

A Survey of Large Language Models in Psychotherapy: Current Landscape and Future Directions

Anonymous ACL submission

Abstract

Mental health remains a critical global challenge, with increasing demand for accessible, effective interventions. Large language models (LLMs) offer promising solutions in psychotherapy by enhancing the assessment, diagnosis, and treatment of mental health conditions through dynamic, context-aware interactions. This survey provides a comprehensive overview of the current landscape of LLM applications in psychotherapy, highlighting the roles of LLMs in symptom detection, severity estimation, cognitive assessment, and therapeutic interventions. We present a novel conceptual taxonomy to organize the psychotherapy process into three core components: assessment, diagnosis, and treatment, and examine the challenges and advancements in each area. The survey also addresses key research gaps, including linguistic biases, limited disorder coverage, and underrepresented therapeutic models. Finally, we discuss future directions to integrate LLMs into a holistic, end-to-end psychotherapy framework, addressing the evolving nature of mental health conditions and fostering more inclusive, personalized care.

1 Introduction

Mental health plays an increasingly critical role in current healthcare and social well-being. The high prevalence of common psychological disorders, such as depression and anxiety, has led to a growing demand for accessible and effective psychological interventions. However, the core of psychotherapy resides in *dynamic, contextual* interpersonal interactions—therapists should continuously assess and adjust their intervention strategies (Wampold and Imel, 2015) based on the patient’s emotional fluctuations, verbal expressions, and social background, fostering a strong therapeutic alliance (Stubbe, 2018) to achieve symptom resilience. This deep and flexible process contrasts

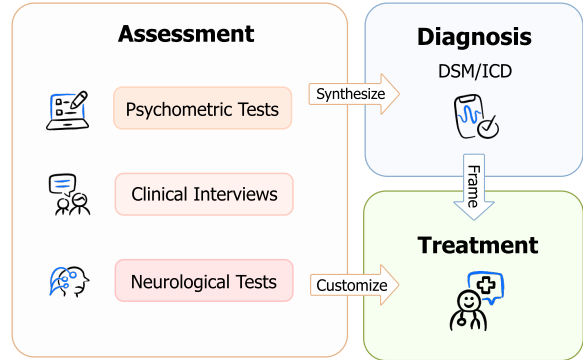


Figure 1: The dynamic and interrelated network among assessment, diagnosis, and treatment in psychotherapy.

sharply with traditional NLP, which is typically limited to static or single-task settings.

Large language models (LLMs) offer a new perspective to addressing this challenge. By leveraging their capability to model extensive context and perform multi-turn reasoning (Wang et al., 2024e; Li et al., 2024b), LLMs can capture rich semantics and emotional signals in dialogues (Ma et al., 2025), enabling end-to-end language understanding and generation. In assessment, LLMs can extract potential symptom cues from vague and fragmented expressions (Tu et al., 2024; Qiu et al., 2024). During diagnosis, they integrate subjective and objective patient information across multiple utterances (Chen et al., 2023a; Ren et al., 2024). In therapeutic interventions, they adapt conversational strategies based on patients’ real-time feedback, enabling more flexible and human-like interactions compared to traditional scripted systems (Lee et al., 2024b,d). As a result, LLMs have the potential to surpass the conventional “discrete label recognition” paradigm, evolving toward a model of continuous, progressive clinical reasoning, enabling seamless connections across *assessment, diagnosis,* and *treatment*, aligning more closely with therapists’ cognitive process and interaction flow.

However, existing research on applying LLMs

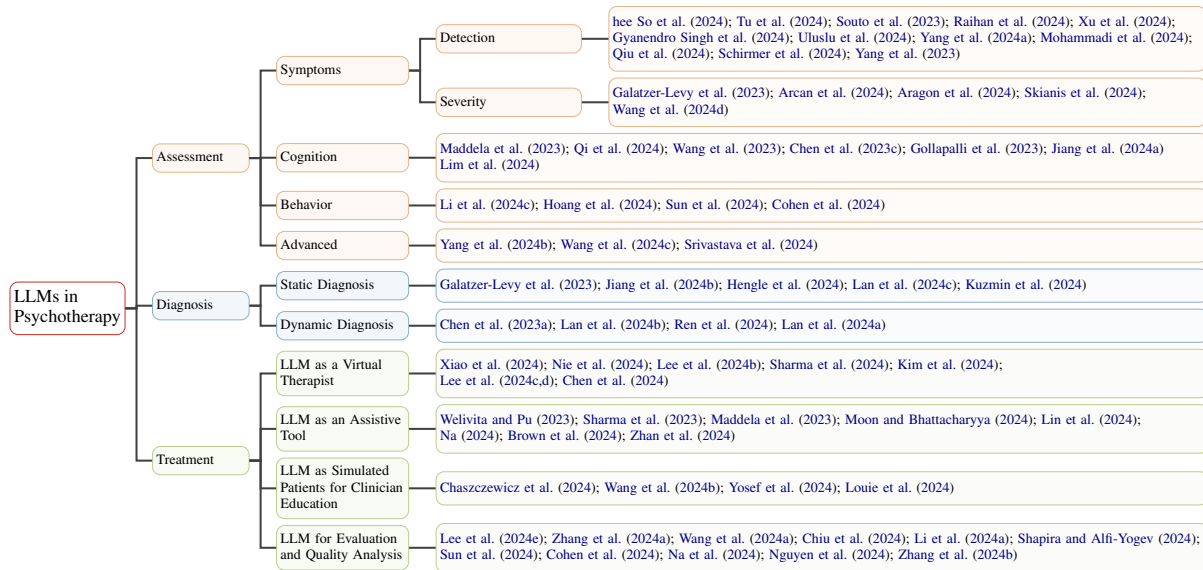


Figure 2: Taxonomy of Research on Large Language Models in Psychotherapy.

in this field remains somewhat *disjointed*. Many studies have utilized LLMs for isolated tasks, such as depression detection (Yang et al., 2023; Souto et al., 2023) or diagnosis (Jiang et al., 2024b), regarding them as superior feature extractors. Another research line has focused on developing mental health counseling chatbots (Chen et al., 2023b; Zhang et al., 2024a); however, these systems remain limited to partial assistance due to insufficient integration with clinical workflows. In other words, although LLMs hold the potential to span the entire continuum from assessment to intervention, they remain limited by the fragmented paradigms of traditional NLP, preventing them from fully leveraging their dynamic, contextual capabilities.

To address these gaps, we introduce the first *taxonomy* that divides the psychotherapy process into three essential dimensions: **Assessment**, **Diagnosis**, and **Treatment** and provides a systematic review of the recent advancements and challenges of LLMs in each stage. We further examine the current landscape from various perspectives, including the coverage of mental disorders, diversity of linguistic resources, alignment with psychotherapy theories, and the types of techniques employed, thereby sketching the overall distribution and characteristics of existing research. Building on this foundation, we discuss key challenges for the future, including issues of technical coherence, resource and language imbalances, and the disconnect between LLM-based approaches and established psychological practices. Through this *comprehensive review and framework*, we aim to offer methodological insights to inform future research

and facilitate the practical integration of intelligent systems across the entire psychotherapy process.

Organization of This Survey. We present the first comprehensive survey of recent advancements in applying LLMs to psychotherapy. We introduce a conceptual taxonomy that organizes psychotherapy into three core components—Assessment, Diagnosis, and Treatment—and details their dynamic interrelations (Section §2). We review how LLMs are applied within these components, highlighting their roles in facilitating assessments, refining diagnostic processes, and enhancing treatment strategies (Section §3). We examine current research trends, including symptom and language coverage as well as the distribution of various models and techniques (Section §4). Finally, we discuss open challenges and outline promising directions for future work (Section §5).

2 Conceptual Taxonomy

To establish a standardized framework for understanding psychotherapy, we propose a hierarchical taxonomy aligned with the American Psychological Association’s tripartite model of psychotherapeutic processes¹. As illustrated in Figure 1, this taxonomy organizes psychotherapy into three core components: (1) Assessment, (2) Diagnosis, and (3) Treatment, with dynamic interconnections². Each component is detailed below.

