

Evaluating AI-driven Psychotherapy: Insights from Large Language Models and Human Expert Comparisons

Anonymous ACL submission

Abstract

The integration of Large Language Models (LLMs), such as GPT-4, has shown great promise in mental health applications for initial assessments based on user-reported symptoms. Traditional assessments often involve subjective evaluations by professional psychologists, leading to inconsistent reproducibility across datasets. To address this, we developed a comprehensive evaluation framework using entropy analysis, keyword frequency analysis, and Latent Dirichlet Allocation (LDA) to quantitatively assess LLM outputs. Our results indicate that LLMs can effectively identify and engage with a range of treatment topics and provide a broader range of treatment opinions than human psychologists. However, LLMs lack depth in their responses, the recommendation generated by LLMs tends to using generalized word instead of using professional words. This study explores the feasibility of LLMs as virtual psychotherapists, highlights their shortcomings in depth, and proposes improved methods for evaluating large model responses. This research provides valuable insights into the potential and challenges of integrating LLMs into mental health practices, paving the way for future research to enhance the effectiveness and reliability of AI-driven therapeutic solutions.

1 Introduction

Psychotherapy, a therapeutic interaction or treatment between a trained professional and a client aimed at addressing psychological issues and improving mental health, is a fundamental component of the mental health cycle. It applies multiple non-invasive methodologies to address psychological problems. Psychotherapy is also considered as a secondary methodology to prevent the recurrence of certain conditions and is often utilized to manage urgent cases of depression (Karrouri et al., 2021). In psychotherapy field, Cognitive Behavioral Therapy (CBT) is recognized to be crucial

in addressing anxiety. This enhances the key position of CBT in helping patients with depression and highlights its importance as both a preventive technique and treatment methodology (Bandelow et al., 2017). Additionally, Non-Directive Support Therapy (NDST) has been applied in psychotherapy treatment methodologies. It provides emotional support and energy for patients in self-exploration and self-development to solve their problems. One research suggests that, compared with traditional methodologies, this psychotherapy approach showed better treatment results in the short term (Cuijpers et al., 2014). Additionally, invasive psychological treatment methodologies have been proven to have similar effectiveness to depression medication treatment during the urgent treatment stage (de Maat et al., 2007). This enhances the key role and benefits of using psychotherapy in treating mental health disorders and managing overall mental health.

The integration of Artificial Intelligence (AI) technologies, such as GPT-4 and other Large Language Models (LLMs), has driven the development of intelligent psychotherapy applications. The primary goal of researchers and institutions is to provide timely and effective treatment recommendations for medical professionals and individuals seeking treatment (Chen et al., 2023; Montagna et al., 2023). Although LLMs have demonstrated the strong ability to analyze natural language and provide diverse feedback quickly (Singhal et al., 2023), the efficiency and reliability of these AI-powered psychotherapy tools in providing accurate diagnoses and recommending effective treatments still remains controversial (Manríquez Roa et al., 2021). This is mainly due to the complexity of the medical field that requires large language models to have the ability to understand the medical context, find appropriate medical knowledge, and reason using authoritative information and clues provided by patients, and this complexity in the medical

field have led to a variety of potential treatments (Singhal et al., 2023). Therefore, assessing the performance of AI-based virtual psychotherapists in the depth and coverage of their therapeutic advice, especially when compared to human professionals, has become a key focus of research.

This study aims to evaluate the effectiveness of LLM-based chatbots in recommending treatment suggestions and their consistency with those proposed by psychotherapists and the depth of the protocol.

- Our (H_0) is that there is no significant difference in the diagnosis and treatment opinions provided by psychologists and LLMs overall, i.e., the quality of the output of the large model is broadly consistent with that of psychologists.

To test this hypothesis, we introduce a novel evaluation framework that applies case studies from the American Psychological Association (APA) as a benchmark to detect the differences between LLMs output and case scenarios through LDA modeling and entropy analysis, so as to comprehensively evaluate their application in the field of psychology.

Our contributions are:

- We propose and implement a framework combining entropy analysis, keyword frequency analysis, and the novel Latent Dirichlet Allocation (LDA) to evaluate the diversity, depth, and applicability of LLMs in generating psychological diagnoses and treatment recommendations. This provides a quantitative way to measure the feasibility of LLM applications in clinical settings and offers a new perspective on evaluating LLM technology in mental health diagnosis and treatment planning.
- Through detailed comparisons and in-depth analysis, we evaluated the differences between LLM-generated treatment recommendations and those made by human psychologists. Our findings suggest that LLM recommendations often lack the detail found in human expert recommendations, highlighting both the strengths and shortcomings of LLMs in generating psychotherapeutic recommendations and providing a balanced perspective on integrating LLM techniques with psychotherapy practice.

- We demonstrate the potential impact of LLMs in increasing access to mental health care by validating their ability to provide mental health related diagnoses and treatment recommendations in evaluating its diversity and depth. Our study highlights the potential for LLMs platforms to improve the accessibility and scalability of psychotherapy services, especially in resource-limited or remote areas. Additionally, we initiate discussions on ethical, practical, and strategic planning considerations to maximize the benefits of AI in mental health practices.

Through rigorous evaluation and comparative analysis utilizing novel Latent Dirichlet Allocation (LDA) modeling, word frequency analysis, and a series of statistical analyses to evaluate the treatment recommendation capability, diversity, and depth based on 10 professional case studies from the American Psychological Association and generated by LLMs, this study highlights the potential and limitations of LLMs in the diagnosis and treatment of mental health, and provides valuable insights and directions for future research and application in this field.

2 Related Work

In recent years, with the advancement of natural language processing technology, Large Language Models (LLMs) like GPT-4 have been widely studied in the field of primary consultation and support in the health field. Michimasa et al. 2024 demonstrated in their experiments that LLMs can exhibit a level of professionalism similar to that of psychologists, with no high-risk, aggressive, or discriminatory responses found in conversations with GPT-4. In addition, Luoma Ke (Ke et al., 2024)'s study also confirmed that LLMs, as a preliminary diagnostic tool in clinical and counseling psychology, can quickly identify potential mental health problems in users, such as depression and anxiety. John (Ayers et al., 2023) evaluated responses from physicians and LLMs, with the results that raters favoring responses from LLMs, and the quality of LLMs outpacing physician responses.

However, although LLMs have generally received neutral and positive feedback in past research evaluations, they exhibit a range of problems. Critics, such as Topol (Meskó and Topol, 2023), point out that the recommendations generated by LLMs were not very reliable because the

182 data used by the large model did not come from a
183 formal bedside conversation, and that the responses
184 from the large model may involve fictitious sources.
185 According to a survey, out of 157 participants, 123
186 used ChatGPT for health queries. Besides, 83 peo-
187 ple believed that the treatment recommendations
188 provided by the large model are more accurate than
189 those provided by traditional online communities.
190 While the study found that people prefer to use
191 LLMs for health consultations, the researchers also
192 expressed concerns that the databases of LLMs
193 need to be updated in a timely manner to ensure the
194 accuracy and reliability of their information. (Xiao
195 et al., 2024).

196 At the same time, Natural Language Process-
197 ing (NLP) technology has been widely used in text
198 analysis tasks, and NLP methods have also shown
199 significant value in the psychological field. For ex-
200 ample, text analysis methods such as Pearson cor-
201 relation coefficients and sentiment analysis have
202 been used to assess the consistency of machine-
203 generated responses with human expert recommen-
204 dations (Danna et al., 2024). In addition, NLP
205 techniques such as TF-IDF and Word2Vec have
206 been applied to data classification for the assess-
207 ment of suicidality (Aldhyani et al., 2022). While
208 these techniques excel in dataset processing, they
209 have traditionally been used primarily for data clas-
210 sification in deep learning, or to predict suicidality
211 and mental illness by analyzing online social me-
212 dia comments. Existing studies have not focused
213 on the application of these methods in assessing
214 the output quality generated by LLMs, revealing
215 potential research gaps and development directions
216 in this field.

217 With the development of AI, especially the inte-
218 gration of LLMs in mental health, finding a way
219 to assess the quality of the output of these mod-
220 els have become particularly urgent (Elyoseph and
221 Levkovich, 2024). The benefit of quantifying the
222 output of LLMs is that it can provide an objective
223 way to evaluate the effectiveness and reliability
224 of these models in real-world applications. Re-
225 search has shown that while LLMs can deal with a
226 wide range of topics, they often lack the depth pro-
227 vided by human experts, a problem that may stem
228 from the phenomenon of knowledge duplication
229 in LLMs (Chen et al., 2023). Therefore, there is
230 a need to explore and establish a new assessment
231 framework to comprehensively assess the capacity
232 of LLMs in terms of mental health diagnosis and
233 treatment recommendations. Such a framework

234 can not only help identify and address the short-
235 comings of LLMs in specific applications, but also
236 facilitate a more effective fusion of AI and human
237 expertise.

