happens"
    ],
    "mental_filter": [
        "only see the negative", "focus on flaws", "ignore the positive", "nothing
        good", "all I notice is the bad", "all problems, no solutions"
    ],
    "disqualifying_the_positive": [
        "that doesn't count", "it was just luck", "anyone could have done it", "I
        don't deserve credit", "it's not a big deal", "they're just being nice"
    ],
    "jumping_to_conclusions_mind_reading": [
        "I know they think I'm a failure", "they must hate me", "they're judging
        me", "everyone sees I'm worthless", "I can tell they don't like me"
    ],
    "jumping_to_conclusions_fortune_telling": [
        "it will never work", "I'm doomed", "this is going to fail", "I just know
        it's going to end badly", "I can't see any good outcome"
    ],
    "magnification_or_minimization": [

```

```

        "it's the end of the world", "this is a total disaster", "it's nothing
special", "downplay my success", "make a mountain out of a molehill", "it's not
that important"
    ],
    "emotional_reasoning": [
        "I feel it, so it must be true", "if I'm scared, there's real danger", "I
feel worthless, so I am worthless", "my emotions are facts", "I feel guilty, so I
must have done something wrong"
    ],
    "should_statements": [
        "I should be able to handle this", "I must do better", "I have to be
perfect", "I shouldn't feel this way", "I ought to be stronger", "I shouldn't make
mistakes"
    ],
    "labeling": [
        "I'm a loser", "I'm worthless", "I'm a failure", "I'm useless", "I'm
stupid", "I'm incompetent"
    ],
    "personalization": [
        "it's all my fault", "I caused this", "I messed up everyone's day", "I'm
responsible for everything", "if only I had done better"
    ],
    "blaming_others": [
        "it's all their fault", "they made this happen", "they ruined everything",
"they are the reason I'm like this", "I wouldn't feel this way if not for them"
    ]
}

# Function to detect cognitive distortions
def detect_cognitive_distortions(conversation):
    detected_distortions = set()
    for turn in conversation['turns']:
        if turn['role'] == 'user': # Only analyze the user's turns
            for distortion, phrases in cognitive_distortions.items():
                if any(re.search(rf'\b{phrase}\b', turn['content'], re.IGNORECASE)
for phrase in phrases):
                    detected_distortions.add(distortion)
    return list(detected_distortions)

# Add detected cognitive distortions to metadata
for conversation in structured_conversations:
    distortions = detect_cognitive_distortions(conversation)
    conversation['metadata']['cognitive_distortions'] = distortions

# Check an example with added cognitive distortion metadata
print("Example conversation with cognitive distortion tags:")
print(json.dumps(structured_conversations[0], indent=4))

print(json.dumps(structured_conversations[974], indent=4))

# Step 5: Save the preprocessed dataset as a JSON file

```

```

# Define output path
output_path = "/content/drive/MyDrive/ACADEMICS/FYP/4.
Implementation/Datasets/SynthDataset for
FTing/Preprocessed_CBT_Dataset_Enriched.json"

# Save structured conversations to JSON
with open(output_path, 'w') as json_file:
    json.dump(structured_conversations, json_file, indent=4)

print(f"Preprocessed dataset saved to: {output_path}")

"""# Dataset Merging"""

/content/drive/MyDrive/ACADEMICS/FYP/4.
Implementation/Datasets/Psych8k/Psych8k_with_distortions.csv
/content/drive/MyDrive/ACADEMICS/FYP/4. Implementation/Datasets/SynthDataset for
FTing/Preprocessed_CBT_Dataset_Enriched.json

import pandas as pd
import json

# Step 1: Load the Psych8k dataset
psych8k_path = "/content/drive/MyDrive/ACADEMICS/FYP/4.
Implementation/Datasets/Psych8k/Psych8k_with_distortions.csv"
psych8k_df = pd.read_csv(psych8k_path)

# Step 2: Convert Psych8k dataset into JSON format of conversations
psych8k_conversations = []
for idx, row in psych8k_df.iterrows():
    conversation = {
        "conversation_id": f"psych8k_{idx}",
        "turns": [
            {"role": "user", "content": row["input"]},
            {"role": "therapist", "content": row["output"]}
        ],
        "metadata": {
            "source": "psych8k",
            "type": "single-turn",
            "turn_count": 2,
            "cognitive_distortions": row["cognitive_distortions"].split(", ") if
isinstance(row["cognitive_distortions"], str) else ["none"]
        }
    }
    psych8k_conversations.append(conversation)

# Step 3: Load the CBT dataset
cbt_path = "/content/drive/MyDrive/ACADEMICS/FYP/4.
Implementation/Datasets/SynthDataset for
FTing/Preprocessed_CBT_Dataset_Enriched.json"
with open(cbt_path, "r") as cbt_file:

```

```

    cbt_conversations = json.load(cbt_file)

# Step 4: Validate and standardize CBT dataset
for conversation in cbt_conversations:
    # Add "none" if cognitive distortions are empty
    if not conversation["metadata"].get("cognitive_distortions"):
        conversation["metadata"]["cognitive_distortions"] = ["none"]
    # Ensure unique conversation IDs
    conversation["conversation_id"] = f"cbt_{conversation['conversation_id']}"

# Step 5: Merge the datasets
merged_conversations = psych8k_conversations + cbt_conversations

# Step 6: Save the merged dataset as a JSON file
output_path = "/content/drive/MyDrive/ACADEMICS/FYP/4.
Implementation/Datasets/Psych8k_and_CBT_Conversation_Dataset.json"
with open(output_path, "w") as output_file:
    json.dump(merged_conversations, output_file, indent=4)

print(f"Merged dataset saved to: {output_path}")

print(merged_conversations[500])

```

Model Training

```

# -*- coding: utf-8 -*-
"""[WORKING1...] Llama3.2_(1B_and_3B)-Conversational.ipynb

Automatically generated by Colab.

Original file is located at
    https://colab.research.google.com/drive/1ocj8GI9imtka4v_Trut2x0Y45MT-utfM

### Installation
"""

# Commented out IPython magic to ensure Python compatibility.
# %%capture
# import os
# if "COLAB_" not in "".join(os.environ.keys()):
#     !pip install unsloth
# else:
#     # Do this only in Colab and Kaggle notebooks! Otherwise use pip install
unsloth
#     !pip install --no-deps bitsandbytes accelerate xformers==0.0.29 peft trl
triton
#     !pip install --no-deps cut_cross_entropy unsloth_zoo
#     !pip install sentencepiece protobuf datasets huggingface_hub hf_transfer

```

```

# !pip install --no-deps unsloth

from google.colab import drive
drive.mount('/content/drive')

"""### Unsloth"""

from unsloth import FastLanguageModel
import torch
max_seq_length = 2048 # Choose any! We auto support RoPE Scaling internally!
dtype = None # None for auto detection. Float16 for Tesla T4, V100, Bfloat16 for
               Ampere+
load_in_4bit = True # Use 4bit quantization to reduce memory usage. Can be False.