¹<https://www.apa.org/topics/psychotherapy>

²Throughout this taxonomy, the terms *Assessment*, *Diagnosis*, and *Treatment* specifically refer to the three core components of psychotherapy.

| | | |
|-----|---|-----|
| 131 | 2.1 Assessment | |
| 132 | Definition. Psychological assessment constitutes | |
| 133 | the systematic collection and interpretation of data | |
| 134 | regarding an individual's cognitive, emotional, and | |
| 135 | behavioral functioning (Cohen et al., 1996; Ka- | |
| 136 | plan and Saccuzzo, 2001). This process employs | |
| 137 | psychometric tests, structured clinical interviews, | |
| 138 | behavioral observations, and collateral information | |
| 139 | to establish a multidimensional profile of psycho- | |
| 140 | logical states (Groth-Marnat, 2009). | |
| 141 | Significance. As the foundational stage of psy- | |
| 142 | chotherapy, assessment provides the empirical ba- | |
| 143 | sis for understanding a client's unique psycholog- | |
| 144 | ical landscape. It enables therapists to identify | |
| 145 | symptom patterns (Phillips et al., 2007), track tem- | |
| 146 | poral changes (Barkham et al., 1993), and con- | |
| 147 | textualize subjective experiences within objective | |
| 148 | frameworks (Groth-Marnat, 2009). The contin- | |
| 149 | uous nature of psychological assessment allows | |
| 150 | for real-time adjustments to therapeutic strate- | |
| 151 | gies (Schiepek et al., 2016), ensuring interventions | |
| 152 | remain responsive to evolving client needs. | |
| 153 | 2.2 Diagnosis | |
| 154 | Definition. Diagnosis represents the analytical | |
| 155 | process of categorizing psychological distress us- | |
| 156 | ing established nosological systems such as the | |
| 157 | DSM-5 (American Psychiatric Association, 2022) | |
| 158 | and ICD-11 (World Health Organization, 2019). | |
| 159 | This involves differentiating normative emotional | |
| 160 | responses from pathological conditions while con- | |
| 161 | sidering cultural (Teo, 2010) and developmen- | |
| 162 | tal (Kawa and Giordano, 2012) variables that influ- | |
| 163 | ence symptom manifestation. | |
| 164 | Significance. Diagnosis serves as the conceptual | |
| 165 | bridge between assessment and treatment, provid- | |
| 166 | ing a structured framework for intervention plan- | |
| 167 | ning (Jensen-Doss and Hawley, 2011). By align- | |
| 168 | ing clinical observations with standardized crite- | |
| 169 | ria, it enhances communication among profession- | |
| 170 | als (Craddock and Mynors-Wallis, 2014) and facil- | |
| 171 | itates evidence-based decision-making (American | |
| 172 | Psychiatric Association, 2006). | |
| 173 | 2.3 Treatment | |
| 174 | Definition. Treatment includes evidence-based | |
| 175 | interventions designed to reduce psychological | |
| 176 | distress and improve functioning (American Psy- | |
| 177 | chiatric Association, 2006). These interventions | |
| 178 | work by building a therapeutic alliance (Elvins and | |
| | Green, 2008), restructuring cognition (Ezawa and | 179 |
| | Hollon, 2023), and modifying behavior (Martin | 180 |
| | and Pear, 2019), all typically grounded in well- | 181 |
| | established theoretical orientations. | 182 |
| | Significance. Treatment transforms the theories | 183 |
| | and information gleaned from assessment and di- | 184 |
| | agnosis into practical interventions (Prochaska and | 185 |
| | Norcross, 2018) that directly address the client's | 186 |
| | psychological distress (Barlow, 2021) and foster | 187 |
| | personal growth (Lambert, 2013). | 188 |
| | 2.4 Interrelations | 189 |
| | The taxonomy's components interact through three | 190 |
| | dynamic processes that define psychotherapy as a | 191 |
| | complex adaptive system: | 192 |
| | Synthesizing (Assessment → Diagnosis) The | 193 |
| | dialectical integration of observational data with | 194 |
| | nosological frameworks enables diagnostic classifi- | 195 |
| | cations to contextualize assessment findings, <i>syn-</i> | 196 |
| | <i>thesizing</i> the patient's various symptoms and be- | 197 |
| | havioral patterns into a diagnostic result (Rencic | 198 |
| | et al., 2016). | 199 |
| | Framing (Diagnosis → Treatment) Diagnosis | 200 |
| | functions as a <i>framing</i> mechanism, integrating com- | 201 |
| | plex and diverse symptoms into a coherent classi- | 202 |
| | fication that establishes a clear blueprint for treat- | 203 |
| | ment (American Psychiatric Association, 2022). | 204 |
| | Customization (Assessment → Treatment) A | 205 |
| | process where treatment plans are continuously | 206 |
| | <i>refined</i> based on assessment results, considering | 207 |
| | individual differences without being constrained | 208 |
| | by diagnostic labels, to enhance therapeutic effec- | 209 |
| | tivenesss (Waszczuk et al., 2017). | 210 |
| | 3 LLMs in Psychotherapy | 211 |
| | 3.1 Assessment | 212 |
| | Symptom Detection leverages LLMs to iden- | 213 |
| | tify mental health conditions including depression, | 214 |
| | anxiety, PTSD, and suicidal ideation, demonstrat- | 215 |
| | ing robust performance and multidimensional ap- | 216 |
| | plicability across diverse scenarios. Yang et al. | 217 |
| | (2023) systematically evaluated GPT-3.5, Instruct- | 218 |
| | GPT3, and LLaMA models across 11 datasets, re- | 219 |
| | vealing that emotion-enhanced chain-of-thought | 220 |
| | prompting improves interpretability yet remains in- | 221 |
| | ferior to specialized supervised methods. hee So | 222 |
| | et al. (2024) achieved 70.8% zero-shot symptom | 223 |
| | retrieval accuracy in Korean psychiatric interviews | 224 |
| | using GPT-4 Turbo, while their fine-tuned GPT-3.5 | 225 |

| Study | Text Granularity | Best Technique | NLP Task | Assessment Focus |
|-------------------------------|------------------------|---------------------|------------|-----------------------|
| <i>Symptom Detection</i> | | | | |
| Yang et al. (2023) | Single Post | Emotion Prompting | BC/MCC/EG | Multiple Symptoms |
| hee So et al. (2024) | Multi-turn Dialogue | Fine-Tuning | MLC/IE/SUM | Multiple Symptoms |
| Tu et al. (2024) | Multi-turn Dialogue | Few-Shot Prompting | MLC/IE/SUM | PTSD |
| Souto et al. (2023) | Single Post | Fine-Tuning | MLC/EG | Depression |
| Raihan et al. (2024) | Single Post | Few-Shot Prompting | MCC | Multiple Symptoms |
| Gyanendro Singh et al. (2024) | Posts From One User | Chain-of-Thought | IE/SUM | Suicidal Ideation |
| Uluslu et al. (2024) | Posts From One User | Role Prompting | IE/SUM | Suicidal Ideation |
| Yang et al. (2024a) | Single Post | Fine-Tuning | BC/MCC/EG | Multiple Symptoms |
| Xu et al. (2024) | Single Post | Fine-Tuning | BC/EG | Multiple Symptoms |
| Mohammadi et al. (2024) | Single Post | Few-Shot Prompting | MLC | Multiple Symptoms |
| Qiu et al. (2024) | Single Post | Fine-Tuning | MLC | Suicidal Ideation |
| Schirmer et al. (2024) | Single Post | Zero-Shot Prompting | BC | PTSD |
| <i>Symptom Severity</i> | | | | |
| Galatzer-Levy et al. (2023) | Multi-turn Dialogue | Zero-Shot Prompting | TR | Depression/PTSD |
| Arcan et al. (2024) | Multi-turn Dialogue | Zero-Shot Prompting | TR | Depression/Anxiety |
| Aragon et al. (2024) | Posts From One User | Zero-Shot Prompting | TR | Depression |
| Wang et al. (2024d) | Posts From One User | Zero-Shot Prompting | TR | Depression |
| Skianis et al. (2024) | Single Post | Zero-Shot Prompting | TR/MCC | Depression/Suicide |
| <i>Cognition</i> | | | | |
| Maddela et al. (2023) | Single Sentence | Few-Shot Prompting | MLC | Cognitive Distortions |
| Qi et al. (2024) | Single Post | Fine-Tuning | MLC | Cognitive Distortions |
| Wang et al. (2023) | Single Sentence | Few-Shot Prompting | MCC | Cognitive Distortions |
| Chen et al. (2023c) | Single-turn Dialogue | Zero-Shot Prompting | BC/MCC/EG | Cognitive Distortions |
| Gollapalli et al. (2023) | Single Post | Zero-Shot Prompting | MLC | Maladaptive Schemas |
| Jiang et al. (2024a) | Single Post | Zero-Shot Prompting | MCC/SUM | Cognitive Pathways |
| Lim et al. (2024) | Single-turn Dialogue | Multi-Agent Debate | MCC | Cognitive Distortions |
| <i>Behavior</i> | | | | |
| Li et al. (2024c) | Single Post | Zero-Shot Prompting | MLC/EG | Interpersonal Risk |
| Hoang et al. (2024) | Sentence From Dialogue | Few-Shot Prompting | MCC | MI-Adherent Behaviors |
| Sun et al. (2024) | Sentence From Dialogue | Zero-Shot Prompting | MCC | MI-Adherent Behaviors |
| Cohen et al. (2024) | Sentence From Dialogue | Zero-Shot Prompting | MCC | MI-Adherent Behaviors |

Table 1: Comparison of Psychological Assessment Studies by Input Characteristics and Methodology. **MLC**: Multi-Label Classification, **IE**: Information Extraction, **SUM**: Summarization, **MCC**: Multi-Class Classification, **BC**: Binary Classification, **TR**: Text Regression, **EG**: Explanation Generation. Studies are categorized through text granularity, optimal technical approach (*Best Technique*), NLP task formulation, and specific assessment focus.

attained 0.817 multi-label classification accuracy. Clinical applications show particular promise, as Tu et al. (2024) leveraged GPT-4 and Llama-2 to automate PTSD assessments through information extraction from 411 interviews, significantly enhancing diagnostic practicality.