238 One methodology can be considered in the
239 framework is the cosine similarity, which can be
240 used to compare similarity of the text written by
241 psychotherapist and LLMs generated text. Cosine
242 similarity is a vector space modeling technique
243 used to quantify the similarity between two docu-
244 ments (Januzaj and Luma, 2022), making it a key
245 tool for text analysis and comparison. This metric
246 calculates the cosine value of the angle between
247 two vectors, representing the position of the text in
248 a multidimensional space, to determine their simi-
249 larity. It has a wide range of applications, especially
250 in the evaluation of text consistency and relevance
251 in automated systems. In the Automated Essay
252 Scoring (AES) system, cosine similarity plays an
253 important role by comparing the text submitted by
254 students with the documents written by experts. By
255 using this method in conjunction with weighted
256 terminology analysis, the AES system achieves
257 a meticulous assessment of textual consistency,
258 demonstrating the effectiveness of the method in
259 an educational setting (Lahitani et al., 2016). In
260 addition, the field of psychology also employs co-
261 sine similarity for diagnostic purposes, facilitating
262 the comparison of the symptoms provided by the
263 patient with the established psychological profile
264 during a virtual consultation. This innovative ap-
265 plication helps doctors reduce their workload as a
266 diagnostic aid by analyzing the user's text input to
267 make a preliminary diagnosis of a patient's mental
268 health (Bhattacharya and Pissurlenkar, 2023).

269 However, when the lengths of the two inputs are
270 different, the output generated by the cosine simi-
271 larity method will be significantly affected, which is
272 not accurate for evaluating the LLMs response and
273 case studies of text of different lengths. Therefore,
274 we introduce entropy analysis to more effectively
275 evaluate the complexity of the results generated by
276 LLMs. Entropy is a measurement derived from
277 information theory that measures uncertainty and
278 randomness within a system. A study using en-
279 tropy to measure the consistency and diversity of
280 Key Audit Matters (KAMs) disclosures in audit re-
281 ports showed that monitoring the entropy of KAMs
282 disclosures can reveal trends and consistency in the
283 evolution of audit practices over time (Lin, 2023).
284 This study suggests that we can evaluate the per-
285 formance of LLMs by measuring the topic distribu-

tion of entropy and further analyze the diversity of LLMs-generated topics.

In our study, we aim to critically assess the effectiveness of LLMs in performing tasks similar to those of virtual psychologists by using APIs such as ChatGPT, as well as mainstream NLP tools, including LDA and entropy analysis.

3 Methods

3.1 Data Source Selection

Our research methodology starts with selecting the appropriate dataset to make the evaluation. We chose a series of formatted case study from APA instead of using non-structural dataset like DAIC-WOZ from USC (Burdisso et al., 2024), mainly because the structured format of the APA is more in line with the capabilities of LLMs. We initially used USC’s DAIC-WOZ dataset, but found that ChatGPT could not track the transcription format of Q&A correspondingly when processing this type of transcription’s data without manually intervened. While we found that manual intervention allowed ChatGPT to follow the Q&A transcription format in the dataset, this intervention method was shown to lead to later human intervention bias in LLM answers, resulting in inaccurate research results (Loya et al., 2023). In contrast, the highly well-formatted APA case studies provide a diverse and comprehensive mental health scenario, and this structured format is more suitable for assessing the diagnostic and treatment recommendation capabilities of LLMs than the DAIC-WOZ dataset. In addition, APA has been mentioned in many psychology research papers and is considered as one of the most authoritative sources of psychological research data (Badr et al., 2023; Sheridan and Carr, 2018).

In our study, we selected 10 case studies from the American Psychological Association (APA), including cases of individuals with depression and Post-traumatic Stress Disorder (PTSD). These cases include the patient’s background, diagnosis, and corresponding treatment plan. All personal information has been anonymized by the APA. The cases cover a diverse range of genders and ages, ensuring a comprehensive evaluation of the treatment recommendations provided by LLMs.

3.2 Entropy Analysis for Topic Distribution

In our study, we used entropy analysis to assess how LLMs divided their attention across different

psychotherapy topics and compared it to human psychologists. Through entropy analysis, we can determine whether the text generated by LLMs is concise or diverse with multiple topics. In order to ensure fair comparison, we have normalized the topic probabilities in the document, and the normalization calculation is as follows:

$$p(t_i) = \frac{n_{t_i}}{\sum_j n_{t_j}}$$

Here, n_{t_i} represents the count of words associated with topic t_i within a document, and $\sum_j n_{t_j}$ is the total word count across all topics in that document. This ensures that the sum of probabilities across topics equals one, facilitating a meaningful entropy calculation.

The entropy for each document’s topic distribution was then computed using the formula:

$$H(T) = - \sum_{i=1}^K p(t_i) \log_2 p(t_i)$$

This equation, where K is the number of topics and $p(t_i)$ denotes the probability of each topic, utilizes the logarithm base 2 to measure entropy in bits, enhancing our understanding of topic distribution’s evenness.

3.3 Prompt Design

In this study, a specific prompt was designed for the LLMs to ensure consistency in the responses across different models. This prompt incorporates a curated list of keywords that are closely related to mental health treatment, ensuring that the treatment recommendations generated are relevant and based on well-established psychological principles.

- **Diagnosis Section:** The prompt includes keywords such as *anxiety*, *depression*, and *panic attacks*. These terms are selected to guide the LLMs to focus on common psychological conditions, facilitating a targeted exploration of potential diagnoses.
- **Treatment Plan Section:** Keywords like *Cognitive Behavioral Therapy (CBT)*, *psychodynamic therapy*, and *humanistic therapy* are included. These therapies represent a range of approaches in psychotherapy, allowing the LLMs to generate diverse and comprehensive treatment plans.

This methodical selection of keywords is informed by recent advancements in AI applications within healthcare, where patients can utilize an LLMs to input relevant keywords or questions, thus accessing a wealth of medical knowledge (Pagad et al., 2022). We used this idea to design the prompt to let LLMs' output become consistent and relevant. The complete prompt utilized in our evaluations is detailed in Appendix A.

3.4 Word Frequency Analysis

We used word frequency analysis to assess the similarity between the treatments described in the APA case study and those generated by LLM. Our study built on the potential LDA topic modeling of Blei (Blei, 2003) and extends the application of natural language processing (NLP) techniques in mental health research outlined by Miner (Miner et al., 2020). We aimed to compare the differences in treatment recommendations between the results generated by LLMs and the demonstration results in the APA case study by quantifying treatment-specific terms in text data. Besides, since one research done by Torous and Keshavan (Torous and Keshavan, 2020) highlights the importance of evaluating digital tools to ensure that these tools meet clinical standards and effectively enhance patient care in the field of mental health. Our another focus in our quantitative assessment framework is the analysis of treatment-related word frequency comparisons between APA case studies and LLM outputs. We wanted to use this approach to assess whether the LLMs was able to generate broader clinical recommendations, while retaining some depth of therapeutic insight. Through this exploration, we aim to uncover the potential of LLMs as a tool for mental health practitioners and the performance of LLMs in the professional field.

3.5 Comparative Analysis Using Latent Dirichlet Allocation (LDA)

3.5.1 Objective of Using LDA

In order to provide a detailed analysis and comparison of the treatment recommendations provided by ChatGPT with those described in (APA) case study, we used the LDA as another important part of our evaluation framework for LLMs. LDA was chosen as our methodological tool based on its effectiveness in identifying potential topics in a large corpus of text, as demonstrated by the groundbreaking study (Hagg et al., 2022; Kotenko et al., 2021). As a result, the application of LDA enables a detailed

and structured comparative analysis, with a particular focus on the thematic differences between the responses generated by ChatGPT and the treatment recommendations described in the case study. This approach allows us to understand ChatGPT's capabilities and limitations in psychotherapy related tasks. Through this analytical perspective, we aim to critically assess the similarity of ChatGPT recommendations with contemporary treatment standards in evaluating the diversity and depth of the responses, thereby contributing to an ongoing conversation about the integration of AI in clinical settings.

3.6 Summary of Analytical Procedure

The comparative analysis is based on a two-stage approach, distilling and examining the essence of the topic through the LDA model of APA case studies and ChatGPT-generated recommendations. Before LDA was applied, extensive text data pre-processing was performed, including tokenization, stop word removal, and invalid word reduction, to optimize the text's topic extraction.