# 4bit pre quantized models we support for 4x faster downloading + no OOMs.
fourbit_models = [
    "unsloth/Meta-Llama-3.1-8B-bnb-4bit",          # Llama-3.1 2x faster
    "unsloth/Meta-Llama-3.1-8B-Instruct-bnb-4bit",
    "unsloth/Meta-Llama-3.1-70B-bnb-4bit",
    "unsloth/Meta-Llama-3.1-405B-bnb-4bit",        # 4bit for 405b!
    "unsloth/Mistral-Small-Instruct-2409",        # Mistral 22b 2x faster!
    "unsloth/mistral-7b-instruct-v0.3-bnb-4bit",
    "unsloth/Phi-3.5-mini-instruct",              # Phi-3.5 2x faster!
    "unsloth/Phi-3-medium-4k-instruct",
    "unsloth/gemma-2-9b-bnb-4bit",
    "unsloth/gemma-2-27b-bnb-4bit",               # Gemma 2x faster!

    "unsloth/Llama-3.2-1B-bnb-4bit",              # NEW! Llama 3.2 models
    "unsloth/Llama-3.2-1B-Instruct-bnb-4bit",
    "unsloth/Llama-3.2-3B-bnb-4bit",
    "unsloth/Llama-3.2-3B-Instruct-bnb-4bit",

    "unsloth/Llama-3.3-70B-Instruct-bnb-4bit" # NEW! Llama 3.3 70B!
] # More models at https://huggingface.co/unsloth

model, tokenizer = FastLanguageModel.from_pretrained(
    model_name = "unsloth/Llama-3.2-3B-Instruct", # or choose "unsloth/Llama-3.2-
1B-Instruct"
    max_seq_length = max_seq_length,
    dtype = dtype,
    load_in_4bit = load_in_4bit,
    # token = "hf_...", # use one if using gated models like meta-llama/Llama-2-7b-
hf
)

"""We now add LoRA adapters so we only need to update 1 to 10% of all
parameters!"""

model = FastLanguageModel.get_peft_model(
    model,
    r = 16, # Choose any number > 0 ! Suggested 8, 16, 32, 64, 128

```

```

target_modules = ["q_proj", "k_proj", "v_proj", "o_proj",
                  "gate_proj", "up_proj", "down_proj",],
lora_alpha = 16,
lora_dropout = 0, # Supports any, but = 0 is optimized
bias = "none",    # Supports any, but = "none" is optimized
# [NEW] "unsloth" uses 30% less VRAM, fits 2x larger batch sizes!
use_gradient_checkpointing = "unsloth", # True or "unsloth" for very long
context
    random_state = 3407,
    use_rslora = False, # We support rank stabilized LoRA
    loftq_config = None, # And LoftQ
)

from unsloth.chat_templates import get_chat_template

tokenizer = get_chat_template(
    tokenizer,
    chat_template = "llama-3.1",
)

def formatting_prompts_func(examples):
    convos = examples["conversations"]
    texts = [tokenizer.apply_chat_template(convo, tokenize = False,
add_generation_prompt = False) for convo in convos]
    return { "text" : texts, }
pass

# Import necessary libraries
import json
from datasets import Dataset
from unsloth.chat_templates import get_chat_template, standardize_sharegpt
from transformers import AutoTokenizer

# File path
json_file_path = "/content/drive/MyDrive/ACADEMICS/FYP/4.
Implementation/Datasets/SynthDataset for Fting/cleaned_cbt_therapy_dataset.json"

# Load JSON and convert to list of dictionaries
def load_json_to_dict(filepath):
    with open(filepath, 'r', encoding='utf-8') as f:
        data = json.load(f)

    processed_data = []
    for conversation in data:
        if "conversations" in conversation and
isinstance(conversation["conversations"], list):
            # Ensure conversation is not empty
            if len(conversation["conversations"]) == 0:
                continue

            # Rename "therapist" role to "assistant" for compatibility

```



```

        conversations = [
            {"role": "assistant" if message["role"] == "therapist" else
message["role"],
            "content": message["content"]}
        for message in conversation["conversations"]
    ]
    processed_data.append({"conversations": conversations})

    return processed_data

# Step 1: Load and process the dataset
data = load_json_to_dict(json_file_path)

# Step 2: Convert to HuggingFace Dataset
dataset = Dataset.from_list(data)

# Step 3: Standardize the Dataset
dataset = standardize_sharegpt(dataset, aliases_for_assistant=["assistant",
"therapist"])

# Step 4: Initialize Tokenizer and Configure Chat Template
tokenizer = AutoTokenizer.from_pretrained("unsloth/Llama-3.2-1B-Instruct") #
Adjust model if needed
chat_template = get_chat_template(tokenizer, chat_template="llama-3.1")

def formatting_prompts_func(examples):
    texts = []

    for conversations in examples["conversations"]:
        # Ensure conversation is valid before applying template
        if not isinstance(conversations, list) or len(conversations) == 0:
            texts.append("") # Append empty string for consistency
            continue

        formatted_text = chat_template.apply_chat_template(
            conversations,
            tokenize=False,
            add_generation_prompt=False
        )
        texts.append(formatted_text)

    return {"text": texts}

# Apply formatting to the dataset
dataset = dataset.map(formatting_prompts_func, batched=True,
remove_columns=["conversations"])

# Final Output
print("Dataset has been processed and is ready for fine-tuning.")
#dataset = load_dataset("mlabonne/FineTome-100k", split = "train")

```

```

from unsloth.chat_templates import get_chat_template

tokenizer = get_chat_template(
    tokenizer,
    chat_template = "llama-3.1",
)

def formatting_prompts_func(examples):
    convos = examples["conversations"]
    texts = [tokenizer.apply_chat_template(convo, tokenize = False,
add_generation_prompt = False) for convo in convos]
    return { "text" : texts, }
pass

# Import necessary libraries
import json
from datasets import Dataset
from unsloth.chat_templates import get_chat_template, standardize_sharegpt
from transformers import AutoTokenizer

# File path
json_file_path = "/content/drive/MyDrive/ACADEMICS/FYP/4.
Implementation/Datasets/SynthDataset for FTing/cleaned_cbt_therapy_dataset.json"

# Load JSON and convert to list of dictionaries
def load_json_to_dict(filepath):
    with open(filepath, 'r', encoding='utf-8') as f:
        data = json.load(f)

    processed_data = []
    for conversation in data:
        if "conversations" in conversation and
isinstance(conversation["conversations"], list):
            if len(conversation["conversations"]) == 0:
                continue # Skip empty conversations

            # Ensure renaming of roles: therapist -> assistant
            conversations = [
                {"role": "assistant" if message["role"] == "therapist" else
message["role"],
                "content": message["content"]}
                for message in conversation["conversations"]
            ]
            processed_data.append({"conversations": conversations})

    return processed_data

# Step 1: Load and process the dataset
data = load_json_to_dict(json_file_path)

# Ensure data is valid before conversion

```

```

if len(data) == 0:
    raise ValueError("Error: No valid data found in JSON. Check the dataset
structure.")