Social media analysis benefits from approaches like Souto et al. (2023)’s interpretable depression detection framework, which demonstrated strong performance across Vicuna-13B and GPT-3.5 environments. Resource development advances include Raihan et al. (2024)’s *MentalHelp* dataset with 14 million instances, validated through GPT-3.5 zero-shot evaluations. For suicidal ideation monitoring, Gyanendro Singh et al. (2024) and Uluslu et al. (2024) achieved state-of-the-art evidence extraction in the CLPsych 2024 shared task through innovative prompting strategies. Open-source ini-

tiatives like *MentaLLaMA* by Yang et al. (2024a) and *Mental-LLM* by Xu et al. (2024) enable multi-symptom detection via instruction-tuned LLaMA variants, though Mohammadi et al. (2024)’s *Well-Dunn* framework reveals persistent gaps in GPT-family models’ explanation consistency.

Cross-lingual adaptations include Qiu et al. (2024)’s *PsyGUARD* system based on fine-tuned CHATGLM2-6B for Chinese suicide risk assessment, while Schirmer et al. (2024) demonstrated domain-specific RoBERTa models outperforming GPT-4 in cross-domain PTSD pattern analysis, highlighting the critical balance between model specialization and interpretability.

Symptom Severity focuses on estimating the level of mental health condition intensity, particularly for depression, anxiety, and PTSD. Clini-

cal evaluations reveal Med-PaLM 2’s zero-shot depression scoring attains clinician-level alignment on interview data (Galatzer-Levy et al., 2023), though with limited PTSD generalizability. When benchmarked against specialized Transformers on DAIC-WOZ dataset (Gratch et al., 2014), ChatGPT and Llama-2 exhibit moderate efficacy (Arcan et al., 2024), suggesting domain-specific architectures retain advantages in structured assessments. Shifting attention to social media data, Aragon et al. (2024) proposed a pipeline that retrieves depression-relevant text, summarizes it according to the Beck Depression Inventory (BDI) (Jackson-Koku, 2016), and then utilizes LLMs to predict symptom severity, achieving performance similar to expert evaluations on certain measures. In a similar vein, Wang et al. (2024d) introduced an explainable depression detection system that leverages multiple open-source LLMs to generate BDI-based answers, reporting near state-of-the-art performance without additional training data. Cross-lingual extensions emerge through Skianis et al. (2024)’s framework enabling severity prediction across 6 languages and 2 mental conditions.

Cognition centers on identifying and understanding maladaptive thinking patterns, such as cognitive distortions and early maladaptive schemas, using LLMs. Maddela et al. (2023) introduced a cognitive distortion dataset and employed a few-shot strategy with GPT-3.5 to generate, classify, and reframe them, while Qi et al. (2024) constructed two Chinese social media benchmarks for cognitive distortion detection and suicidal risk assessment, demonstrating that fine-tuned LLMs are more closely than zero-/few-shot methods to supervised baselines. In a related effort, Wang et al. (2023) released the C2D2 dataset containing 7,500 Chinese sentences with distorted thinking patterns. Expanding on detection methods, Chen et al. (2023c) proposed a Diagnosis of Thought prompting approach for GPT-4 and ChatGPT, which breaks down patient utterances into factual versus subjective content and supports the generation of interpretable diagnostic reasoning. Beyond cognitive distortions, Gollapalli et al. (2023) investigated zero-shot approaches with GPT-3.5 to identify early maladaptive schemas in mental health forums, highlighting challenges in label interpretability and prompt sensitivity. Complementarily, Jiang et al. (2024a) presented a hierarchical classification and summarization pipeline to extract cognitive pathways from

Chinese social media text, underscoring GPT-4’s strong performance albeit with occasional hallucinations. Finally, Lim et al. (2024) introduced a multi-agent debate framework for cognitive distortion classification, reporting substantial gains in both accuracy and specificity by synthesizing multiple LLM opinions before forming a final verdict.

Behavior highlights how user actions—or in the case of Motivational Interviewing (MI), language itself—can serve as a measurable indicator of one’s readiness for change. For instance, Li et al. (2024c) introduced the MAIMS framework, employing mental scales in a zero-shot setting to identify interpersonal risk factors on social media, thereby enhancing both interpretability and accuracy. In clinical dialogues, Hoang et al. (2024) demonstrated how LLMs can automatically detect a client’s motivational direction (e.g., change versus sustain talk) and commitment level, offering valuable insights for MI-based interventions. Extending such analyses to bilingual settings, Sun et al. (2024) proposed the BiMISC dataset and prompt strategies that enable LLMs to code MI behaviors across multiple languages with expert-level performance. Lastly, Cohen et al. (2024) presented MI-TAGS for automated annotation of global MI scores, illustrating how context-sensitive modeling can approximate human annotations in psychotherapy transcripts.

Advanced research has evolved beyond foundational assessment tasks to emphasize novel methodological paradigms, bias mitigation, and domain-specific summarization frameworks. For instance, Yang et al. (2024b) introduced *PsychoGAT*—an interactive, game-based approach that transforms standardized psychometric instruments into engaging narrative experiences, improving psychometric reliability, construct validity, and user satisfaction when measuring constructs such as depression, cognitive distortions, and personality traits. In parallel, Wang et al. (2024c) systematically investigated potential biases in various LLMs across multiple mental health datasets, revealing that even high-performing models exhibit unfairness related to demographic factors. The authors proposed fairness-aware prompts to substantially reduce such biases without sacrificing predictive accuracy. Furthermore, Srivastava et al. (2024) presented the *PIECE* framework, which adopts a planning-based approach to domain-aligned counseling summarization, structuring and filtering conversation content before integrating domain knowledge.

3.2 Diagnosis

Static Diagnosis is based on a fixed set of data, typically derived from complete dialogues or social media posts. Galatzer-Levy et al. (2023) highlighted the effectiveness of Med-PaLM 2 in psychiatric condition assessment from patient interviews and clinical descriptions without specialized training. Similarly, Jiang et al. (2024b) showcased LLMs’ superior performance on depression and anxiety detection on Russian datasets, particularly with noisy or small datasets. Hengle et al. (2024) evaluated PLMs and LLMs on multi-label classification in depression and anxiety, underscoring the ongoing challenges in applying LMs to mental health diagnostics. Besides, Lan et al. (2024c) introduced *DORIS*, a depression detection system integrating text embeddings with LLMs, utilizing symptom features, post-history, and mood course representations to make diagnostic predictions and generate explanatory outputs. Kuzmin et al. (2024) developed *ADOS-Copilot* for ASD diagnosis through diagnostic dialogues, employing In-context Enhancement, Interpretability Augmentation, and Adaptive Fusion based on real-world ADOS-2 clinic scenarios.

Dynamic Diagnosis involves real-time evaluation based on ongoing, interactive conversations between the patient and LLM, enabling more personalized and contextually relevant insights. Chen et al. (2023a) simulated psychiatrist-patient interactions with ChatGPT, in which the doctor chatbot focused on role, tasks, empathy, and questioning strategies, while the patient chatbot emphasized symptoms, language style, emotions, and resistance behaviors. Lan et al. (2024b) introduced the *Symptom-related and Empathy-related Ontology (SEO)*, grounded in DSM-5 and Helping Skills Theory, for depression diagnosis dialogues. Ren et al. (2024) dissected the doctor-patient relationship into psychologist’s empathy and proactive guidance and introduced *WundtGPT* that integrated these elements. Lan et al. (2024a) further presented the *AMC*, a self-improving conversational agent system for depression diagnosis through simulated dialogues between patient and psychiatrist agents.

3.3 Treatment

LLM as a Virtual Therapist centers on leveraging LLMs to directly engage in therapeutic conversations, often adopting multi-turn dialogues that incorporate recognized psychotherapeutic frame-

works. For instance, Xiao et al. (2024) proposed *HealMe* to facilitate cognitive reframing and empathetic support in line with established psychotherapy principles. Likewise, Nie et al. (2024) introduced *CaiTI*, a system embedded in everyday smart devices that conducts assessments of users’ daily functioning and delivers psychotherapeutic interventions through adaptive dialogue flows. In a similar vein, Lee et al. (2024b) presented *CoCoA*, specializing in identifying and resolving cognitive distortions via dynamic memory mechanisms and CBT-based strategies, while Sharma et al. (2024) proposed a step-by-step approach guiding users to execute self-guided cognitive restructuring through multiple interactive sessions. Beyond standard CBT protocols, Kim et al. (2024) focused on aiding psychiatric patients in journaling their experiences, thereby offering richer clinical insights, whereas Lee et al. (2024c) developed a multi-round CBT dataset to refine LLMs for direct counseling-like interactions. Additionally, multi-agent frameworks like *MentalAgora* (Lee et al., 2024d) highlighted personalized mental health support by integrating multiple specialized agents, and Chen et al. (2024) further explored “mixed chain-of-psychotherapies” to combine various therapeutic methods, aiming to enhance the emotional support and customization delivered by chatbot interactions.