The analysis process is as follows: First, the Subject Heading Distribution Analysis involves identifying and visualizing the most important words within the topics extracted from APA case studies and ChatGPT outputs. By examining word distribution, the main thematic focus of each source is elucidated, thereby assessing the consistency and differences in treatment topics. Next, the Document-topic ratio assessment quantifies the representation of each topic in a single document, facilitating a fine-grained comparison of topic prevalence between the original case study and ChatGPT recommendations. This stage uses heat map visualization to display the topic distribution pattern, highlighting the similarities and differences in theme emphasis. Following this, the Compare Topic-Word Relationship Exploration uses a heat map to further dissect the relationship between key terms and their related topics in the two datasets. This step is essential for assessing the depth and specificity of ChatGPT's treatment recommendations relative to the established treatment modalities documented in the APA case study. Finally, the Entropy-based variability assessment employs entropy measurements to assess the variability of topic distributions in LLMs and artificially generated text. This analysis quantifies the diversity of topics covered by each source, providing insights into the comprehensiveness of treatment recommendations and concerns.

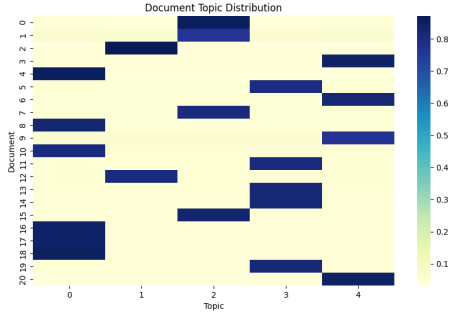


Figure 3: Word frequency heatmap analysis of treatment plans generated from LLMs by using APA study cases

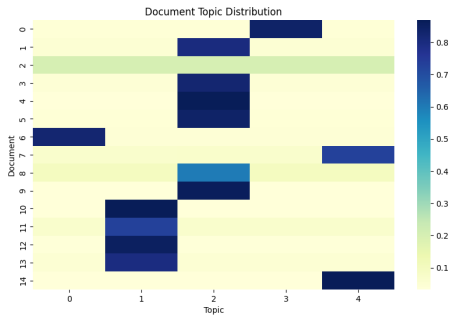


Figure 4: Word frequency heatmap analysis of treatment plans based on APA study cases

Figure 7 presents a box-plot comparing the entropy values of original text and LLM-generated text across different case studies, providing a visual representation of the data.

4.2.1 Mann-Whitney U test

To statistically evaluate the difference in entropy between the original text and the LLM-generated text, we also performed the Mann-Whitney U test, which is a non-parametric test suitable for comparing differences between two independent samples. The results are shown in the table 1 and show that there is no statistically significant difference between the two sets of text.

Table 1: Mann-Whitney U test result

measurement	value
U statistics	31.0
P value	0.162

In our study, the Null Hypothesis (H_0) is that there is no significant difference in treatment recommendations between large language models (LLMs) and case studies overall. The Mann-



Figure 5: Entropy heat-map analysis of treatment plans generated from LLMs by using APA study cases



Figure 6: Entropy heat-map analysis of treatment plans based on APA study cases

Whitney U test of entropy showed a U statistic of 31.0 and a P-value of 0.162, suggesting that the difference in topic distribution between the text generated by LLMs and the text generated by human psychologists was not statistically significant, which supported our (H_0) that LLMs and human psychologists' recommendations were similar in diversity and uniformity.

5 Results

5.1 Interpretation of Topic-Word Frequency Analysis

The LDA analysis of the case studies uncovered a diverse range of topics associated with PTSD and depression, including treatment methods, patient living environments, and social factors such as school and social circles. These topics reveal

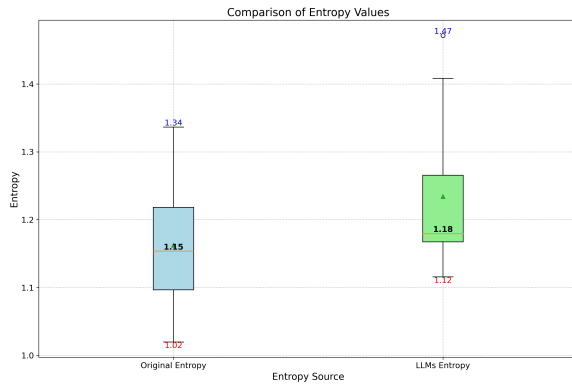


Figure 7: Comparative Entropy Analysis of LLMs-Generated and Expert-Designed Treatment Plans in Psychological Case Studies

both consistency and differentiation between the treatment recommendations generated by LLMs and those provided by human experts. While LLMs effectively identified general treatment topics like "medications" and "symptoms," they often lacked the depth and specificity evident in human-generated recommendations. For instance, human psychologists frequently mentioned specific therapies such as CBT, whereas LLMs tended to use broader terms like "treatment."

The document-topic distribution analysis further highlighted significant differences in the depth of engagement between LLMs and human experts. Human psychologists provided detailed and more professional terms, such as "CBT," whereas LLM-generated responses were more general. This suggests that, while LLMs can cover a wide range of relevant topics, they do not engage with the same level of depth, complexity, and detail as human experts. These findings align with the entropy analysis results, reinforcing the ongoing disparity in professionalism between large language models and human experts.

Based on the analysis of Figures 5 and 6, distinct differences were observed in the topic-word associations and entropy values between the treatment recommendations provided by human psychologists and those generated by LLMs. The heatmaps demonstrate that human-generated texts have more concentrated keyword relevance within specific topics, resulting in lower entropy values and indicating a more focused and detailed discussion. In contrast, LLM-generated texts display a broader but less focused distribution of keywords, leading to higher entropy values. This dispersion suggests that LLMs cover a wider array of topics but with less depth and

specificity. For example, in Topic 3, the keywords in LLM-generated texts are more evenly spread across terms like "academic," "week," and "initially," reflecting a general approach rather than a detailed examination.

To statistically validate these findings, we tested the Null Hypothesis (H_0) that there is no significant difference in treatment recommendations between large language models (LLMs) and human experts. Using the Mann-Whitney U test on entropy values, we obtained a U statistic of 31.0 with a p-value of 0.162. This supports the null hypothesis, indicating no significant difference in the overall uniformity of topic distribution between LLMs and human experts.

Overall, the heatmap and entropy analysis highlight the need for further refinement of LLMs to enhance their ability to provide detailed and specific treatment recommendations, aiming to achieve a balance and depth similar to that of human psychotherapists. These observations underscore the ongoing need to improve LLMs for more effective therapeutic applications.

6 Conclusion

In this study, we have conducted a comprehensive analysis comparing Large Language Models (LLMs) with human psychologists in providing treatment recommendations for depression. Employing Latent Dirichlet Allocation (LDA) and entropy analysis, we found that while LLMs exhibit comparable diversity and uniformity in generating treatment recommendations, they lack the specificity and depth of human experts. LLMs effectively cover a wide range of topics but do not engage with the nuanced details that characterize human-generated recommendations. Despite this, the uniformity and diversity in LLM-generated recommendations suggest significant potential for their application in mental health care. However, further improvements are necessary to ensure consistent and in-depth performance across therapeutic scenarios. This study provides valuable insights into the potential and challenges of integrating LLMs into mental health practices while providing a new methodology in evaluating text-based LLMs generated response, paving the way for future research to enhance the effectiveness and reliability of AI-driven therapeutic solutions.

7 Limitations

Although this study provides important insights into AI-based approaches to the comparison of virtual psychotherapists with human professionals, there are several limitations to be concerned about. First, the generalizability of our results may be limited due to the specificity of the case studies used from APA sources and the datasets on which LLMs were trained, and the conclusions of the study are not broadly representative due to the small amount of data. In addition, since we did not collect the latest study cases, and the LLMs was trained based on massive amounts of data, this may lead to the limitation of our analysis conclusions due to the fact that some of the treatments in our study cases are not the latest mainstream treatments. At the same time, our analysis relied on textual data through LDA models, which also limited our ability to consider non-verbal cues and clinical intuitions inherent in human treatment. In addition, the ability of large model systems to interpret complex human emotions and clinical contexts remains lacking, which may affect the depth of therapeutic interventions recommended by these systems. Ethical issues regarding privacy and data sensitivity, as well as the enormous computational demands for deploying LLMs in clinical settings, also pose significant challenges.

It is worth noting that in some cases during our LDA analysis, the Entropy of the subject distribution (Entropy) appeared to be greater than 1. This may be due to outliers in the data preprocessing steps (such as word frequency calculation, TF-IDF transformation, etc.), which further affects the output of the model. These outliers may be amplified in subsequent processing steps, resulting in a probability value greater than 1. Although these numerical errors usually do not significantly affect the overall performance of the model and the final analysis results, the handling of these anomalies needs to be considered in further research.