# Step 2: Convert to Hugging Face Dataset
dataset = Dataset.from_list(data)

# **Ensure 'conversations' exists before standardization**
if "conversations" not in dataset.column_names:
    raise KeyError("Error: 'conversations' column is missing from dataset. Check
preprocessing steps.")

# Step 3: Standardize the Dataset (Ensuring therapist -> assistant)
dataset = standardize_sharegpt(dataset, aliases_for_assistant=["assistant",
"therapist"])

# Step 4: Initialize Tokenizer and Configure Chat Template
tokenizer = AutoTokenizer.from_pretrained("unsloth/Llama-3.2-1B-Instruct") #
Adjust model if needed
chat_template = get_chat_template(tokenizer, chat_template="llama-3.1")

def formatting_prompts_func(examples):
    texts = []

    for conversations in examples["conversations"]:
        if not isinstance(conversations, list) or len(conversations) == 0:
            texts.append("") # Ensure consistency
            continue

        formatted_text = chat_template.apply_chat_template(
            conversations,
            tokenize=False,
            add_generation_prompt=False
        )
        texts.append(formatted_text)

    return {"text": texts}

# Apply formatting to the dataset
formatted_dataset = dataset.map(formatting_prompts_func, batched=True,
remove_columns=["conversations"])

# Final Output
print("Dataset has been processed and is ready for fine-tuning.")

from unsloth.chat_templates import standardize_sharegpt
dataset = standardize_sharegpt(dataset)
dataset = dataset.map(formatting_prompts_func, batched = True,)

print(dataset)

```

```

dataset[5]["text"]

from trl import SFTTrainer
from transformers import TrainingArguments, DataCollatorForSeq2Seq
from unsloth import is_bfloat16_supported

trainer = SFTTrainer(
    model = model,
    tokenizer = tokenizer,
    train_dataset = dataset,
    dataset_text_field = "text",
    max_seq_length = max_seq_length,
    data_collator = DataCollatorForSeq2Seq(tokenizer = tokenizer),
    dataset_num_proc = 2,
    packing = False, # Can make training 5x faster for short sequences.
    args = TrainingArguments(
        per_device_train_batch_size = 2,
        gradient_accumulation_steps = 4,
        warmup_steps = 5,
        # num_train_epochs = 1, # Set this for 1 full training run.
        max_steps = 60,
        learning_rate = 2e-4, #1e-5, 2e-4
        fp16 = not is_bfloat16_supported(),
        bf16 = is_bfloat16_supported(),
        logging_steps = 1,
        optim = "adamw_8bit",
        weight_decay = 0.01,
        lr_scheduler_type = "linear",
        seed = 3407,
        output_dir = "outputs",
        report_to = "none", # Use this for WandB etc
    ),
)

"""We verify masking is actually done:"""

from unsloth.chat_templates import train_on_responses_only
trainer = train_on_responses_only(
    trainer,
    instruction_part = "<|start_header_id|>user<|end_header_id|>\n\n",
    response_part = "<|start_header_id|>assistant<|end_header_id|>\n\n",
)

space = tokenizer(" ", add_special_tokens = False).input_ids[0]
tokenizer.decode([space if x == -100 else x for x in
trainer.train_dataset[5]["labels"]])

"""We can see the System and Instruction prompts are successfully masked!"""

# @title Show current memory stats
gpu_stats = torch.cuda.get_device_properties(0)

```

```

start_gpu_memory = round(torch.cuda.max_memory_reserved() / 1024 / 1024 / 1024, 3)
max_memory = round(gpu_stats.total_memory / 1024 / 1024 / 1024, 3)
print(f"GPU = {gpu_stats.name}. Max memory = {max_memory} GB.")
print(f"{start_gpu_memory} GB of memory reserved.")

trainer_stats = trainer.train()

# @title Show final memory and time stats
used_memory = round(torch.cuda.max_memory_reserved() / 1024 / 1024 / 1024, 3)
used_memory_for_lora = round(used_memory - start_gpu_memory, 3)
used_percentage = round(used_memory / max_memory * 100, 3)
lora_percentage = round(used_memory_for_lora / max_memory * 100, 3)
print(f"{trainer_stats.metrics['train_runtime']} seconds used for training.")
print(
    f"{round(trainer_stats.metrics['train_runtime']/60, 2)} minutes used for
training."
)
print(f"Peak reserved memory = {used_memory} GB.")
print(f"Peak reserved memory for training = {used_memory_for_lora} GB.")
print(f"Peak reserved memory % of max memory = {used_percentage} %.")
print(f"Peak reserved memory for training % of max memory = {lora_percentage} %.")

"""<a name="Inference"></a>
### Inference
Let's run the model! You can change the instruction and input - leave the output
blank!

**[NEW] Try 2x faster inference in a free Colab for Llama-3.1 8b Instruct
[here] (https://colab.research.google.com/drive/1T-YBVfnphoVc8E2E854qF3jdia2Ll2W2?usp=sharing)**

We use `min_p = 0.1` and `temperature = 1.5`. Read this
[Tweet] (https://x.com/menhguin/status/1826132708508213629) for more information on
why.
"""

from unsloth.chat_templates import get_chat_template

tokenizer = get_chat_template(
    tokenizer,
    chat_template = "llama-3.1",
)
FastLanguageModel.for_inference(model) # Enable native 2x faster inference

messages = [
    {"role": "user", "content": "I am feeling really stressed lately because of a
project i am working on"},
]
inputs = tokenizer.apply_chat_template(
    messages,
    tokenize = True,

```

```

        add_generation_prompt = True, # Must add for generation
        return_tensors = "pt",
    ).to("cuda")

outputs = model.generate(input_ids = inputs, max_new_tokens = 64, use_cache = True,
                        temperature = 1.5, min_p = 0.1)
tokenizer.batch_decode(outputs)

""" You can also use a `TextStreamer` for continuous inference - so you can see the
generation token by token, instead of waiting the whole time!"""