LLM as an Assistive Tool refrains from providing a holistic therapy role but instead offers targeted support such as rewriting suboptimal counselor responses, generating controlled reappraisal prompts, or aiding clinicians in specific tasks. For example, Welivita and Pu (2023) proposed to rewrite responses that violate MI principles into MI-adherent forms, ensuring more consistent therapeutic dialogue. Meanwhile, Sharma et al. (2023) and Madela et al. (2023) focused on generating single-turn reframes of negative thoughts—often anchored in cognitive distortions—through controlled language attributes. On the detection side, Moon and Bhat-tacharyya (2024) built a multimodal pipeline to identify depression and provide CBT-style replies, albeit with an emphasis on technological assistance rather than full-fledged therapy. In the Chinese context, Lin et al. (2024) combined cognitive distortion detection with “positive reconstruction,” demonstrating a single-round rewrite approach for negative or distorted statements, while Na (2024) showcased a structured Q&A format that offers professional yet succinct CBT-based responses. From a

knowledge-distillation angle, [Brown et al. \(2024\)](#) demonstrated how smaller models could replicate GPT-4’s MI-style reflective statements, and [Zhan et al. \(2024\)](#) introduced a lighter-weight framework *RESORT* to guide smaller LLMs toward effective cognitive reappraisal prompts, thus enabling broader accessibility of self-help tools.

LLM as Simulated Patients for Clinician Education pivots toward generating synthetic yet realistic patient behaviors or multi-level feedback to train or support mental health practitioners. For instance, [Chaszczewicz et al. \(2024\)](#) leveraged LLMs to deliver multi-tier feedback on novice peer counselors’ conversational skills, significantly reducing the need for continuous expert oversight. Similarly, [Wang et al. \(2024b\)](#) using LLM-driven patient simulations that help trainees practice CBT core skills in a controlled, repeatable setup. In the realm of assessing therapy quality, [Yosef et al. \(2024\)](#) showcased a digital patient system to evaluate MI sessions, employing AI-generated transcripts to differentiate novice, intermediate, and expert therapeutic skill levels. Complementarily, [Louie et al. \(2024\)](#) offered *Roleplay-doh*, a pipeline wherein domain experts craft specialized principles that guide LLM-based role-playing agents, thereby providing customizable training for new therapists.

LLM for Evaluation and Quality Analysis targets the appraisal of therapy dialogue, counselor techniques, and treatment processes, typically without delivering direct interventions to clients. For instance, [Lee et al. \(2024e\)](#) augmented crisis counseling outcome prediction by fusing annotated counseling strategies with LLM-derived features, achieving substantially improved accuracy. In the Chinese context, [Zhang et al. \(2024a\)](#) introduced *CPsyCoun*, employing reports-based dialogue reconstruction and automated evaluation to verify counseling realism and professionalism. Beyond single-session analyses, [Wang et al. \(2024a\)](#) used simulated clients to assess perceived therapy outcomes, while [Chiu et al. \(2024\)](#) created the *BOLT* framework for systematically comparing LLM-based therapy behaviors with high- and low-quality human sessions. Further extending to online counseling, [Li et al. \(2024a\)](#) proposed an LLM-based approach to measure therapeutic alliance, whereas [Shapira and Alfi-Yogev \(2024\)](#) delineated therapist self-disclosure classification as a new NLP task. In the MI domain, [Sun et al. \(2024\)](#) and [Cohen et al. \(2024\)](#) collected bilingual transcripts to sys-

tematically annotate therapist–client exchanges for behavior coding and global scores, respectively. Additionally, multi-session perspectives emerge in [Na et al. \(2024\)](#), who proposed *IPAEval* to track long-term progress from the client’s viewpoint, and [Nguyen et al. \(2024\)](#) analyzed conversation redirection and its impact on patient–therapist alliance over multiple sessions. Finally, [Iftikhar et al. \(2024\)](#) and [Zhang et al. \(2024b\)](#) explored the disparities between LLM- and human-led CBT sessions, highlighting gaps such as empathy and cultural nuance while also introducing *CBT-Bench* to probe LLMs’ deeper psychotherapeutic competencies.

4 Current Landscape

Our survey encompasses a total of 69 studies in the field of LLMs in psychotherapy. Specifically, 33 studies address assessment, 9 focus on diagnosis, and 32 concentrate on treatment, with 5 studies overlapping across these dimensions. Approximately 74% of the studies employed commercial large language models, while about 77% used prompt-based techniques. This distribution highlights an imbalance in research focus across different stages of the psychotherapy process and reflects a heavy reliance on commercial models and prompt technologies.

Figure 3 presents a comprehensive analysis of the current research landscape in this field. Panel (a) reveals a significant linguistic bias in existing studies, with English-language corpora dominates. While there are limited studies involving Korean and Dutch languages, this highlights a substantial gap in multilingual research approaches. Panel (b) quantitatively demonstrates the distribution of mental health research focuses. Mental disorder-related studies constitute 32% of the total research corpus (represented by the orange outer ring). Within this subset, depression-focused research accounts for 50% of mental disorder studies, followed by anxiety-related research. This distribution indicates a concerning imbalance, where common conditions receive disproportionate attention while more complex disorders, such as bipolar disorder, remain understudied. The analysis of psychotherapy theories in panel (c) uncovers another critical gap in the field. Only 32.8% of the studies incorporate psychotherapy theories in their methodological approach. Notably, emerging therapeutic frameworks, such as humanistic therapy, are particularly under-represented in current research applications.

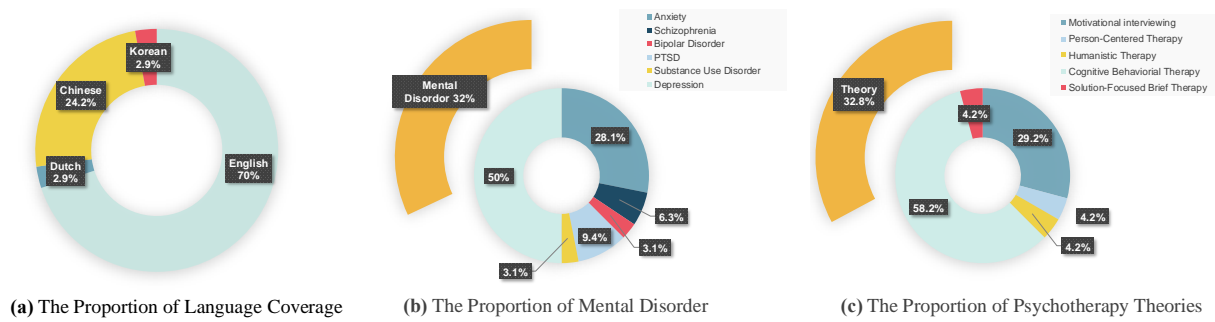


Figure 3: Distribution analysis of the current landscape.

5 Future Directions

Integrative Psychotherapy Framework. While many existing studies focus on a single dimension of psychotherapy, real-world practice involves a continuous process that spans assessment, diagnosis, and intervention (Waszczuk et al., 2017). Moreover, these stages typically unfold over multiple sessions, necessitating iterative, multi-turn interactions that incorporate the evolving context of each patient. Future work could therefore aim to develop an end-to-end conversational framework that seamlessly spans from initial evaluation to personalized intervention. By maintaining a system grounded on ongoing, context-sensitive engagement, models could dynamically update assessments and diagnoses over time, ultimately providing more responsive and individualized care.

Addressing Evolving and Multifaceted Nature of Psychotherapy. Psychotherapy commonly involves shifting symptoms, comorbidities, and nuanced patient experiences, making static or single-label predictions insufficient. Models should integrate multi-label and temporal data to capture how symptoms and emotional states evolve, while avoiding the pitfalls of incomplete symptom detection. For instance, focusing solely on the depressive features of a bipolar patient could lead to an inaccurate diagnosis if the manic phase is overlooked (Lee et al., 2024a). Furthermore, current research suggests that LLMs often struggle with multi-label tasks (hee So et al., 2024; Mohammadi et al., 2024), highlighting the need for improved model architectures and algorithms that can better account for these complexities.

Resource Infrastructure and Open-Source Tools. Current research heavily relies on commercial closed-source models, lacking reproducible open-source evaluation methods and multilingual data. Notably, developing multilingual datasets should

not solely rely on translating English resources, as psychological research indicates that cultural context plays a critical role in mental health. English-based translations cannot fully substitute for culturally specific data from other languages (Watters, 2010; Abdelkadir et al., 2024).

Broadening Scope of Disorders and Therapeutic Approaches. Most studies to date have concentrated on prevalent conditions such as depression and anxiety, leaving complex or less common disorders underexplored. Additionally, research tends to focus on a limited range of therapeutic modalities—primarily cognitive behavioral therapy. Future work could broaden both the range of disorders and the variety of therapeutic approaches, such as humanistic (Schneider and Krug, 2010) and dialectical behavior therapy (Lynch et al., 2006), to better reflect clinical realities (Norcross et al., 2022). Such an expansion could deepen the theoretical underpinnings of LLM-based psychotherapy tools and enhance the quality and relevance of digital interventions.

6 Conclusion

LLMs hold significant promise for revolutionizing psychotherapy by enhancing assessment, diagnosis, and treatment processes through dynamic, context-sensitive interactions. Despite the progress made, key challenges such as linguistic biases, limited disorder coverage, and underrepresented therapeutic models persist. Future research should focus on creating integrative, multi-turn systems that span the entire psychotherapy process while addressing the evolving nature of mental health conditions. Expanding resources, embracing diverse therapeutic approaches, and improving model architectures will be crucial in making LLM-driven psychotherapy tools more effective, inclusive, and adaptable.