Moreover, LLMs are highly dependent on the variety and richness of the input data they are trained on. In cases where training data lack demographic diversity or contain biased information, this can lead to skewed or biased AI-generated recommendations and diagnoses. Therefore, while LLMs can significantly expand access to mental health services, the underlying biases in training datasets can limit the appropriateness and effectiveness of the recommendations provided, especially for under-

represented groups. Addressing these data biases is essential to ensure equitable mental health support across diverse populations. This aspect highlights the need for continual updates and the urgency of having a professional clinical dataset to mitigate biases and improve the accuracy and fairness of AI-driven mental health interventions. Further studies should focus on developing robust methods for continuous data validation and enhancement, introducing more comprehensive textual data analysis methods based on quantitative approaches, as well as the implementing comprehensive ethical frameworks to govern AI usage in mental health settings.

8 Acknowledge

Our paper underwent a thorough review by the University of Washington Human Subjects Division (HSD). On November 28, 2023, the HSD determined that our proposed activity does not involve human subjects as defined by federal and state regulations. Consequently, review and approval by the University of Washington IRB is not required.

This determination applies specifically to the activities described in our application (IRB ID: STUDY00019126). We appreciate the guidance provided by the HSD.

References

- Theyazn H. H. Aldhyani, Saleh Nagi Alsubari, Ali Saleh Alshebami, Hasan Alkahtani, and Zeyad A. T. Ahmed. 2022. [Detecting and Analyzing Suicidal Ideation on Social Media Using Deep Learning and Machine Learning Models](#). *International Journal of Environmental Research and Public Health*, 19(19):12635.
- John W. Ayers, Adam Poliak, Mark Dredze, Eric C. Leas, Zechariah Zhu, Jessica B. Kelley, Dennis J. Faix, Aaron M. Goodman, Christopher A. Longhurst, Michael Hogarth, and Davey M. Smith. 2023. [Comparing physician and artificial intelligence chatbot responses to patient questions posted to a public social media forum](#). *JAMA Internal Medicine*, 183(6):589–596.
- Yara Badr, Fares Al-Shargie, Usman Tariq, Fabio Babiloni, Fadwa Al Mughairbi, and Hasan Al-Nashash. 2023. [Classification of Mental Stress using Dry EEG Electrodes and Machine Learning](#). In *2023 Advances in Science and Engineering Technology International Conferences (ASET)*, pages 1–5. ISSN: 2831-6878.
- Borwin Bandelow, Sophie Michaelis, and Dirk Wedekind. 2017. [Treatment of anxiety disorders](#). *Dialogues in Clinical Neuroscience*, 19(2):93–107.

773	Basabdatta Sen Bhattacharya and Vibhav Sinai Pisurlenkar. 2023. Assistive Chatbots for healthcare: a succinct review . <i>arXiv preprint</i> . ArXiv:2308.04178 [cs].	827
774		828
775		829
776		830
777	David M Blei. 2003. Latent Dirichlet Allocation. <i>Journal of Machine Learning Research</i> .	831
778		832
779	Sergio Burdisso, Ernesto Reyes-Ramírez, Esaú Villatoro-Tello, Fernando Sánchez-Vega, Pastor López-Monroy, and Petr Motlicek. 2024. Daic-woz: On the validity of using the therapist’s prompts in automatic depression detection from clinical interviews . <i>Preprint</i> , arXiv:2404.14463.	833
780		834
781		835
782		836
783		837
784		838
785	Siyuan Chen, Mengyue Wu, Kenny Q. Zhu, Kunyao Lan, Zhiling Zhang, and Lyuchun Cui. 2023. Llm-empowered chatbots for psychiatrist and patient simulation: Application and evaluation . <i>Preprint</i> , arXiv:2305.13614.	839
786		840
787		841
788		842
789		843
790	Pim Cuijpers, Eirini Karyotaki, Erica Weitz, Gerhard Andersson, Steven D. Hollon, and Annemieke van Straten. 2014. The effects of psychotherapies for major depression in adults on remission, recovery and improvement: A meta-analysis . <i>Journal of Affective Disorders</i> , 159:118–126.	844
791		845
792		846
793		847
794		848
795		849
796	Zheng Danna, Liu Danyang, Lapata Mirella, and Pan Jeff Z. 2024. TrustScore: Reference-Free Evaluation of LLM Response Trustworthiness .	850
797		851
798		852
799	Saskia M. de Maat, Jack Dekker, Robert A. Schoevers, and Frans de Jonghe. 2007. Relative efficacy of psychotherapy and combined therapy in the treatment of depression: A meta-analysis . <i>European Psychiatry</i> , 22(1):1–8.	853
800		854
801		855
802		856
803		857
804	Zohar Elyoseph and Inbar Levkovich. 2024. Comparing the perspectives of generative AI, mental health experts, and the general public on schizophrenia recovery: Case vignette study . <i>JMIR Ment Health</i> , 11:53043.	858
805		859
806		860
807		861
808		862
809	Lauryn J Hagg, Stephanie S Merkouris, Gypsy A O’Dea, Lauren M Francis, Christopher J Greenwood, Matthew Fuller-Tyszkiewicz, Elizabeth M Westrupp, Jacqui A Macdonald, and George J Youssef. 2022. Examining Analytic Practices in Latent Dirichlet Allocation Within Psychological Science: Scoping Review . <i>Journal of Medical Internet Research</i> , 24(11):e33166.	863
810		864
811		865
812		866
813		867
814		868
815		869
816		870
817	Ylber Januzaj and Artan Luma. 2022. Cosine similarity – a computing approach to match similarity between higher education programs and job market demands based on maximum number of common words . <i>International Journal of Emerging Technologies in Learning (iJET)</i> , 17:258–268.	871
818		872
819		873
820		874
821		875
822		876
823	Rabie Karrouri, Zakaria Hammani, Roukaya Benjeloun, and Yassine Otheman. 2021. Major depressive disorder: Validated treatments and future challenges . <i>World Journal of Clinical Cases</i> , 9(31):9350–9367.	877
824		878
825		879
826		880
	Luoma Ke, Song Tong, Peng Cheng, and Kaiping Peng. 2024. Exploring the Frontiers of LLMs in Psychological Applications: A Comprehensive Review . <i>arXiv preprint</i> . ArXiv:2401.01519 [cs].	881
		882
	Igor Kotenko, Yash Sharma, and Alexander Branitskiy. 2021. Predicting the Mental State of the Social Network Users based on the Latent Dirichlet Allocation and fastText . In <i>2021 11th IEEE International Conference on Intelligent Data Acquisition and Advanced Computing Systems: Technology and Applications (IDAACS)</i> , volume 1, pages 191–195. ISSN: 2770-4254.	883
		884
	Alfirna Rizqi Lahitani, Adhistya Erna Permanasari, and Noor Akhmad Setiawan. 2016. Cosine similarity to determine similarity measure: Study case in online essay assessment . In <i>2016 4th International Conference on Cyber and IT Service Management</i> , pages 1–6, Bandung, Indonesia. IEEE.	885
		886
	Jintian Lin. 2023. Measuring and analyzing the information entropy value of key Audit matters (KAMs) disclosure at the system and reporting scale . <i>Heliyon</i> , 10(1):e23255.	887
		888
	Manikanta Loya, Divya Anand Sinha, and Richard Futrell. 2023. Exploring the Sensitivity of LLMs’ Decision-Making Capabilities: Insights from Prompt Variation and Hyperparameters . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 3711–3716. ArXiv:2312.17476 [cs].	889
		890
	Tania Manríquez Roa, Nikola Biller-Andorno, and Manuel Trachsel. 2021. <i>The Ethics of Artificial Intelligence in Psychotherapy</i> , chapter The Ethics of Artificial Intelligence in Psychotherapy. Oxford University Press.	891
		892
	Bertalan Meskó and Eric J. Topol. 2023. The imperative for regulatory oversight of large language models (or generative AI) in healthcare . <i>npj Digital Medicine</i> , 6(1):1–6. Publisher: Nature Publishing Group.	893
		894
	Inaba Michimasa, Ukiyo Mariko, and Takamizo Keiko. 2024. Can Large Language Models be Used to Provide Psychological Counselling? An Analysis of GPT-4-Generated Responses Using Role-play Dialogues .	895
		896
	Adam S. Miner, Liliana Laranjo, and A. Baki Kocaballi. 2020. Chatbots in the fight against the COVID-19 pandemic . <i>npj Digital Medicine</i> , 3(1):1–4. Number: 1 Publisher: Nature Publishing Group.	897
		898
	Sara Montagna, Stefano Ferretti, Lorenz Cuno Klopfenstein, Antonio Florio, and Martino Francesco Pengo. 2023. Data Decentralisation of LLM-Based Chatbot Systems in Chronic Disease Self-Management . In <i>Proceedings of the 2023 ACM Conference on Information Technology for Social Good</i> , pages 205–212, Lisbon Portugal. ACM.	899
		900
	Naveen S Pagad, Pradeep N, Khalid K. Almuzaini, Manish Maheshwari, Durgaprasad Gangodkar, Piyush Shukla, and Musah Alhassan. 2022. Clinical	901
		902

883 text data categorization and feature extraction using
884 medical-fissure algorithm and neg-seq algo-
885 rithm. *Computational Intelligence and Neuroscience*,
886 2022(1):5759521.