FastLanguageModel.for_inference(model) # Enable native 2x faster inference

messages = [
    {"role": "user", "content": "Continue the fibonnaci sequence: 1, 1, 2, 3, 5,
8,"},
]
inputs = tokenizer.apply_chat_template(
    messages,
    tokenize = True,
    add_generation_prompt = True, # Must add for generation
    return_tensors = "pt",
).to("cuda")

from transformers import TextStreamer
text_streamer = TextStreamer(tokenizer, skip_prompt = True)
_ = model.generate(input_ids = inputs, streamer = text_streamer, max_new_tokens =
128,
                  use_cache = True, temperature = 1.5, min_p = 0.1)

"""<a name="Save"></a>
### Saving, loading finetuned models
To save the final model as LoRA adapters, either use Huggingface's `push_to_hub`
for an online save or `save_pretrained` for a local save.

**[NOTE]** This ONLY saves the LoRA adapters, and not the full model. To save to
16bit or GGUF, scroll down!
"""

model.push_to_hub_merged(
    "rikzyj/CBT_Therapist_LLMA3_2_3B", # Replace with your Hugging Face
username/model_id
    tokenizer,
    save_method="merged_16bit",
    token="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" # Replace with your Hugging
Face token
)

from huggingface_hub import HfApi

# Assuming 'model' is your fine-tuned model and 'tokenizer' is your tokenizer

```

```

api = HfApi()

# Save the model and tokenizer locally using SafeTensors
model.save_pretrained("my_model", safe_serialization=True)
tokenizer.save_pretrained("my_model", safe_serialization=True)

# Create a new repository on Hugging Face
api.create_repo(repo_id="rikzyj/CBT_Therapist_LLMA3_2_1B_SafeTensors",
                repo_type="model",
                private=False,
                token="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")

# Push the SafeTensors files to the new repository
api.upload_folder(
    folder_path="my_model",
    repo_id="rikzyj/CBT_Therapist_LLMA3_2_1B_SafeTensors", # New repo ID
    repo_type="model",
    token="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" # Your Hugging Face token
)

model.save_pretrained("lora_model") # Local saving
tokenizer.save_pretrained("lora_model")
model.push_to_hub("rikzyj/CBT_Therapist_LLMA3_2_1B", token =
"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") # Online saving
# tokenizer.push_to_hub("your_name/lora_model", token = "...") # Online saving

"""Now if you want to load the LoRA adapters we just saved for inference, set
`False` to `True`:"""

if False:
    from unsloth import FastLanguageModel
    model, tokenizer = FastLanguageModel.from_pretrained(
        model_name = "lora_model", # YOUR MODEL YOU USED FOR TRAINING
        max_seq_length = max_seq_length,
        dtype = dtype,
        load_in_4bit = load_in_4bit,
    )
    FastLanguageModel.for_inference(model) # Enable native 2x faster inference

messages = [
    {"role": "user", "content": "Describe a tall tower in the capital of France."},
]
inputs = tokenizer.apply_chat_template(
    messages,
    tokenize = True,
    add_generation_prompt = True, # Must add for generation
    return_tensors = "pt",
).to("cuda")

from transformers import TextStreamer
text_streamer = TextStreamer(tokenizer, skip_prompt = True)

```

```
_ = model.generate(input_ids = inputs, streamer = text_streamer, max_new_tokens =
128,
                    use_cache = True, temperature = 1.5, min_p = 0.1)
```

"""You can also use Hugging Face's `AutoModelForPeftCausalLM`. Only use this if you do not have `unsloth` installed. It can be hopelessly slow, since `4bit` model downloading is not supported, and Unsloth's ****inference is 2x faster****."""

```
if False:
    # I highly do NOT suggest - use Unsloth if possible
    from peft import AutoPeftModelForCausalLM
    from transformers import AutoTokenizer
    model = AutoPeftModelForCausalLM.from_pretrained(
        "lora_model", # YOUR MODEL YOU USED FOR TRAINING
        load_in_4bit = load_in_4bit,
    )
    tokenizer = AutoTokenizer.from_pretrained("lora_model")
```

"""### Saving to float16 for VLLM

We also support saving to `float16` directly. Select `merged_16bit` for float16 or `merged_4bit` for int4. We also allow `lora` adapters as a fallback. Use `push_to_hub_merged` to upload to your Hugging Face account! You can go to <https://huggingface.co/settings/tokens> for your personal tokens.

```
# Merge to 16bit
if False: model.save_pretrained_merged("model", tokenizer, save_method =
"merged_16bit",)
if False: model.push_to_hub_merged("hf/model", tokenizer, save_method =
"merged_16bit", token = "")

# Merge to 4bit
if False: model.save_pretrained_merged("model", tokenizer, save_method =
"merged_4bit",)
if False: model.push_to_hub_merged("hf/model", tokenizer, save_method =
"merged_4bit", token = "")

# Just LoRA adapters
if False: model.save_pretrained_merged("model", tokenizer, save_method = "lora",)
if False: model.push_to_hub_merged("hf/model", tokenizer, save_method = "lora",
token = "")
```

"""### GGUF / llama.cpp Conversion

To save to `GGUF` / `llama.cpp`, we support it natively now! We clone `llama.cpp` and we default save it to `q8_0`. We allow all methods like `q4_k_m`. Use `save_pretrained_gguf` for local saving and `push_to_hub_gguf` for uploading to HF.

Some supported quant methods (full list on our [Wiki page](<https://github.com/unslothai/unsloth/wiki#gguf-quantization-options>)):
* `q8_0` - Fast conversion. High resource use, but generally acceptable.


```

* `q4_k_m` - Recommended. Uses Q6_K for half of the attention.wv and
feed_forward.w2 tensors, else Q4_K.
* `q5_k_m` - Recommended. Uses Q6_K for half of the attention.wv and
feed_forward.w2 tensors, else Q5_K.

[**NEW**] To finetune and auto export to Ollama, try our [Ollama
notebook] (https://colab.research.google.com/drive/1WZDi7APtQ9VsvOrQSSC5DDtxq159j8iZ?usp=sharing)
"""

# Save to 8bit Q8_0
if False: model.save_pretrained_gguf("model", tokenizer,)
# Remember to go to https://huggingface.co/settings/tokens for a token!
# And change hf to your username!
if False: model.push_to_hub_gguf("hf/model", tokenizer, token = "")

# Save to 16bit GGUF
if False: model.save_pretrained_gguf("model", tokenizer, quantization_method =
"f16")
if False: model.push_to_hub_gguf("hf/model", tokenizer, quantization_method =
"f16", token = "")

# Save to q4_k_m GGUF
if False: model.save_pretrained_gguf("model", tokenizer, quantization_method =
"q4_k_m")
if False: model.push_to_hub_gguf("hf/model", tokenizer, quantization_method =
"q4_k_m", token = "")

# Save to multiple GGUF options - much faster if you want multiple!
if False:
    model.push_to_hub_gguf(
        "hf/model", # Change hf to your username!
        tokenizer,
        quantization_method = ["q4_k_m", "q8_0", "q5_k_m",],
        token = "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", # Get a token at
https://huggingface.co/settings/tokens
    )

```

Model Evaluation