641 **Limitations**

642 This survey paper, while comprehensive for LLMs
643 in psychotherapy, has several limitations: 1) The
644 studies reviewed primarily focus on the applica-
645 tion of LLMs in psychotherapy, and there may be
646 relevant research in adjacent fields or interdis-
647 ciplinary domains that was not included. 2) Due to
648 the rapidly evolving nature of this area, some re-
649 cent advancements may not be captured. The scope
650 of this survey is limited to the available literature
651 and may overlook emerging trends or unpublished
652 findings. 3) The review primarily examines studies
653 in English, which could introduce a bias towards
654 research from English-speaking countries, poten-
655 tially overlooking important cultural perspectives.
656 4) While we provide a taxonomy of LLM appli-
657 cations in psychotherapy, this framework may not
658 fully encompass the complexity of real-world clin-
659 ical settings or the diverse range of therapeutic
660 approaches currently in practice.

661 **References**

662 Nuredin Ali Abdelkadir, Charles Zhang, Ned Mayo, and
663 Stevie Chancellor. 2024. [Diverse perspectives, divergent models: Cross-cultural evaluation of depression detection on Twitter](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers)*, pages 672–680, Mexico City, Mexico. Association for Computational Linguistics.

664 American Psychiatric Association. 2006. Evidence-
665 based practice in psychology. *The American Psychologist*, 61(4):271–285.

666 American Psychiatric Association. 2022. *Diagnostic and Statistical Manual of Mental Disorders*, 5th ed. text revision edition. American Psychiatric Publishing, Arlington, VA.

667 Mario Aragon, Javier Parapar, and David E Losada. 2024. [Delving into the depths: Evaluating depression severity through BDI-biased summaries](#). In *Proceedings of the 9th Workshop on Computational Linguistics and Clinical Psychology (CLPsych 2024)*, pages 12–22, St. Julians, Malta. Association for Computational Linguistics.

668 Mihael Arcan, David-Paul Niland, and Fionn Delahunty. 2024. [An assessment on comprehending mental health through large language models](#). *Preprint*, arXiv:2401.04592.

669 Michael Barkham, William B Stiles, and David A
670 Shapiro. 1993. [The shape of change in psychotherapy: longitudinal assessment of personal problems](#). *Journal of consulting and clinical psychology*, 61(4):667.

694 David H Barlow. 2021. *Clinical handbook of psycho-
695 logical disorders: A step-by-step treatment manual*.
696 Guilford publications.

697 Andrew Brown, Jiading Zhu, Mohamed Abdelwahab,
698 Alec Dong, Cindy Wang, and Jonathan Rose. 2024. [Generation, distillation and evaluation of motivational interviewing-style reflections with a foundational language model](#). In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1241–1252, St. Julian’s, Malta. Association for Computational Linguistics.

706 Alicja Chaszczewicz, Raj Shah, Ryan Louie, Bruce
707 Arnow, Robert Kraut, and Diyi Yang. 2024. [Multi-level feedback generation with large language models for empowering novice peer counselors](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4130–4161, Bangkok, Thailand. Association for Computational Linguistics.

714 Siyuan Chen, Cong Ming, Zhiling Zhang, Yanyi Chen,
715 Kenny Q. Zhu, and Mengyue Wu. 2024. [Mixed chain-of-psychotherapies for emotional support chatbot](#). *Preprint*, arXiv:2409.19533.

718 Siyuan Chen, Mengyue Wu, Kenny Q. Zhu, Kun-
719 yao Lan, Zhiling Zhang, and Lyuchun Cui. 2023a. [Llm-empowered chatbots for psychiatrist and patient simulation: Application and evaluation](#). *Preprint*, arXiv:2305.13614.

723 Yirong Chen, Xiaofen Xing, Jingkai Lin, Huimin Zheng,
724 Zhenyu Wang, Qi Liu, and Xiangmin Xu. 2023b. [SoulChat: Improving LLMs’ empathy, listening, and comfort abilities through fine-tuning with multi-turn empathy conversations](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 1170–1183, Singapore. Association for Computational Linguistics.

731 Zhiyu Chen, Yujie Lu, and William Wang. 2023c. [Empowering psychotherapy with large language models: Cognitive distortion detection through diagnosis of thought prompting](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 4295–4304, Singapore. Association for Computational Linguistics.

738 Yu Ying Chiu, Ashish Sharma, Inna Wanyin Lin, and
739 Tim Althoff. 2024. [A computational framework for behavioral assessment of llm therapists](#). *Preprint*, arXiv:2401.00820.

742 Ben Cohen, Moreah Zisquit, Stav Yosef, Doron Fried-
743 man, and Kfir Bar. 2024. [Motivational interviewing transcripts annotated with global scores](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 11642–11657, Torino, Italia. ELRA and ICCL.

749 Ronald Jay Cohen, Mark E Swerdlik, and Suzanne M
750 Phillips. 1996. *Psychological testing and assessment*: 750

| | | |
|-----|--|-----|
| 751 | <i>An introduction to tests and measurement.</i> Mayfield Publishing Co. | 807 |
| 752 | | 808 |
| 753 | Nick Craddock and Laurence Mynors-Wallis. 2014. Psychiatric diagnosis: impersonal, imperfect and important. <i>The British Journal of Psychiatry</i> , 204(2):93–95. | 809 |
| 754 | | 810 |
| 755 | | 811 |
| 756 | | 812 |
| 757 | Rachel Elvins and Jonathan Green. 2008. The conceptualization and measurement of therapeutic alliance: An empirical review. <i>Clinical psychology review</i> , 28(7):1167–1187. | 813 |
| 758 | | 814 |
| 759 | | 815 |
| 760 | | 816 |
| 761 | Iony D Ezawa and Steven D Hollon. 2023. Cognitive restructuring and psychotherapy outcome: A meta-analytic review. <i>Psychotherapy</i> , 60(3):396. | 817 |
| 762 | | 818 |
| 763 | | |
| 764 | Isaac R. Galatzer-Levy, Daniel McDuff, Vivek Natara- jan, Alan Karthikesalingam, and Matteo Malgar- oli. 2023. The capability of large language mod- els to measure psychiatric functioning. <i>Preprint</i> , arXiv:2308.01834. | 819 |
| 765 | | 820 |
| 766 | | 821 |
| 767 | | 822 |
| 768 | | 823 |
| 769 | Sujatha Gollapalli, Beng Ang, and See-Kiong Ng. 2023. Identifying Early Maladaptive Schemas from mental health question texts. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 11832–11843, Singapore. Association for Computational Linguistics. | 824 |
| 770 | | 825 |
| 771 | | 826 |
| 772 | | 827 |
| 773 | | 828 |
| 774 | | 829 |
| 775 | Jonathan Gratch, Ron Artstein, Gale Lucas, Giota Stra- tou, Stefan Scherer, Angela Nazarian, Rachel Wood, Jill Boberg, David DeVault, Stacy Marsella, David Traum, Skip Rizzo, and Louis-Philippe Morency. 2014. The distress analysis interview corpus of human and computer interviews. In <i>Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)</i> , pages 3123– 3128, Reykjavik, Iceland. European Language Re- sources Association (ELRA). | 830 |
| 776 | | 831 |
| 777 | | 832 |
| 778 | | 833 |
| 779 | | 834 |
| 780 | | 835 |
| 781 | | 836 |
| 782 | | 837 |
| 783 | | 838 |
| 784 | | 839 |
| 785 | Gary Groth-Marnat. 2009. <i>Handbook of psychological assessment.</i> John Wiley & Sons. | 840 |
| 786 | | 841 |
| 787 | | 842 |
| 788 | | 843 |
| 789 | Loitongbam Gyanendro Singh, Junyu Mao, Rudra Mu- talik, and Stuart E. Middleton. 2024. Extracting and summarizing evidence of suicidal ideation in social media contents using large language models. In <i>Pro- ceedings of the 9th Workshop on Computational Lin- guistics and Clinical Psychology (CLPsych 2024)</i> , pages 218–226, St. Julians, Malta. Association for Computational Linguistics. | 844 |
| 790 | | 845 |
| 791 | | 846 |
| 792 | | 847 |
| 793 | | 848 |
| 794 | | 849 |
| 795 | Jae hee So, Joonhwan Chang, Eunji Kim, Junho Na, JiYeon Choi, Jy yong Sohn, Byung-Hoon Kim, and Sang Hui Chu. 2024. Aligning large language mod- els for enhancing psychiatric interviews through symptom delineation and summarization. <i>Preprint</i> , arXiv:2403.17428. | 850 |
| 796 | | 851 |
| 797 | | 852 |
| 798 | | 853 |
| 799 | | 854 |
| 800 | | 855 |
| 801 | Amey Hengle, Atharva Kulkarni, Shantanu Patankar, Madhumitha Chandrasekaran, Sneha D'Silva, Jemima Jacob, and Rashmi Gupta. 2024. Still not quite there! evaluating large language models for comorbid mental health diagnosis. <i>Preprint</i> , arXiv:2410.03908. | 856 |
| 802 | | 857 |
| 803 | | 858 |
| 804 | | 859 |
| 805 | | 860 |
| 806 | | 861 |
| | | 862 |
| | Van Hoang, Eoin Rogers, and Robert Ross. 2024. How can client motivational language inform psychother- apy agents? In <i>Proceedings of the 9th Workshop on Computational Linguistics and Clinical Psychol- ogy (CLPsych 2024)</i> , pages 23–40, St. Julians, Malta. Association for Computational Linguistics. | 807 |
| | | 808 |
| | | 809 |
| | | 810 |
| | | 811 |
| | | 812 |
| | Zainab Iftikhar, Sean Ransom, Amy Xiao, and Jeff Huang. 2024. Therapy as an nlp task: Psycholo- gists' comparison of llms and human peers in cbt. <i>Preprint</i> , arXiv:2409.02244. | 813 |
| | | 814 |
| | | 815 |
| | | 816 |
| | Gordon Jackson-Koku. 2016. Beck depression inven- tory. <i>Occupational medicine</i> , 66(2):174–175. | 817 |
| | | 818 |
| | Amanda Jensen-Doss and Kristin M Hawley. 2011. Un- derstanding clinicians' diagnostic practices: Atti- tudes toward the utility of diagnosis and standard- ized diagnostic tools. <i>Administration and Policy in Mental Health and Mental Health Services Research</i> , 38:476–485. | 819 |
| | | 820 |
| | | 821 |
| | | 822 |
| | | 823 |
| | | 824 |
| | Meng Jiang, Yi Jing Yu, Qing Zhao, Jianqiang Li, Changwei Song, Hongzhi Qi, Wei Zhai, Dan Luo, Xiaoqin Wang, Guanghui Fu, and Bing Xiang Yang. 2024a. Ai-enhanced cognitive behavioral therapy: Deep learning and large language models for ex- tracting cognitive pathways from social media texts. <i>Preprint</i> , arXiv:2404.11449. | 825 |
| | | 826 |
| | | 827 |
| | | 828 |
| | | 829 |
| | | 830 |
| | | 831 |
| | Yi Jiang, Qingyang Shen, Shuzhong Lai, Shunyu Qi, Qian Zheng, Lin Yao, Yueming Wang, and Gang Pan. 2024b. Copiloting diagnosis of autism in real clinical scenarios via llms. <i>Preprint</i> , arXiv:2410.05684. | 832 |
| | | 833 |
| | | 834 |
| | | 835 |
| | Robert M Kaplan and Dennis P Saccuzzo. 2001. <i>Psy- chological testing: Principles, applications, and is- sues.</i> Wadsworth/Thomson Learning. | 836 |
| | | 837 |
| | | 838 |
| | Shadia Kawa and James Giordano. 2012. A brief his- toricity of the diagnostic and statistical manual of mental disorders: Issues and implications for the fu- ture of psychiatric canon and practice. <i>Philosophy, Ethics, and Humanities in Medicine</i> , 7(1):2. | 839 |
| | | 840 |
| | | 841 |
| | | 842 |
| | | 843 |
| | Taewan Kim, Seolyeong Bae, Hyun Ah Kim, Su-Woo Lee, Hwajung Hong, Chanmo Yang, and Young-Ho Kim. 2024. Mindfuldiary: Harnessing large language model to support psychiatric patients' journaling. In <i>Proceedings of the CHI Conference on Human Fac- tors in Computing Systems, CHI '24</i> , New York, NY, USA. Association for Computing Machinery. | 844 |
| | | 845 |
| | | 846 |
| | | 847 |
| | | 848 |
| | | 849 |
| | | 850 |
| | Gleb Kuzmin, Petr Strepetov, Maksim Stankevich, Artem Shelmanov, and Ivan Smirnov. 2024. Men- tal disorders detection in the era of large language models. <i>Preprint</i> , arXiv:2410.07129. | 851 |
| | | 852 |
| | | 853 |
| | | 854 |
| | Michael J Lambert. 2013. <i>Bergin and Garfield's hand- book of psychotherapy and behavior change.</i> John Wiley & Sons. | 855 |
| | | 856 |
| | | 857 |
| | Kunyao Lan, Bingrui Jin, Zichen Zhu, Siyuan Chen, Shu Zhang, Kenny Q. Zhu, and Mengyue Wu. 2024a. Depression diagnosis dialogue simulation: Self-improving psychiatrist with tertiary memory. <i>Preprint</i> , arXiv:2409.15084. | 858 |
| | | 859 |
| | | 860 |
| | | 861 |
| | | 862 |