887 Domhnall Sheridan and Michael Carr. 2018. *Mens*
888 *sana: An Investigation into the Relationship between*
889 *Psychological Traits and Academic Success of First*
890 *Year Engineering Students*. In *2018 3rd International*
891 *Conference of the Portuguese Society for Engineering*
892 *Education (CISPEE)*, pages 1–5.

893 Karan Singhal, Shekoofeh Azizi, Tao Tu, S. Sara
894 Mahdavi, Jason Wei, Hyung Won Chung, Nathan
895 Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen
896 Pfohl, Perry Payne, Martin Seneviratne, Paul Gam-
897 ble, Chris Kelly, Abubakr Babiker, Nathanael Schärli,
898 Aakanksha Chowdhery, Philip Mansfield, Dina
899 Demner-Fushman, Blaise Agüera y Arcas, Dale Web-
900 ster, Greg S. Corrado, Yossi Matias, Katherine Chou,
901 Juraj Gottweis, Nenad Tomasev, Yun Liu, Alvin Ra-
902 jkomar, Joelle Barral, Christopher Semturs, Alan
903 Karthikesalingam, and Vivek Natarajan. 2023. *Large*
904 *language models encode clinical knowledge*. *Nature*,
905 620(7972):172–180. Number: 7972 Publisher:
906 Nature Publishing Group.

907 John Torous and Matcheri Keshavan. 2020. *COVID-19,*
908 *mobile health and serious mental illness*. *Schizophre-*
909 *nia Research*, 218:36–37.

910 Yunpeng Xiao, Kyrie Zhixuan Zhou, Yueqing Liang,
911 and Kai Shu. 2024. *Understanding the concerns and*
912 *choices of public when using large language models*
913 *for healthcare*. *arXiv preprint*. ArXiv:2401.09090
914 [cs].

915 A Generating Psychologically-Informed 916 Treatment Recommendations Based on 917 Detailed Case Information

918 **Introduction:** I will provide you with essential
919 information about a client needing psychological
920 consultation, which may include but is not limited
921 to age, gender, symptoms, past diagnostic informa-
922 tion, current life circumstances, treatment goals,
923 and expectations. Assuming you are a highly pro-
924 fessional psychotherapist, you are required to give
925 me psychological counseling treatment recommen-
926 dations.

927 Detailed Case Information:

928 Age and Gender: Provide the client’s age and
929 gender. Primary Symptoms: Describe the client’s
930 main psychological or emotional symptoms, such
931 as anxiety, depression, panic attacks, etc. Diag-
932 nostic History: Outline any formal diagnoses the
933 client has received in the past and the outcomes
934 of any treatments they underwent. Current Life
935 Circumstances: Describe the client’s family envi-
936 ronment, work or school situation, and social activ-

ities. Treatment Goals and Expectations: Specify
937 the concrete expectations and goals of the client
938 and their family for the treatment, including prob-
939 lems they hope to resolve and areas of life they
940 wish to improve. 941

942 Treatment Recommendation Requirements

943 **Diagnosis:** You need to give a diagnosis based
944 on the information I provided to you. You also
945 need to give a detail reason of why you give this
946 diagnosis.

947 **Psychotherapy Plan:** Theoretical Framework
948 Selection: Based on the client’s symptoms and di-
949 agnosis, choose an appropriate psychotherapy the-
950 oretical framework, such as Cognitive Behavioral
951 Therapy (CBT), psychodynamic therapy, humanis-
952 tic therapy, etc. Specific Therapeutic Techniques:
953 Detail the therapeutic techniques to be used, such
954 as exposure therapy, emotional restructuring, psy-
955 choeducation, etc.

956 **Medication Recommendations (if applicable):**
957 Suggest possible pharmacological treatments, not-
958 ing recommended types of medication, suggested
959 dosages, and potential side effects.

960 **Supportive Therapy Measures:** Recommend
961 supportive therapy measures such as group therapy,
962 family therapy, or other community resources to
963 enhance the effectiveness of the primary treatment
964 plan.

965 **Lifestyle and Behavioral Advice:** Provide rec-
966 ommendations for lifestyle adjustments that im-
967 prove overall health and psychological state, in-
968 cluding regular physical activity, healthy eating,
969 and good sleep habits.

970 **Monitoring and Adjustment:** Describe the pro-
971 posed evaluation and monitoring plan to regularly
972 check the effectiveness of the treatment and adjust
973 the treatment plan as needed.

974 **Output Format Requirements:** Please provide
975 the treatment plan in a report format, where each
976 section is clearly titled and thoroughly described.
977 Language and Expression:

978 **Use precise professional terminology, ensur-**
979 **ing that language is clear, rigorous, yet empa-**
980 **thetic and understanding toward the client.**

981 Ethical Considerations

982 B Study Cases and ChatGPT-Generated 983 Responses Used in the Evaluation

984 The responses were generated on 5/7/2024 and
985 5/8/2024 by OpenAI’s large language model, GPT-
986 4. The links to the original chat history with Chat-

987 GPT were listed below:
988 Study Case 1:
989 [https://chat.openai.com/share/23f1b1f1-021b-](https://chat.openai.com/share/23f1b1f1-021b-4d82-be34-dd71ba6d1348)
990 [4d82-be34-dd71ba6d1348](https://chat.openai.com/share/23f1b1f1-021b-4d82-be34-dd71ba6d1348)
991
992 Study Case 2:
993 [https://chat.openai.com/share/91c6bbca-cd58-](https://chat.openai.com/share/91c6bbca-cd58-4510-8894-55ad9d773112)
994 [4510-8894-55ad9d773112](https://chat.openai.com/share/91c6bbca-cd58-4510-8894-55ad9d773112)
995
996 Study Case 3:
997 [https://chat.openai.com/share/73dda11f-afeb-](https://chat.openai.com/share/73dda11f-afeb-4fb4-b5e1-4faa14cb4c72)
998 [4fb4-b5e1-4faa14cb4c72](https://chat.openai.com/share/73dda11f-afeb-4fb4-b5e1-4faa14cb4c72)
999
1000 Study Case 4:
1001 [https://chat.openai.com/share/a721c2a6-1405-](https://chat.openai.com/share/a721c2a6-1405-4bb9-a1dd-ea80beb78a9a)
1002 [4bb9-a1dd-ea80beb78a9a](https://chat.openai.com/share/a721c2a6-1405-4bb9-a1dd-ea80beb78a9a)
1003
1004 Study Case 5:
1005 [https://chat.openai.com/share/b94e1d7f-7f52-](https://chat.openai.com/share/b94e1d7f-7f52-49a3-b0ab-9e14c74cbf1a)
1006 [49a3-b0ab-9e14c74cbf1a](https://chat.openai.com/share/b94e1d7f-7f52-49a3-b0ab-9e14c74cbf1a)
1007
1008 Study Case 6:
1009 [https://chat.openai.com/share/78605bcd-473d-](https://chat.openai.com/share/78605bcd-473d-4e83-b8c0-467700083251)
1010 [4e83-b8c0-467700083251](https://chat.openai.com/share/78605bcd-473d-4e83-b8c0-467700083251)
1011
1012 Study Case 7:
1013 [https://chat.openai.com/share/78605bcd-473d-](https://chat.openai.com/share/78605bcd-473d-4e83-b8c0-467700083251)
1014 [4e83-b8c0-467700083251](https://chat.openai.com/share/78605bcd-473d-4e83-b8c0-467700083251)
1015
1016 Study Case 8:
1017 [https://chat.openai.com/share/80dc4e95-07dd-](https://chat.openai.com/share/80dc4e95-07dd-4bae-a5ec-2c20d12b41fc)
1018 [4bae-a5ec-2c20d12b41fc](https://chat.openai.com/share/80dc4e95-07dd-4bae-a5ec-2c20d12b41fc)
1019
1020 Study Case 9:
1021 [https://chat.openai.com/share/41571272-52f1-](https://chat.openai.com/share/41571272-52f1-4954-9431-87eb04a3cd16)
1022 [4954-9431-87eb04a3cd16](https://chat.openai.com/share/41571272-52f1-4954-9431-87eb04a3cd16)
1023
1024 Study Case 10:
1025 [https://chat.openai.com/share/7993a913-b157-](https://chat.openai.com/share/7993a913-b157-4374-9618-fc34470c8bad)
1026 [4374-9618-fc34470c8bad](https://chat.openai.com/share/7993a913-b157-4374-9618-fc34470c8bad)
1027