```

# -*- coding: utf-8 -*-
"""Experiment Model Evaluation.ipynb

```

Automatically generated by Colab.

Original file is located at

<https://colab.research.google.com/drive/1kcnbdcaWFzQKCDzRjF5e7zLdxsDBBKtU>

```

# Similarity Scores
"""

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sentence_transformers import SentenceTransformer
from sklearn.metrics.pairwise import cosine_similarity
import re

# Step 1: Load the CSV data
print("Loading data...")
df = pd.read_csv("/content/Experiment Summary - Sheet2.csv")
print(f"Loaded {len(df)} conversations")

print(df)

# Step 2: Clean text function
def clean_text(text):
    if not isinstance(text, str):
        return ""
    # Remove strange characters and normalize whitespace
    text = re.sub(r'[^\\w\\s.,!?\\'"]-', ' ', text)
    text = re.sub(r'\\s+', ' ', text).strip()
    return text

# Step 3: Load the sentence transformer model
print("Loading sentence transformer model...")
model = SentenceTransformer('all-MiniLM-L6-v2')

# Step 4: Calculate similarities
print("Calculating similarities...")
# Get LLM column names (all columns after the therapist_response)
llm_columns = df.columns[2:].tolist()

# Create dataframe to store results
results = pd.DataFrame(index=df.index)

# Process each row
for idx, row in df.iterrows():
    # Get therapist response
    therapist_text = clean_text(row['therapist_response'])
    if not therapist_text:
        continue

    # Get therapist embedding
    therapist_embedding = model.encode(therapist_text)

    # Calculate similarity for each LLM response
    for col in llm_columns:

```

```

    llm_text = clean_text(row[col])
    if not llm_text:
        results.loc[idx, col] = np.nan
        continue

    # Get LLM embedding
    llm_embedding = model.encode(llm_text)

    # Calculate cosine similarity
    similarity = cosine_similarity(
        therapist_embedding.reshape(1, -1),
        llm_embedding.reshape(1, -1)
    )[0][0]

    # Store result
    results.loc[idx, col] = similarity

# Step 5: Create visualizations
print("Creating visualizations...")

# 5.1: Average similarity bar chart
plt.figure(figsize=(10, 6))
avg_similarities = results.mean().sort_values(ascending=False)
avg_similarities.plot(kind='bar', color='skyblue')
plt.title('Average Similarity to Therapist Responses')
plt.ylabel('Cosine Similarity')
plt.xlabel('Model')
plt.xticks(rotation=45)
plt.tight_layout()
plt.savefig('avg_similarities.png')

# 5.2: Similarity distributions (boxplot)
plt.figure(figsize=(10, 6))
sns.boxplot(data=results)
plt.title('Distribution of Similarities to Therapist Responses')
plt.ylabel('Cosine Similarity')
plt.xlabel('Model')
plt.xticks(rotation=45)
plt.tight_layout()
plt.savefig('similarity_distributions.png')

# Step 6: Output results
print("\nAverage similarity scores:")
for model, score in avg_similarities.items():
    print(f"{model}: {score:.4f}")

# Save raw similarity scores to CSV
results.to_csv('similarity_scores.csv')
print("\nResults saved to 'similarity_scores.csv'")
print("Visualizations saved as 'avg_similarities.png' and 'similarity_distributions.png'")

```

```

"""# Empathy Scores"""

!pip install pandas numpy matplotlib seaborn sentence-transformers scikit-learn
transformers torch

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sentence_transformers import SentenceTransformer
from sklearn.metrics.pairwise import cosine_similarity
from transformers import AutoModelForSequenceClassification, AutoTokenizer
import torch
import re

# Step 1: Load the CSV data
print("Loading data...")
df = pd.read_csv("/content/Experiment Summary - Sheet2.csv")
print(f"Loaded {len(df)} conversations")

# Step 2: Clean text function
def clean_text(text):
    if not isinstance(text, str):
        return ""
    # Remove strange characters and normalize whitespace
    text = re.sub(r'[^\\w\\s.,!?\\'"]', ' ', text)
    text = re.sub(r'\\s+', ' ', text).strip()
    return text

# Step 3: Load models
print("Loading models...")
# Model for semantic similarity
similarity_model = SentenceTransformer('all-MiniLM-L6-v2')

# Model for empathy detection
# Using emotion-based model as a proxy for empathy detection
empathy_model_name = "bhadresh-savani/distilbert-base-uncased-emotion"
empathy_tokenizer = AutoTokenizer.from_pretrained(empathy_model_name)
empathy_model =
AutoModelForSequenceClassification.from_pretrained(empathy_model_name)

# Function to evaluate empathy
def evaluate_empathy(text):
    """Return empathy score based on emotional content"""
    if not text:
        return np.nan

    # Prepare inputs
    inputs = empathy_tokenizer(text, return_tensors="pt", truncation=True,
max_length=512)

```

```

# Get prediction
with torch.no_grad():
    outputs = empathy_model(**inputs)

# Get scores (emotion probabilities)
scores = torch.nn.functional.softmax(outputs.logits, dim=-1)

# Emotions in this model: ['anger', 'fear', 'joy', 'love', 'sadness',
'surprise']
# For empathy, we focus on 'love' (index 3) which is closest to empathy/caring
empathy_score = scores[0][3].item()

return empathy_score

# Step 4: Calculate similarities and empathy scores
print("Calculating similarities and empathy scores...")
# Get LLM column names (all columns after the therapist_response)
llm_columns = df.columns[2:].tolist()

# Create dataframes to store results
similarity_results = pd.DataFrame(index=df.index)
empathy_results = pd.DataFrame(index=df.index)

# First calculate empathy score for therapist responses
empathy_scores_therapist = []
for idx, row in df.iterrows():
    therapist_text = clean_text(row['therapist_response'])
    empathy_score = evaluate_empathy(therapist_text)
    empathy_scores_therapist.append(empathy_score)

# Add therapist empathy scores
empathy_results['therapist_response'] = empathy_scores_therapist

# Process each row for LLM responses
for idx, row in df.iterrows():
    # Get therapist response
    therapist_text = clean_text(row['therapist_response'])
    if not therapist_text:
        continue

    # Get therapist embedding
    therapist_embedding = similarity_model.encode(therapist_text)

    # Calculate similarity and empathy for each LLM response
    for col in llm_columns:
        llm_text = clean_text(row[col])
        if not llm_text:
            similarity_results.loc[idx, col] = np.nan
            empathy_results.loc[idx, col] = np.nan
            continue

```

```

# Calculate similarity
llm_embedding = similarity_model.encode(llm_text)
similarity = cosine_similarity(
    therapist_embedding.reshape(1, -1),
    llm_embedding.reshape(1, -1)
)[0][0]
similarity_results.loc[idx, col] = similarity