| | | |
|-----|---|-----|
| 863 | Kunyao Lan, Cong Ming, Binwei Yao, Lu Chen, and Mengyue Wu. 2024b. Towards reliable and empathetic depression-diagnosis-oriented chats . <i>Preprint</i> , arXiv:2404.05012. | 920 |
| 864 | | 921 |
| 865 | | |
| 866 | | |
| 867 | Xiaochong Lan, Yiming Cheng, Li Sheng, Chen Gao, and Yong Li. 2024c. Depression detection on social media with large language models . <i>Preprint</i> , arXiv:2403.10750. | 922 |
| 868 | | 923 |
| 869 | | 924 |
| 870 | | 925 |
| 871 | Daeun Lee, Hyolim Jeon, Sejung Son, Chaewon Park, Ji hyun An, Seungbae Kim, and Jinyoung Han. 2024a. Detecting bipolar disorder from misdiagnosed major depressive disorder with mood-aware multi-task learning . In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 4954–4970, Mexico City, Mexico. Association for Computational Linguistics. | 926 |
| 872 | | 927 |
| 873 | | 928 |
| 874 | | |
| 875 | | |
| 876 | | |
| 877 | | |
| 878 | | |
| 879 | | |
| 880 | | |
| 881 | Suyeon Lee, Jieun Kang, Harim Kim, Kyoung-Mee Chung, Dongha Lee, and Jinyoung Yeo. 2024b. Co-coa: Cbt-based conversational counseling agent using memory specialized in cognitive distortions and dynamic prompt . <i>Preprint</i> , arXiv:2402.17546. | 929 |
| 882 | | 930 |
| 883 | | 931 |
| 884 | | 932 |
| 885 | | |
| 886 | Suyeon Lee, Sunghwan Kim, Minju Kim, Dongjin Kang, Dongil Yang, Harim Kim, Minseok Kang, Dayi Jung, Min Hee Kim, Seungbeen Lee, Kyoung-Mee Chung, Youngjae Yu, Dongha Lee, and Jinyoung Yeo. 2024c. Cactus: Towards psychological counseling conversations using cognitive behavioral theory . <i>Preprint</i> , arXiv:2407.03103. | 933 |
| 887 | | 934 |
| 888 | | 935 |
| 889 | | 936 |
| 890 | | 937 |
| 891 | | 938 |
| 892 | | 939 |
| 893 | Yeonji Lee, Sangjun Park, Kyunghyun Cho, and JinYeong Bak. 2024d. Mentalagora: A gateway to advanced personalized care in mental health through multi-agent debating and attribute control . <i>Preprint</i> , arXiv:2407.02736. | 940 |
| 894 | | |
| 895 | | |
| 896 | | |
| 897 | | |
| 898 | Younghun Lee, Dan Goldwasser, and Laura Schwab Reese. 2024e. Towards understanding counseling conversations: Domain knowledge and large language models . In <i>Findings of the Association for Computational Linguistics: EACL 2024</i> , pages 2032–2047, St. Julian’s, Malta. Association for Computational Linguistics. | 941 |
| 899 | | 942 |
| 900 | | 943 |
| 901 | | 944 |
| 902 | | 945 |
| 903 | | |
| 904 | | |
| 905 | Anqi Li, Yu Lu, Nirui Song, Shuai Zhang, Lizhi Ma, and Zhenzhong Lan. 2024a. Automatic evaluation for mental health counseling using llms . <i>Preprint</i> , arXiv:2402.11958. | 946 |
| 906 | | 947 |
| 907 | | 948 |
| 908 | | 949 |
| 909 | Ruosun Li, Zimu Wang, Son Tran, Lei Xia, and Xinya Du. 2024b. Meqa: A benchmark for multi-hop event-centric question answering with explanations . In <i>Advances in Neural Information Processing Systems</i> , volume 37, pages 126835–126862. Curran Associates, Inc. | 950 |
| 910 | | 951 |
| 911 | | 952 |
| 912 | | |
| 913 | | |
| 914 | | |
| 915 | Wenyu Li, Yinuo Zhu, Xin Lin, Ming Li, Ziyue Jiang, and Ziqian Zeng. 2024c. Zero-shot explainable mental health analysis on social media by incorporating mental scales . In <i>Companion Proceedings of the ACM on Web Conference 2024, WWW ’24</i> , page 959–962, New York, NY, USA. Association for Computing Machinery. | 953 |
| 916 | | 954 |
| 917 | | 955 |
| 918 | | 956 |
| 919 | | 957 |
| | | 958 |
| | | 959 |
| | | 960 |
| | | 961 |
| | | 962 |
| | | 963 |
| | | 964 |
| | | 965 |
| | | 966 |
| | | 967 |
| | | 968 |
| | | 969 |
| | | 970 |
| | | 971 |
| | | 972 |
| | | 973 |
| | | 974 |
| | | 975 |