C Entropy Comparison Table

1028

Table 2: Entropy Comparison

Case	Orig. Entropy	LLM Entropy	% Diff.	Abs. Diff.
Case Study 1	1.163029	1.167993	-0.43	0.004964
Case Study 2	1.019811	1.250628	-22.63	0.230817
Case Study 3	1.266897	1.130842	10.74	0.136054
Case Study 4	1.065192	1.115729	-4.74	0.050537
Case Study 5	1.144286	1.167156	-2.00	0.022870
Case Study 6	1.094543	1.471529	-34.44	0.376986
Case Study 7	1.208523	1.408343	-16.53	0.199819
Case Study 8	1.336413	1.189755	10.97	0.146658
Case Study 9	1.102466	1.270382	-15.23	0.167915
Case Study 10	1.221561	1.168837	4.32	0.052724
Mean Entropy	1.162272	1.234119		
Standard Deviation	0.097259	0.119291		

D Entropy Comparison Table - Box plot

Table 3: Entropy Comparison Box plot

Scores	Median	IQR	Q1	Q3	Min	Max
Entropy (Psychotherapist)	1.15	0.12	1.10	1.22	1.02	1.34
Entropy (Large Language Model)	1.18	0.10	1.17	1.27	1.12	1.47

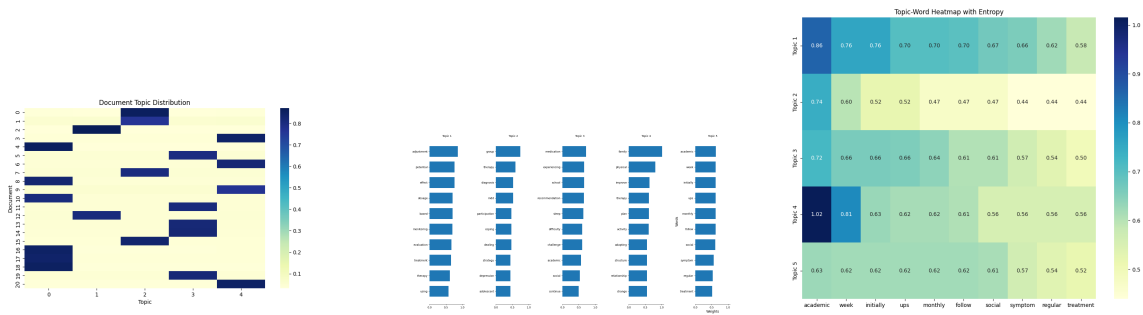


Figure 8: LLMs Analysis Results for Case Study 1