# Calculate empathy
empathy_score = evaluate_empathy(llm_text)
empathy_results.loc[idx, col] = empathy_score

# Step 5: Create visualizations
print("Creating visualizations...")

# 5.1: Average similarity bar chart
plt.figure(figsize=(10, 6))
avg_similarities = similarity_results.mean().sort_values(ascending=False)
avg_similarities.plot(kind='bar', color='skyblue')
plt.title('Average Similarity to Therapist Responses')
plt.ylabel('Cosine Similarity')
plt.xlabel('Model')
plt.xticks(rotation=45)
plt.tight_layout()
plt.savefig('avg_similarities.png')

# 5.2: Average empathy scores bar chart
plt.figure(figsize=(10, 6))
avg_empathy = empathy_results.mean().sort_values(ascending=False)
avg_empathy.plot(kind='bar', color='lightgreen')
plt.axhline(y=empathy_results['therapist_response'].mean(), color='r', linestyle='-',
            label=f'Therapist Average\n({empathy_results["therapist_response"].mean():.3f})')
plt.title('Average Empathy Scores')
plt.ylabel('Empathy Score')
plt.xlabel('Model')
plt.xticks(rotation=45)
plt.legend()
plt.tight_layout()
plt.savefig('avg_empathy.png')

# 5.3: Empathy vs Similarity scatter plot
plt.figure(figsize=(10, 8))
for col in llm_columns:
    avg_sim = similarity_results[col].mean()
    avg_emp = empathy_results[col].mean()
    plt.scatter(avg_sim, avg_emp, s=100, label=col)

plt.title('Empathy vs Similarity by Model')

```

```

plt.xlabel('Average Similarity to Therapist')
plt.ylabel('Average Empathy Score')
plt.grid(alpha=0.3)
plt.legend(bbox_to_anchor=(1.05, 1), loc='upper left')
plt.tight_layout()
plt.savefig('empathy_vs_similarity.png')

# Save raw scores to CSV
similarity_results.to_csv('similarity_scores.csv')
empathy_results.to_csv('empathy_scores.csv')

print("\nResults saved to CSV files")
print("Visualizations saved as PNG files")

"""# NEW"""

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sentence_transformers import SentenceTransformer
from sklearn.metrics.pairwise import cosine_similarity
from transformers import AutoModelForSequenceClassification, AutoTokenizer
import torch
import re

# Step 1: Load the CSV data
print("Loading data...")
df = pd.read_csv("/content/Experiment Summary - Sheet2.csv")
print(f"Loaded {len(df)} conversations")

# Step 2: Clean text function
def clean_text(text):
    if not isinstance(text, str):
        return ""
    # Remove strange characters and normalize whitespace
    text = re.sub(r'[^\\w\\s.,!?\\'"]-', ' ', text)
    text = re.sub(r'\\s+', ' ', text).strip()
    return text

# Step 3: Load models
print("Loading models...")
# Model for semantic similarity
similarity_model = SentenceTransformer('all-MiniLM-L6-v2')

# Model for empathy detection
empathy_model_name = "bhadresh-savani/distilbert-base-uncased-emotion"
empathy_tokenizer = AutoTokenizer.from_pretrained(empathy_model_name)
empathy_model =
AutoModelForSequenceClassification.from_pretrained(empathy_model_name)

```

```

# Function to evaluate empathy
def evaluate_empathy(text):
    """Return empathy score based on emotional content"""
    if not text:
        return np.nan

    # Prepare inputs
    inputs = empathy_tokenizer(text, return_tensors="pt", truncation=True,
max_length=512)

    # Get prediction
    with torch.no_grad():
        outputs = empathy_model(**inputs)

    # Get scores (emotion probabilities)
    scores = torch.nn.functional.softmax(outputs.logits, dim=-1)

    # Emotions in this model: ['anger', 'fear', 'joy', 'love', 'sadness',
'surprise']
    # For empathy, we focus on 'love' (index 3) which is closest to empathy/caring
    empathy_score = scores[0][3].item()

    return empathy_score

# Step 4: Calculate empathy scores for all responses (including therapist)
print("Calculating empathy scores...")
# Get LLM column names (all columns after the therapist_response)
llm_columns = df.columns[2:].tolist()
all_columns = ['therapist_response'] + llm_columns

# Create dataframe to store empathy results
empathy_results = pd.DataFrame(index=df.index)

# Calculate empathy for all responses
for col in all_columns:
    print(f"Processing empathy for {col}...")
    empathy_scores = []
    for idx, row in df.iterrows():
        text = clean_text(row[col])
        empathy_score = evaluate_empathy(text)
        empathy_scores.append(empathy_score)
    empathy_results[col] = empathy_scores

# Step 5: Calculate similarity between therapist and LLM responses
print("Calculating similarities...")
similarity_results = pd.DataFrame(index=df.index)

for idx, row in df.iterrows():
    # Get therapist response
    therapist_text = clean_text(row['therapist_response'])
    if not therapist_text:

```



```

        continue

# Get therapist embedding
therapist_embedding = similarity_model.encode(therapist_text)

# Calculate similarity for each LLM response
for col in llm_columns:
    llm_text = clean_text(row[col])
    if not llm_text:
        similarity_results.loc[idx, col] = np.nan
        continue

    # Calculate similarity
    llm_embedding = similarity_model.encode(llm_text)
    similarity = cosine_similarity(
        therapist_embedding.reshape(1, -1),
        llm_embedding.reshape(1, -1)
    )[0][0]
    similarity_results.loc[idx, col] = similarity

# Step 6: Calculate average scores
avg_empathy = empathy_results.mean().sort_values(ascending=False)
avg_similarity = similarity_results.mean().sort_values(ascending=False)

# Step 7: Create visualizations
print("Creating visualizations...")

# 7.1: Average similarity bar chart
plt.figure(figsize=(12, 6))
avg_similarity.plot(kind='bar', color='skyblue')
plt.title('Average Similarity to Therapist Responses')
plt.ylabel('Cosine Similarity')
plt.xlabel('Model')
plt.xticks(rotation=45)
plt.tight_layout()
plt.savefig('avg_similarities.png')

# 7.2: Average empathy scores bar chart
plt.figure(figsize=(12, 6))
avg_empathy.plot(kind='bar', color='lightgreen')
plt.title('Average Empathy Scores')
plt.ylabel('Empathy Score')
plt.xlabel('Response Source')
plt.xticks(rotation=45)
plt.tight_layout()
plt.savefig('avg_empathy.png')

# 7.3: Empathy scores comparison (therapist vs LLMs)
plt.figure(figsize=(12, 6))
therapist_avg = empathy_results['therapist_response'].mean()
llm_empathy_avgs = empathy_results[llm_columns].mean().sort_values(ascending=False)

```