| | | | |
|------|--|---|------|
| 976 | | | |
| 977 | | | |
| 978 | | <i>the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)</i> , pages 2930–2940, | |
| 979 | | Torino, Italia. ELRA and ICCL. | |
| 980 | Hongbin Na, Tao Shen, Shumao Yu, and Ling | | |
| 981 | Chen. 2024. Multi-session client-centered treatment outcome evaluation in psychotherapy . <i>Preprint</i> , | | |
| 982 | arXiv:2410.05824. | | |
| 983 | | | |
| 984 | Vivian Nguyen, Sang Min Jung, Lillian Lee, Thomas D. | | |
| 985 | Hull, and Cristian Danescu-Niculescu-Mizil. 2024. | | |
| 986 | Taking a turn for the better: Conversation redirection throughout the course of mental-health therapy . | | |
| 987 | <i>Preprint</i> , arXiv:2410.07147. | | |
| 988 | | | |
| 989 | Jingping Nie, Hanya Shao, Yuang Fan, Qijia Shao, | | |
| 990 | Haoxuan You, Matthias Preindl, and Xiaofan Jiang. | | |
| 991 | 2024. Llm-based conversational ai therapist for daily functioning screening and psychotherapeutic intervention via everyday smart devices . <i>Preprint</i> , | | |
| 992 | arXiv:2403.10779. | | |
| 993 | | | |
| 994 | | | |
| 995 | John C Norcross, Rory A Pfund, and Danielle M Cook. | | |
| 996 | 2022. The predicted future of psychotherapy: A | | |
| 997 | decennial e-delphi poll. <i>Professional Psychology: Research and Practice</i> , 53(2):109. | | |
| 998 | | | |
| 999 | Michael Robert Phillips, Qijie Shen, Xiehe Liu, Sonya | | |
| 1000 | Pritzker, David Streiner, Ken Conner, and Gonghuan | | |
| 1001 | Yang. 2007. Assessing depressive symptoms in persons who die of suicide in mainland china . <i>Journal of Affective Disorders</i> , 98(1-2):73–82. | | |
| 1002 | | | |
| 1003 | | | |
| 1004 | James O Prochaska and John C Norcross. 2018. <i>Systems of psychotherapy: A transtheoretical analysis</i> . | | |
| 1005 | Oxford University Press. | | |
| 1006 | | | |
| 1007 | Hongzhi Qi, Qing Zhao, Jianqiang Li, Changwei Song, | | |
| 1008 | Wei Zhai, Dan Luo, Shuo Liu, Yi Jing Yu, Fan Wang, | | |
| 1009 | Huijing Zou, Bing Xiang Yang, and Guanghui Fu. | | |
| 1010 | 2024. Supervised learning and large language model benchmarks on mental health datasets: Cognitive distortions and suicidal risks in chinese social media . | | |
| 1011 | <i>Preprint</i> , arXiv:2309.03564. | | |
| 1012 | | | |
| 1013 | | | |
| 1014 | Huachuan Qiu, Lizhi Ma, and Zhenzhong Lan. 2024. | | |
| 1015 | Psyguard: An automated system for suicide detection and risk assessment in psychological counseling . | | |
| 1016 | <i>Preprint</i> , arXiv:2409.20243. | | |
| 1017 | | | |
| 1018 | Nishat Raihan, Sadiya Sayara Chowdhury Puspo, | | |
| 1019 | Shafkat Farabi, Ana-Maria Bucur, Tharindu Ranas- | | |
| 1020 | inghe, and Marcos Zampieri. 2024. MentalHelp: A multi-task dataset for mental health in social media . In <i>Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)</i> , | | |
| 1021 | pages 11196–11203, Torino, Italia. ELRA and ICCL. | | |
| 1022 | | | |
| 1023 | | | |
| 1024 | | | |
| 1025 | | | |
| 1026 | Chenyu Ren, Yazhou Zhang, Daihai He, and Jing Qin. | | |
| 1027 | 2024. Wundtgppt: Shaping large language models to be an empathetic, proactive psychologist . <i>Preprint</i> , | | |
| 1028 | arXiv:2406.15474. | | |
| 1029 | | | |
| | | Joseph Rencic, Steven J Durning, Eric Holmboe, and | 1030 |
| | | Larry D Gruppen. 2016. Understanding the assess- | 1031 |
| | | ment of clinical reasoning. <i>Assessing competence</i> | 1032 |
| | | <i>in professional performance across disciplines and</i> | 1033 |
| | | <i>professions</i> , pages 209–235. | 1034 |
| | | | |
| | | Günter Schiepek, Wolfgang Aichhorn, Martin Gruber, | 1035 |
| | | Guido Strunk, Egon Bachler, and Benjamin Aas. | 1036 |
| | | 2016. Real-time monitoring of psychotherapeutic | 1037 |
| | | processes: concept and compliance. <i>Frontiers in</i> | 1038 |
| | | <i>psychology</i> , 7:604. | 1039 |
| | | | |
| | | Miriam Schirmer, Tobias Leemann, Gjergji Kasneci, | 1040 |
| | | Jürgen Pfeffer, and David Jurgens. 2024. The language of trauma: Modeling traumatic event descriptions across domains with explainable ai . <i>Preprint</i> , | 1041 |
| | | arXiv:2408.05977. | 1042 |
| | | | 1043 |
| | | | 1044 |
| | | Kirk J Schneider and Orah T Krug. 2010. <i>Existential-</i> | 1045 |
| | | <i>humanistic therapy</i> . American Psychological Associ- | 1046 |
| | | ation Washington, DC. | 1047 |
| | | | |
| | | Natalie Shapira and Tal Alfi-Yogev. 2024. Therapist self-disclosure as a natural language processing task . | 1048 |
| | | In <i>Proceedings of the 9th Workshop on Computational Linguistics and Clinical Psychology (CLPsych 2024)</i> , | 1049 |
| | | pages 61–73, St. Julians, Malta. Association | 1050 |
| | | for Computational Linguistics. | 1051 |
| | | | 1052 |
| | | | 1053 |
| | | Ashish Sharma, Kevin Rushton, Inna Lin, David Wad- | 1054 |
| | | den, Khendra Lucas, Adam Miner, Theresa Nguyen, | 1055 |
| | | and Tim Althoff. 2023. Cognitive reframing of negative thoughts through human-language model interaction . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , | 1056 |
| | | pages 9977–10000, Toronto, | 1057 |
| | | Canada. Association for Computational Linguistics. | 1058 |
| | | | 1059 |
| | | | 1060 |
| | | | 1061 |
| | | Ashish Sharma, Kevin Rushton, Inna Wanyin Lin, | 1062 |
| | | Theresa Nguyen, and Tim Althoff. 2024. Facilitating self-guided mental health interventions through human-language model interaction: A case study of cognitive restructuring . In <i>Proceedings of the CHI Conference on Human Factors in Computing Systems, CHI '24</i> , New York, NY, USA. Association for | 1063 |
| | | Computing Machinery. | 1064 |
| | | | 1065 |
| | | | 1066 |
| | | | 1067 |
| | | | 1068 |
| | | | 1069 |
| | | Konstantinos Skianis, John Pavlopoulos, and A. Seza | 1070 |
| | | Doğruöz. 2024. Severity prediction in mental health: Llm-based creation, analysis, evaluation of a novel multilingual dataset . <i>Preprint</i> , arXiv:2409.17397. | 1071 |
| | | | 1072 |
| | | | 1073 |
| | | Eliseo Bao Souto, Anxo Pérez, and Javier Parapar. 2023. | 1074 |
| | | Explainable depression symptom detection in social media . <i>Preprint</i> , arXiv:2310.13664. | 1075 |
| | | | 1076 |
| | | Aseem Srivastava, Smriti Joshi, Tanmoy Chakraborty, | 1077 |
| | | and Md Shad Akhtar. 2024. Knowledge planning in large language models for domain-aligned counseling summarization . <i>Preprint</i> , arXiv:2409.14907. | 1078 |
| | | | 1079 |
| | | | 1080 |
| | | Dorothy E Stubbe. 2018. The therapeutic alliance: | 1081 |
| | | The fundamental element of psychotherapy. <i>Focus</i> , | 1082 |
| | | 16(4):402–403. | 1083 |