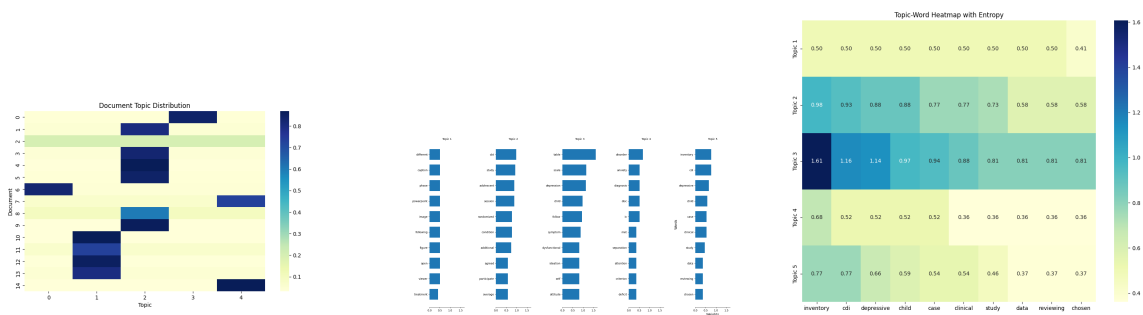


Figure 9: Original Analysis Results for Case Study 1

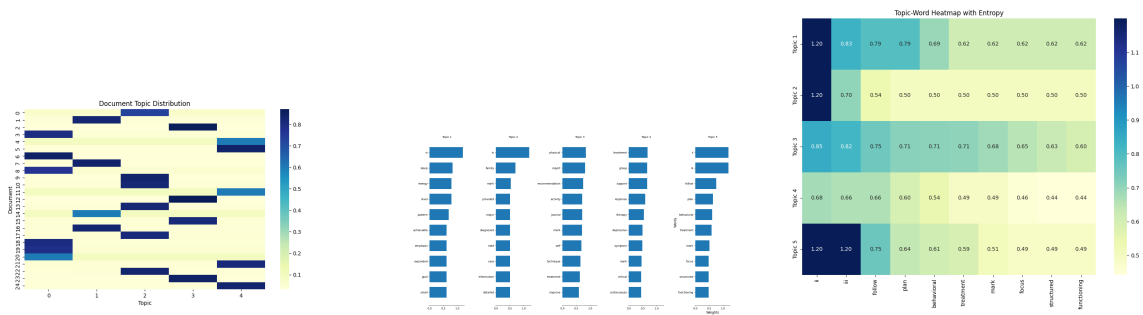


Figure 10: LLMs Analysis Results for Case Study 2

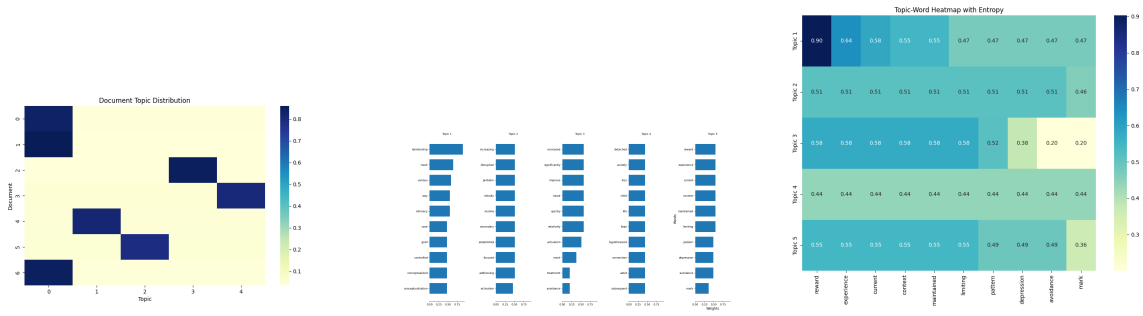


Figure 11: Original Analysis Results for Case Study 2

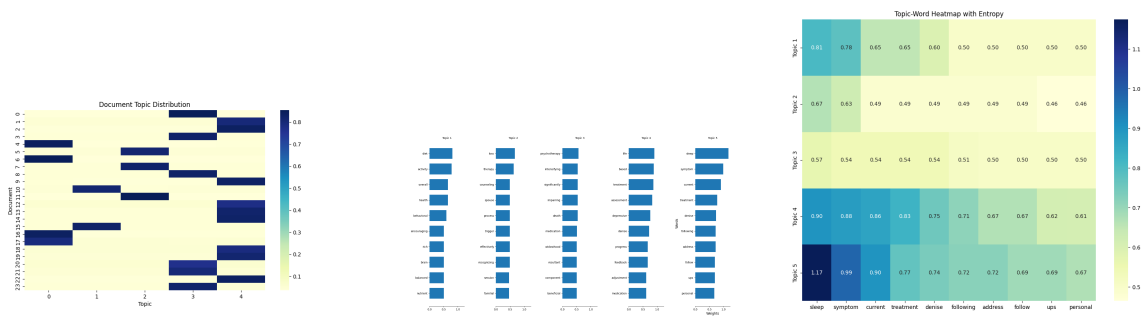


Figure 12: LLMs Analysis Results for Case Study 3

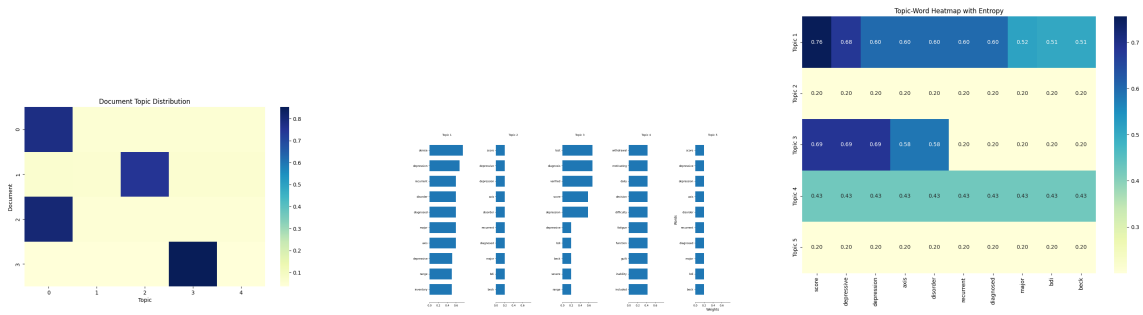


Figure 13: Original Analysis Results for Case Study 3

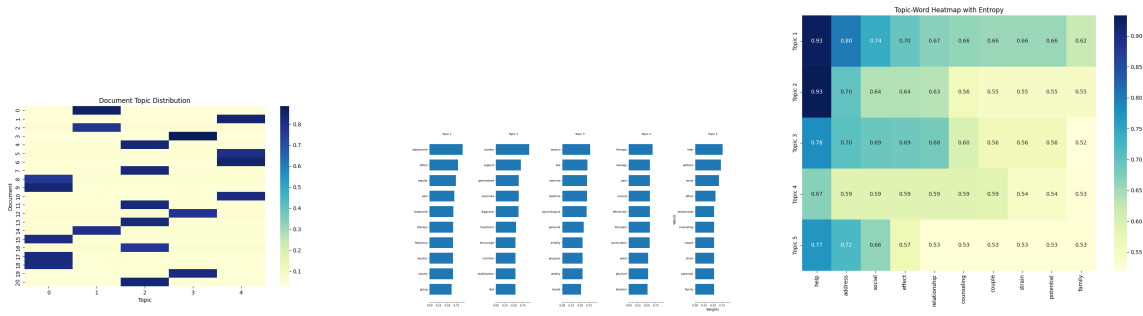


Figure 14: LLMs Analysis Results for Case Study 4

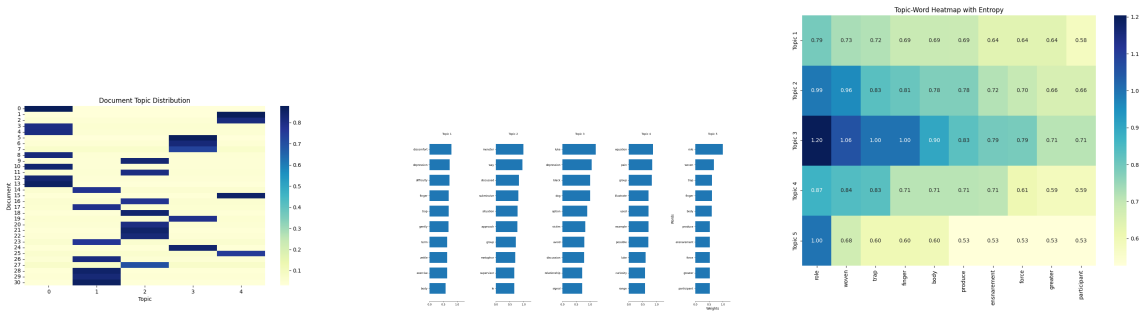


Figure 15: Original Analysis Results for Case Study 4

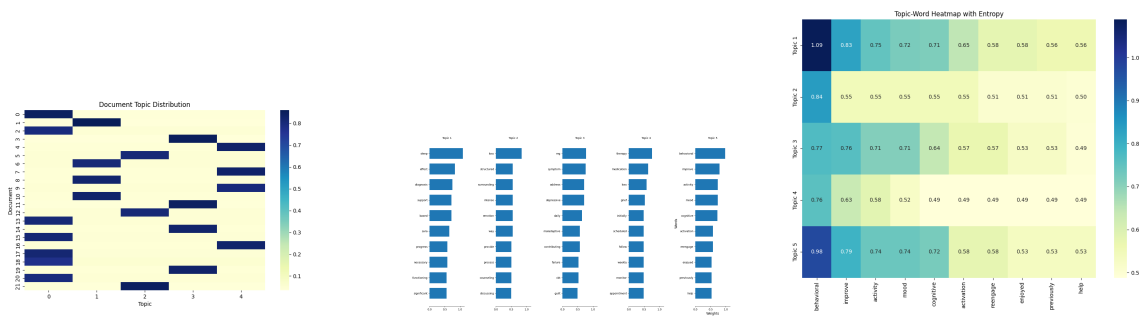


Figure 16: LLMs Analysis Results for Case Study 5

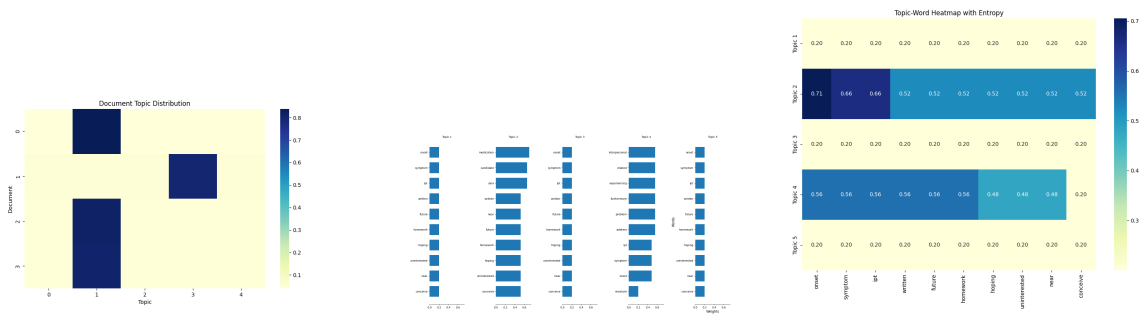


Figure 17: Original Analysis Results for Case Study 5

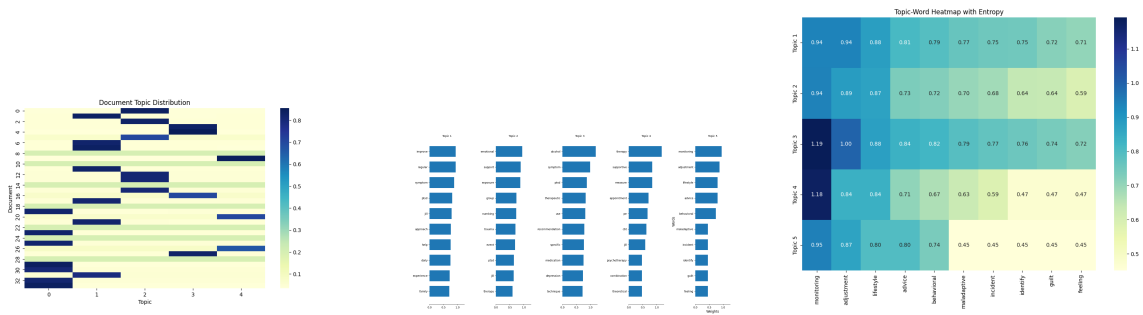


Figure 18: LLMs Analysis Results for Case Study 6