```

# Create comparison dataframe
comparison_df = pd.DataFrame({
    'Empathy Score': [therapist_avg] + llm_empathy_avgs.tolist()
}, index=['Therapist'] + llm_empathy_avgs.index.tolist())

comparison_df.plot(kind='bar', color=['red'] + ['lightblue'] *
len(llm_empathy_avgs))
plt.title('Therapist vs LLM Empathy Scores')
plt.ylabel('Average Empathy Score')
plt.xlabel('Response Source')
plt.axhline(y=therapist_avg, color='red', linestyle='--', alpha=0.7,
            label=f'Therapist Average ({therapist_avg:.3f})')
plt.xticks(rotation=45)
plt.legend()
plt.tight_layout()
plt.savefig('therapist_vs_llm_empathy.png')

# 7.4: Empathy vs Similarity scatter plot
plt.figure(figsize=(10, 8))
for col in llm_columns:
    avg_sim = similarity_results[col].mean()
    avg_emp = empathy_results[col].mean()
    plt.scatter(avg_sim, avg_emp, s=100, label=col)

plt.title('Empathy vs Similarity by Model')
plt.xlabel('Average Similarity to Therapist')
plt.ylabel('Average Empathy Score')
plt.grid(alpha=0.3)
plt.legend(bbox_to_anchor=(1.05, 1), loc='upper left')
plt.tight_layout()
plt.savefig('empathy_vs_similarity.png')

# Save raw scores to CSV
similarity_results.to_csv('similarity_scores.csv')
empathy_results.to_csv('empathy_scores.csv')

# Step 8: Generate text output with average scores
print("\n=== SIMILARITY SCORES ===")
for model, score in avg_similarity.items():
    print(f"{model}: {score:.4f}")

print("\n=== EMPATHY SCORES ===")
for source, score in avg_empathy.items():
    print(f"{source}: {score:.4f}")

print("\n=== THERAPIST VS LLM EMPATHY COMPARISON ===")
print(f"Real therapist average empathy score: {therapist_avg:.4f}")
for model, score in llm_empathy_avgs.items():
    diff = score - therapist_avg
    diff_percent = (diff / therapist_avg) * 100

```

```

    print(f"{model}: {score:.4f} ({diff:+.4f}, {diff_percent:+.2f}% compared to
therapist)")

# Create a text file with the results
with open('evaluation_results.txt', 'w') as f:
    f.write("=== EVALUATION RESULTS ===\n\n")

    f.write("=== SIMILARITY SCORES ===\n")
    for model, score in avg_similarity.items():
        f.write(f"{model}: {score:.4f}\n")

    f.write("\n=== EMPATHY SCORES ===\n")
    for source, score in avg_empathy.items():
        f.write(f"{source}: {score:.4f}\n")

    f.write("\n=== THERAPIST VS LLM EMPATHY COMPARISON ===\n")
    f.write(f"Real therapist average empathy score: {therapist_avg:.4f}\n")
    for model, score in llm_empathy_avgs.items():
        diff = score - therapist_avg
        diff_percent = (diff / therapist_avg) * 100
        f.write(f"{model}: {score:.4f} ({diff:+.4f}, {diff_percent:+.2f}% compared
to therapist)\n")

print("\nDetailed results saved to 'evaluation_results.txt'")

"""# Model Benchmarking"""

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sentence_transformers import SentenceTransformer
from sklearn.metrics.pairwise import cosine_similarity
from transformers import AutoModelForSequenceClassification, AutoTokenizer
import torch
import re

# Step 1: Load the CSV data
print("Loading data...")
df = pd.read_csv("/content/Model Benchmarking.csv")
print(f"Loaded {len(df)} conversations")

# Step 2: Clean text function
def clean_text(text):
    if not isinstance(text, str):
        return ""
    # Remove strange characters and normalize whitespace
    text = re.sub(r'[^\\w\\s.,!?\\'"]-', ' ', text)
    text = re.sub(r'\\s+', ' ', text).strip()
    return text

```

```

# Step 3: Load models
print("Loading models...")
# Model for semantic similarity
similarity_model = SentenceTransformer('all-MiniLM-L6-v2')

# Model for empathy detection
empathy_model_name = "bhadresh-savani/distilbert-base-uncased-emotion"
empathy_tokenizer = AutoTokenizer.from_pretrained(empathy_model_name)
empathy_model =
AutoModelForSequenceClassification.from_pretrained(empathy_model_name)

# Function to evaluate empathy
def evaluate_empathy(text):
    """Return empathy score based on emotional content"""
    if not text:
        return np.nan

    # Prepare inputs
    inputs = empathy_tokenizer(text, return_tensors="pt", truncation=True,
max_length=512)

    # Get prediction
    with torch.no_grad():
        outputs = empathy_model(**inputs)

    # Get scores (emotion probabilities)
    scores = torch.nn.functional.softmax(outputs.logits, dim=-1)

    # Emotions in this model: ['anger', 'fear', 'joy', 'love', 'sadness',
'surprise']
    # For empathy, we focus on 'love' (index 3) which is closest to empathy/caring
    empathy_score = scores[0][3].item()

    return empathy_score

# Step 4: Calculate empathy scores for all responses (including therapist)
print("Calculating empathy scores...")
# Get model column names (all columns after the therapist_response)
model_columns = df.columns[2:].tolist()
all_columns = ['therapist_response'] + model_columns

# Create dataframe to store empathy results
empathy_results = pd.DataFrame(index=df.index)

# Calculate empathy for all responses
for col in all_columns:
    print(f"Processing empathy for {col}...")
    empathy_scores = []
    for idx, row in df.iterrows():
        text = clean_text(row[col])
        empathy_score = evaluate_empathy(text)

```

```

        empathy_scores.append(empathy_score)
    empathy_results[col] = empathy_scores

# Step 5: Calculate similarity between therapist and model responses
print("Calculating similarities...")
similarity_results = pd.DataFrame(index=df.index)

for idx, row in df.iterrows():
    # Get therapist response
    therapist_text = clean_text(row['therapist_response'])
    if not therapist_text:
        continue

    # Get therapist embedding
    therapist_embedding = similarity_model.encode(therapist_text)

    # Calculate similarity for each model response
    for col in model_columns:
        model_text = clean_text(row[col])
        if not model_text:
            similarity_results.loc[idx, col] = np.nan
            continue

        # Calculate similarity
        model_embedding = similarity_model.encode(model_text)
        similarity = cosine_similarity(
            therapist_embedding.reshape(1, -1),
            model_embedding.reshape(1, -1)
        )[0][0]
        similarity_results.loc[idx, col] = similarity

# Step 6: Calculate average scores
avg_empathy = empathy_results.mean().sort_values(ascending=False)
avg_similarity = similarity_results.mean().sort_values(ascending=False)