| | | | |
|------|---|--|------|
| 1084 | Xin Sun, Jiahuan Pei, Jan de Wit, Mohammad Alian- | <i>Psychology (CLPsych 2024)</i> , pages 108–126, St. | 1141 |
| 1085 | nejadi, Emiel Kraemer, Jos T.P. Dobber, and Jos A. | Julians, Malta. Association for Computational | 1142 |
| 1086 | Bosch. 2024. Eliciting motivational interviewing | Linguistics. | 1143 |
| 1087 | skill codes in psychotherapy with LLMs: A bilin- | | |
| 1088 | gual dataset and analytical study . In <i>Proceedings of</i> | Zimu Wang, Lei Xia, Wei Wang, and Xinya Du. | 1144 |
| 1089 | <i>the 2024 Joint International Conference on Compu-</i> | 2024e. Document-level causal relation extraction | 1145 |
| 1090 | <i>tational Linguistics, Language Resources and Eval-</i> | with knowledge-guided binary question answering . | 1146 |
| 1091 | <i>uation (LREC-COLING 2024)</i> , pages 5609–5621, | In <i>Findings of the Association for Computational</i> | 1147 |
| 1092 | Torino, Italia. ELRA and ICCL. | <i>Linguistics: EMNLP 2024</i> , pages 16944–16955, Mi- | 1148 |
| | | ami, Florida, USA. Association for Computational | 1149 |
| 1093 | Alan R. Teo. 2010. A new form of social withdrawal | Linguistics. | 1150 |
| 1094 | in japan: a review of hikikomori . <i>International Jour-</i> | | |
| 1095 | <i>nal of Social Psychiatry</i> , 56(2):178–185. PMID: | Monika A. Waszczuk, Mark Zimmerman, Camilo Rug- | 1151 |
| 1096 | 19567455. | gero, Kaiqiao Li, Annmarie MacNamara, Anna Wein- | 1152 |
| | | berg, Greg Hajcak, David Watson, and Roman Kotov. | 1153 |
| 1097 | Sichang Tu, Abigail Powers, Natalie Merrill, Negar | 2017. What do clinicians treat: Diagnoses or symp- | 1154 |
| 1098 | Fani, Sierra Carter, Stephen Doogan, and Jinho D. | toms? the incremental validity of a symptom-based, | 1155 |
| 1099 | Choi. 2024. Automating ptsd diagnostics in clinical | dimensional characterization of emotional disorders | 1156 |
| 1100 | interviews: Leveraging large language models for | in predicting medication prescription patterns . <i>Com-</i> | 1157 |
| 1101 | trauma assessments . <i>Preprint</i> , arXiv:2405.11178. | <i>prehensive Psychiatry</i> , 79:80–88. Advances in Trans- | 1158 |
| | | diagnostic Psychopathology Research. | 1159 |
| 1102 | Ahmet Yavuz Uluslu, Andrianos Michail, and Simon | | |
| 1103 | Clematide. 2024. Utilizing large language models to | Ethan Watters. 2010. <i>Crazy Like Us: The Globalization</i> | 1160 |
| 1104 | identify evidence of suicidality risk through analysis | <i>of the American Psyche</i> . Free Press. | 1161 |
| 1105 | of emotionally charged posts . In <i>Proceedings of</i> | | |
| 1106 | <i>the 9th Workshop on Computational Linguistics and</i> | Anuradha Welivita and Pearl Pu. 2023. Boosting dis- | 1162 |
| 1107 | <i>Clinical Psychology (CLPsych 2024)</i> , pages 264–269, | tress support dialogue responses with motivational | 1163 |
| 1108 | St. Julians, Malta. Association for Computational | interviewing strategy . In <i>Findings of the Associa-</i> | 1164 |
| 1109 | Linguistics. | <i>tion for Computational Linguistics: ACL 2023</i> , pages | 1165 |
| | | 5411–5432, Toronto, Canada. Association for Com- | 1166 |
| 1110 | Bruce E Wampold and Zac E Imel. 2015. <i>The great</i> | putational Linguistics. | 1167 |
| 1111 | <i>psychotherapy debate: The evidence for what makes</i> | | |
| 1112 | <i>psychotherapy work</i> . Routledge. | World Health Organization. 2019. International clas- | 1168 |
| 1113 | Bichen Wang, Pengfei Deng, Yanyan Zhao, and Bing | sification of diseases, eleventh revision (icd-11) . | 1169 |
| 1114 | Qin. 2023. C2D2 dataset: A resource for the cog- | Licensed under Creative Commons Attribution- | 1170 |
| 1115 | nitive distortion analysis and its impact on mental | NoDerivatives 3.0 IGO licence (CC BY-ND 3.0 | 1171 |
| 1116 | health . In <i>Findings of the Association for Compu-</i> | IGO). | 1172 |
| 1117 | <i>tational Linguistics: EMNLP 2023</i> , pages 10149– | | |
| 1118 | 10160, Singapore. Association for Computational | Mengxi Xiao, Qianqian Xie, Ziyang Kuang, Zhicheng | 1173 |
| 1119 | Linguistics. | Liu, Kailai Yang, Min Peng, Weiguang Han, and | 1174 |
| | | Jimin Huang. 2024. Healme: Harnessing cognitive | 1175 |
| 1120 | Jiashuo Wang, Yang Xiao, Yanran Li, Changhe Song, | reframing in large language models for psychother- | 1176 |
| 1121 | Chunpu Xu, Chenhao Tan, and Wenjie Li. 2024a. To- | apy . <i>Preprint</i> , arXiv:2403.05574. | 1177 |
| 1122 | wards a client-centered assessment of llm therapists | | |
| 1123 | by client simulation . <i>Preprint</i> , arXiv:2406.12266. | Xuhai Xu, Bingsheng Yao, Yuanzhe Dong, Saadia | 1178 |
| | | Gabriel, Hong Yu, James Hendler, Marzyeh Ghas- | 1179 |
| 1124 | Ruiyi Wang, Stephanie Milani, Jamie C. Chiu, Jiayin | semi, Anind K. Dey, and Dakuo Wang. 2024. Mental- | 1180 |
| 1125 | Zhi, Shaun M. Eack, Travis Labrum, Samuel M. | llm: Leveraging large language models for mental | 1181 |
| 1126 | Murphy, Nev Jones, Kate Hardy, Hong Shen, Fei | health prediction via online text data . <i>Proc. ACM</i> | 1182 |
| 1127 | Fang, and Zhiyu Zoey Chen. 2024b. PATIENT-Ψ: | <i>Interact. Mob. Wearable Ubiquitous Technol.</i> , 8(1). | 1183 |
| 1128 | Using Large Language Models to Simulate Patients | | |
| 1129 | for Training Mental Health Professionals . <i>Preprint</i> , | Kailai Yang, Shaoxiong Ji, Tianlin Zhang, Qianqian | 1184 |
| 1130 | arXiv:2405.19660. | Xie, Ziyang Kuang, and Sophia Ananiadou. 2023. To- | 1185 |
| | | wards interpretable mental health analysis with large | 1186 |
| 1131 | Yuqing Wang, Yun Zhao, Sara Alessandra Keller, Anne | language models . In <i>Proceedings of the 2023 Con-</i> | 1187 |
| 1132 | de Hond, Marieke M. van Buchem, Malvika Pillai, | <i>ference on Empirical Methods in Natural Language</i> | 1188 |
| 1133 | and Tina Hernandez-Boussard. 2024c. Unveiling and | <i>Processing</i> , pages 6056–6077, Singapore. Associa- | 1189 |
| 1134 | mitigating bias in mental health analysis with large | <i>tion for Computational Linguistics</i> . | 1190 |
| 1135 | language models . <i>Preprint</i> , arXiv:2406.12033. | | |
| | | Kailai Yang, Tianlin Zhang, Ziyang Kuang, Qianqian Xie, | 1191 |
| 1136 | Yuxi Wang, Diana Inkpen, and Prasadith | Jimin Huang, and Sophia Ananiadou. 2024a. Mental- | 1192 |
| 1137 | Kirinde Gamaarachchige. 2024d. Explainable | lama: Interpretable mental health analysis on social | 1193 |
| 1138 | depression detection using large language models | media with large language models . In <i>Proceedings</i> | 1194 |
| 1139 | on social media data . In <i>Proceedings of the 9th</i> | <i>of the ACM on Web Conference 2024</i> , WWW '24, | 1195 |
| 1140 | <i>Workshop on Computational Linguistics and Clinical</i> | page 4489–4500, New York, NY, USA. Association | 1196 |
| | | for Computing Machinery. | 1197 |

- 1198 Qisen Yang, Zekun Wang, Honghui Chen, Shenzhi
1199 Wang, Yifan Pu, Xin Gao, Wenhao Huang, Shiji
1200 Song, and Gao Huang. 2024b. [Llm agents for psy-](#)
1201 [chology: A study on gamified assessments](#). *Preprint*,
1202 [arXiv:2402.12326](#).
- 1203 Stav Yosef, Moreah Zisquit, Ben Cohen, Anat
1204 Klomek Brunstein, Kfir Bar, and Doron Friedman.
1205 2024. [Assessing motivational interviewing sessions](#)
1206 [with AI-generated patient simulations](#). In *Proceed-*
1207 *ings of the 9th Workshop on Computational Linguis-*
1208 *tics and Clinical Psychology (CLPsych 2024)*, pages
1209 1–11, St. Julians, Malta. Association for Computa-
1210 tional Linguistics.
- 1211 Hongli Zhan, Allen Zheng, Yoon Kyung Lee, Jina Suh,
1212 Junyi Jessy Li, and Desmond C. Ong. 2024. [Large](#)
1213 [language models are capable of offering cognitive](#)
1214 [reappraisal, if guided](#). *Preprint*, [arXiv:2404.01288](#).
- 1215 Chenhao Zhang, Renhao Li, Minghuan Tan, Min Yang,
1216 Jingwei Zhu, Di Yang, Jiahao Zhao, Guancheng Ye,
1217 Chengming Li, and Xiping Hu. 2024a. [Cpsycoun:](#)
1218 [A report-based multi-turn dialogue reconstruction](#)
1219 [and evaluation framework for chinese psychological](#)
1220 [counseling](#). *Preprint*, [arXiv:2405.16433](#).
- 1221 Mian Zhang, Xianjun Yang, Xinlu Zhang, Travis
1222 Labrum, Jamie C. Chiu, Shaun M. Eack, Fei Fang,
1223 William Yang Wang, and Zhiyu Zoey Chen. 2024b.
1224 [Cbt-bench: Evaluating large language models on](#)
1225 [assisting cognitive behavior therapy](#). *Preprint*,
1226 [arXiv:2410.13218](#).