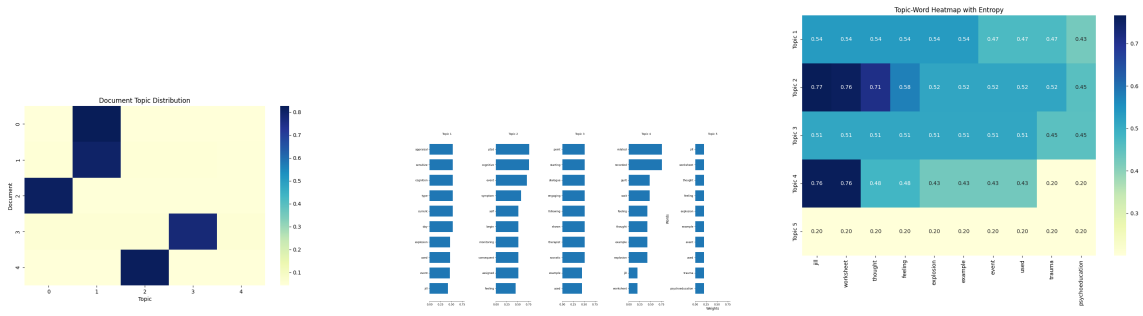


Figure 19: Original Analysis Results for Case Study 6

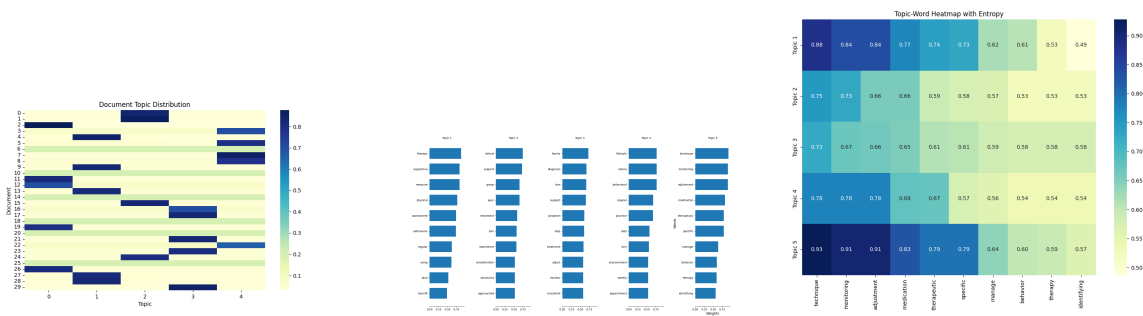


Figure 20: LLMs Analysis Results for Case Study 7

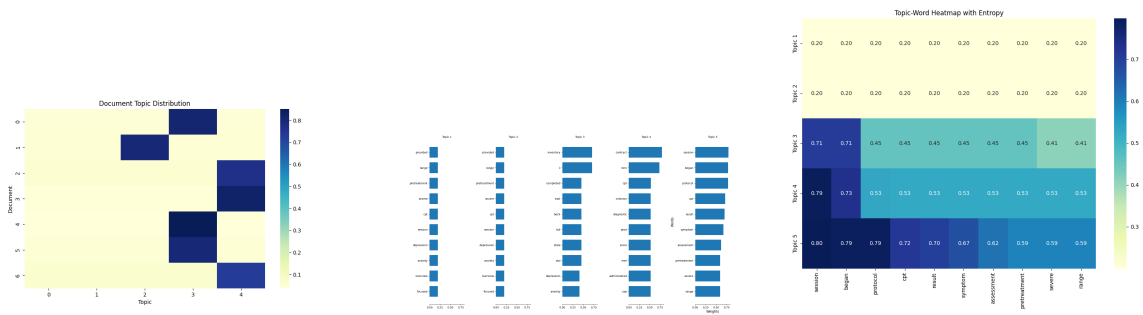


Figure 21: Original Analysis Results for Case Study 7

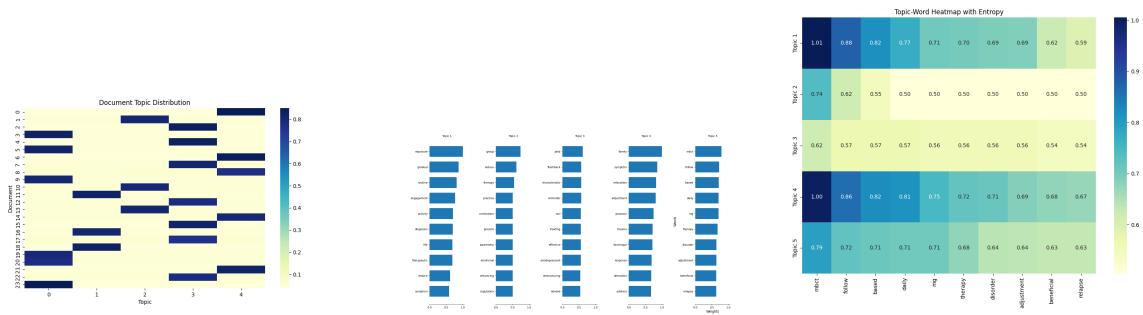


Figure 22: LLMs Analysis Results for Case Study 8

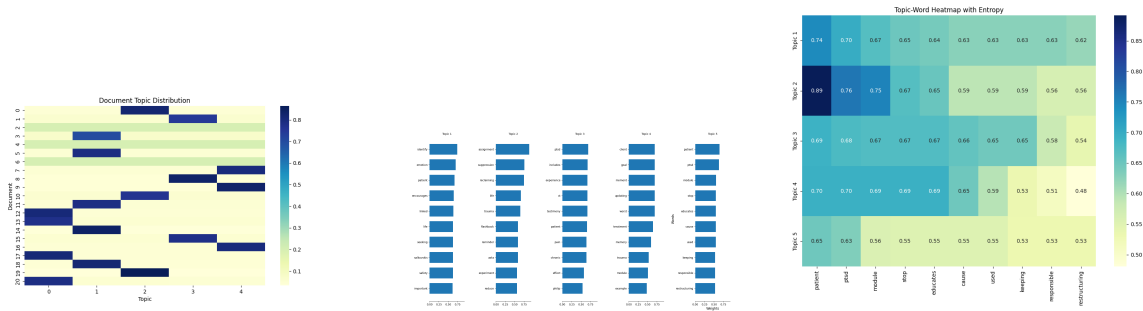


Figure 23: Original Analysis Results for Case Study 8

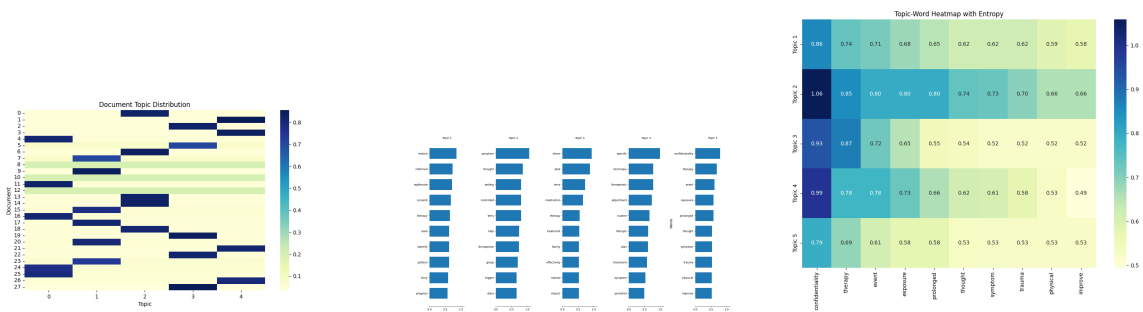


Figure 24: LLMs Analysis Results for Case Study 9

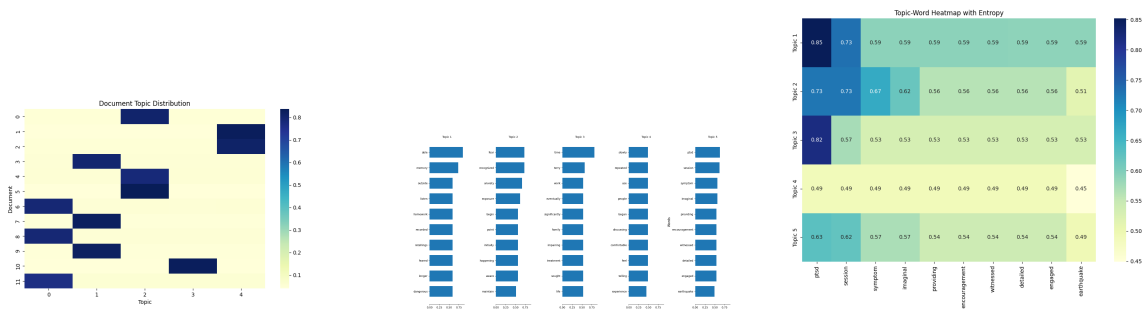


Figure 25: Original Analysis Results for Case Study 9

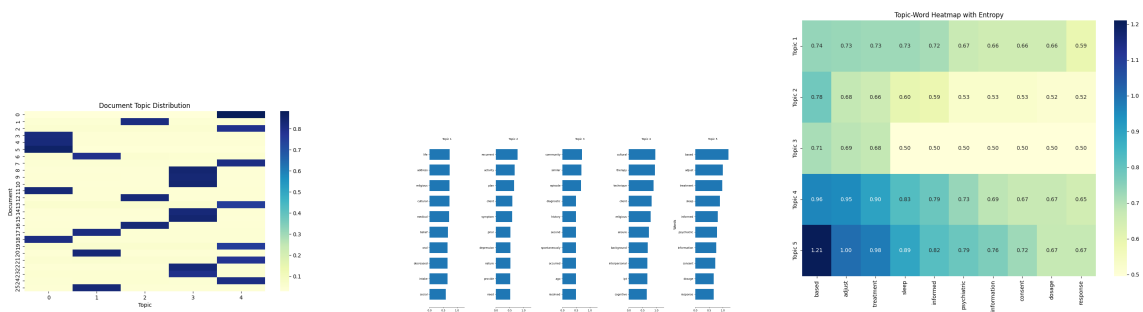


Figure 26: LLMs Analysis Results for Case Study 10

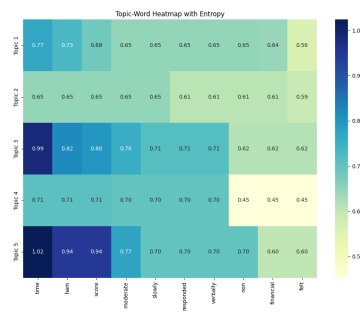
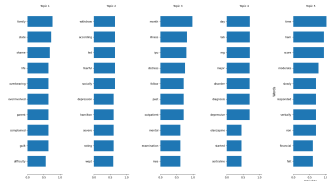
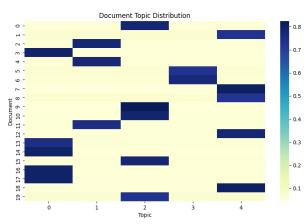


Figure 27: Original Analysis Results for Case Study 10