# Step 7: Create visualizations
print("Creating visualizations...")

# 7.1: Highlight MirrorAI in the similarity chart
plt.figure(figsize=(12, 6))
colors = ['red' if x == 'MirrorAI' else 'skyblue' for x in avg_similarity.index]
avg_similarity.plot(kind='bar', color=colors)
plt.title('Average Similarity to Therapist Responses')
plt.ylabel('Cosine Similarity')
plt.xlabel('Model')
plt.xticks(rotation=45)
plt.tight_layout()
plt.savefig('model_similarity.png')

# 7.2: Empathy scores comparison (therapist vs models)
plt.figure(figsize=(12, 6))

```

```

therapist_avg = empathy_results['therapist_response'].mean()

# Use different colors to highlight MirrorAI and therapist
colors = []
for x in avg_empathy.index:
    if x == 'therapist_response':
        colors.append('darkgreen')
    elif x == 'MirrorAI':
        colors.append('red')
    else:
        colors.append('lightblue')

avg_empathy.plot(kind='bar', color=colors)
plt.title('Average Empathy Scores by Model')
plt.ylabel('Empathy Score')
plt.xlabel('Response Source')
plt.axhline(y=therapist_avg, color='darkgreen', linestyle='--', alpha=0.7,
            label=f'Therapist Average ({therapist_avg:.3f})')
plt.xticks(rotation=45)
plt.legend()
plt.tight_layout()
plt.savefig('model_empathy.png')

# 7.3: Empathy vs Similarity scatter plot
plt.figure(figsize=(10, 8))
# Add therapist point for reference (100% similarity with itself)
plt.scatter(1.0, empathy_results['therapist_response'].mean(),
            s=150, color='darkgreen', label='Therapist')

# Add model points
for col in model_columns:
    avg_sim = similarity_results[col].mean()
    avg_emp = empathy_results[col].mean()
    color = 'red' if col == 'MirrorAI' else 'blue'
    plt.scatter(avg_sim, avg_emp, s=100, color=color, label=col)

plt.title('Empathy vs Similarity by Model')
plt.xlabel('Average Similarity to Therapist')
plt.ylabel('Average Empathy Score')
plt.grid(alpha=0.3)
plt.legend(bbox_to_anchor=(1.05, 1), loc='upper left')
plt.tight_layout()
plt.savefig('empathy_vs_similarity.png')

# Save raw scores to CSV
similarity_results.to_csv('model_similarity_scores.csv')
empathy_results.to_csv('model_empathy_scores.csv')

# Step 8: Create comparison dataframe
comparison_df = pd.DataFrame({
    'Similarity': [similarity_results[col].mean() for col in model_columns],

```

```

        'Empathy': [empathy_results[col].mean() for col in model_columns],
    })
comparison_df.index = model_columns

# Add therapist empathy score for reference
therapist_empathy = empathy_results['therapist_response'].mean()

# Step 9: Generate text output with average scores and rankings
print("\n=== SIMILARITY SCORES ===")
for model, score in avg_similarity.items():
    print(f"{model}: {score:.4f}")

print("\n=== EMPATHY SCORES ===")
for source, score in avg_empathy.items():
    print(f"{source}: {score:.4f}")

print("\n=== MODEL VS THERAPIST EMPATHY COMPARISON ===")
print(f"Real therapist average empathy score: {therapist_avg:.4f}")
for model in model_columns:
    diff = empathy_results[model].mean() - therapist_avg
    diff_percent = (diff / therapist_avg) * 100
    print(f"{model}: {empathy_results[model].mean():.4f} ({diff:+.4f},
{diff_percent:+.2f}% compared to therapist)")

# Create a text file with the results
with open('model_evaluation_results.txt', 'w') as f:
    f.write("=== MODEL EVALUATION RESULTS ===\n\n")

    f.write("=== SIMILARITY SCORES (higher is better) ===\n")
    for model, score in avg_similarity.items():
        f.write(f"{model}: {score:.4f}\n")

    f.write("\n=== EMPATHY SCORES ===\n")
    for source, score in avg_empathy.items():
        f.write(f"{source}: {score:.4f}\n")

    f.write("\n=== MODEL VS THERAPIST EMPATHY COMPARISON ===\n")
    f.write(f"Real therapist average empathy score: {therapist_avg:.4f}\n")
    for model in model_columns:
        diff = empathy_results[model].mean() - therapist_avg
        diff_percent = (diff / therapist_avg) * 100
        f.write(f"{model}: {empathy_results[model].mean():.4f} ({diff:+.4f},
{diff_percent:+.2f}% compared to therapist)\n")

    # Add overall ranking based on combined metrics
    f.write("\n=== OVERALL MODEL RANKING ===\n")
    # Normalize similarity and empathy scores (0-1 scale)
    normalized_similarity = (comparison_df['Similarity'] -
comparison_df['Similarity'].min()) / (comparison_df['Similarity'].max() -
comparison_df['Similarity'].min())

```

```

    normalized_empathy = (comparison_df['Empathy'] -
comparison_df['Empathy'].min()) / (comparison_df['Empathy'].max() -
comparison_df['Empathy'].min())

    # Calculate combined score (equal weighting)
    comparison_df['Combined Score'] = (normalized_similarity + normalized_empathy)
/ 2

    # Sort by combined score
    ranked_models = comparison_df.sort_values('Combined Score', ascending=False)

    for rank, (model, row) in enumerate(ranked_models.iterrows(), 1):
        f.write(f'{rank}. {model} - Similarity: {row['Similarity']:.4f}, Empathy:
{row['Empathy']:.4f}, Combined: {row['Combined Score']:.4f}\n')

    # MirrorAI specific analysis
    f.write("\n=== MIRRORAI SPECIFIC ANALYSIS ===\n")
    mirror_rank = ranked_models.index.get_loc('MirrorAI') + 1
    mirror_similarity = similarity_results['MirrorAI'].mean()
    mirror_empathy = empathy_results['MirrorAI'].mean()

    f.write(f"MirrorAI Rank: {mirror_rank} out of {len(model_columns)} models\n")
    f.write(f"Similarity to therapist: {mirror_similarity:.4f}\n")
    f.write(f"Empathy score: {mirror_empathy:.4f}\n")

    # Compare to average of other models
    other_models = [m for m in model_columns if m != 'MirrorAI']
    avg_others_similarity = similarity_results[other_models].mean().mean()
    avg_others_empathy = empathy_results[other_models].mean().mean()

    sim_diff = (mirror_similarity - avg_others_similarity) / avg_others_similarity
* 100
    emp_diff = (mirror_empathy - avg_others_empathy) / avg_others_empathy * 100

    f.write(f"Similarity comparison: {sim_diff:+.2f}% versus average of other
models\n")
    f.write(f"Empathy comparison: {emp_diff:+.2f}% versus average of other
models\n")

print("\nDetailed results saved to 'model_evaluation_results.txt'")

```