# MIRROR: Multimodal Cognitive Reframing Therapy for Rolling with Resistance

**Important:** We explore how vision-language models support digital CBT, but they should NOT replace professional psychological treatment.

Anonymous ACL submission

### Abstract

Recent studies have explored the use of large language models (LLMs) in psychotherapy; however, text-based cognitive behavioral therapy (CBT) models often struggle with client resistance, which can weaken therapeutic alliance. To address this, we propose a multimodal approach that incorporates nonverbal cues, which allows the AI therapist to better align its responses with the client's negative emotional state. Specifically, we introduce a new synthetic dataset, MIRROR (Multimodal Interactive Rolling with Resistance), which is a novel synthetic dataset that pairs each client's statements with corresponding facial images. Using this dataset, we train baseline vision language models (VLMs) so that they can analyze facial cues, infer emotions, and generate empathetic responses to effectively manage client resistance. These models are then evaluated in terms of both their counseling skills as a therapist, and the strength of therapeutic alliance in the presence of client resistance. Our results demonstrate that MIRROR significantly enhances the AI therapist's ability to handle resistance, which outperforms existing textbased CBT approaches. Human expert evaluations further confirm the effectiveness of our approach in managing client resistance and fostering therapeutic alliance.

# 1 Introduction

800

011

012

014

018

040

043

Cognitive reframing is a central part of cognitive behavioral therapy (CBT), which helps individuals replace negative and intrusive thoughts with more rational and balanced ones. Towards this objective, large language models (LLMs) have recently shown great promise and are increasingly being explored in psychotherapy (Ziems et al., 2022; Maddela et al., 2023a; Sharma et al., 2023; Qu et al., 2023; Yang et al., 2023, 2024; Xiao et al., 2024; Na, 2024; Lee et al., 2024a). As such, these systems have actually been utilized in real-world applications as effective adjunct tools in psychotherapy,



Figure 1: Text-based therapists have limitations in interpreting nonverbal cues, as they cannot perceive behaviors such as sighs or posture shifts, which can lead to premature problem-solving rather than addressing deeper emotions.

providing meaningful support for individuals with mental disorders such as depression and anxiety (Fitzpatrick et al., 2017; Haque and Rubya, 2023; Mehta et al., 2021)<sup>1</sup>.

045

046

047

053

060

061

Despite this progress, the existing text-based CBT model struggles to detect and respond to client resistance (Wang et al., 2025), which is a common therapeutic challenge that involves the client's reluctance or opposition to change. This resistance often stems from the directive nature of CBT, where structured interventions may unintentionally provoke discomfort or defensiveness (Patterson and Chamberlain, 1994; Moyers and Martin, 2006; Constantino et al., 2017; Westra and Norouzian, 2018; Hara, 2020). Left unaddressed, resistance can diminish therapeutic alliance and reduce treatment efficacy. It is crucial to note that such resistance is frequently conveyed through non-

<sup>&</sup>lt;sup>1</sup>A comprehensive review of related work is provided in Appendix B.

verbal cues like facial expressions, sighs, or posture shifts. Due to this property, pure-text-based models fail to perceive resistance, which leads to premature advice-giving rather than addressing deeper emotional needs (Figure 1). Addressing this limitation thus requires multimodal integration. However, collecting real multimodal psychotherapy data to train models to identify such multimodal cues, introduces severe privacy risks as sessions often involve deeply personal disclosures, including trauma, mental illness, and other confidential experiences.

062

063

064

067

097

100

102

103

104

107

108

109

110

111

112

In this work, we propose a multimodal approach to cognitive reframing that integrates both textual and nonverbal information to better detect and manage client resistance. We introduce MIRROR (Multimodal Interactive Rolling with Resistance), which is a synthetic dataset designed to simulate real therapeutic interactions. Specifically, MIRROR features generated dialogues between clients and therapists, annotated with client facial expressions reflecting three distinct types of resistance. We leverage LLMs to generate realistic session content, synthesize corresponding facial cues, and apply rigorous filtering to ensure quality and safety. This dataset enables the development of vision-language models (VLMs) tailored to CBT scenarios, where emotional alignment and alliance are essential. In addition, we propose *emotional captioning*, a novel reasoning method designed to explicitly interpret and respond to the client's emotional state.

> We evaluate our approach using a VLM that is trained on the MIRROR dataset and enhanced with *planning* and *emotional captioning*. Compared to existing LLMs and VLMs, our model demonstrates superior performance across therapist skill assessment, alliance building, and applicability to real counseling scenarios. The results highlight the importance of multimodal approaches in managing client resistance and improving CBT outcomes.

Our contributions are summarized as follows:

- We explore a multimodal cognitive reframing for coping with client resistance, and present MIRROR, which features turn-level client facial expressions across diverse resistance types.
- We establish baseline models on the MIRROR dataset and propose an *emotional captioning* method, which helps VLMs generate emotionally aligned, vision-aware therapeutic responses.

To further support research in this area, we will publicly release our code and dataset.

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

# 2 **Problem Definition**

Our goal is to enhance the AI therapist's ability to manage client resistance by integrating both verbal and nonverbal cues through a multimodal approach. To guide the development and evaluation of such models, we define two key assessment dimensions that reflect essential aspects of effective therapy:

- Therapist Skills Assessment: Evaluates the AI therapist's competence in two key categories of general counseling skills and CBT-specific techniques.
- **Client Alliance Assessment:** Focuses on the AI therapist's ability to establish a strong therapeutic bond, which is critical for reducing resistance and promoting positive outcomes.

# 3 MIRROR: Multimodal Interactive Rolling with Resistance Dataset

As illustrated in Figure 2, the MIRROR dataset is constructed through three main steps, which is followed by a comprehensive quality and safety validation process. Through dataset synthesis, we generate over 3,000 multimodal counseling dialogues, with each client turn annotated with a facial expression image that captures the client's emotional state<sup>2</sup>.

# 3.1 Step 1: Multimodal Dialogue Design

To build the multimodal dialogue design for MIR-ROR, we combine facial and textual data from two sources: CelebA (Liu et al., 2015) for facial expressions and CACTUS (Lee et al., 2024a) for textbased cognitive reframing therapy. While CACTUS is originally a text-only dialogue dataset, we only extract its underlying structured profiles, which includes client intake forms, thinking traps, counseling plans, and CBT techniques.

In order to assign a facial identity to each client, we pair every CACTUS profile with a CelebA image based on gender and age predictions from the DeepFace library (Serengil and Ozpinar, 2021).

We further augment each client profile with four distinct resistance types: cognitive, emotional, behavioral, and non-resistant, following the taxonomy proposed by Beal III et al. (2013). Rather than assigning a single resistance label to each profile, we

<sup>&</sup>lt;sup>2</sup>All used prompts are provided in Appendix J.



Figure 2: Overview of the MIRROR dataset construction. The pipeline consists of three main stages: Multimodal Dialogue Design (§3.1), Counseling Screenplay Generation (§3.2), and Facial Expression Synthesis (§3.3).

generate four variants per client, each conditioned 159 on a different resistance type. This results in four 160 variants per client, allowing the model to encounter 161 diverse resistance behaviors in the same therapeu-162 tic context and ensuring class balance across the 163 dataset. This process yields a complete multimodal 164 165 dialogue setup for each client, where a structured CACTUS profile, a facial identity, and a specified 166 resistance type are jointly configured. The resulting 167 design supports therapeutically grounded dialogue generation based on client context and CBT plan. 169

# 3.2 Step 2: Counseling Screenplay Generation

We synthesize counseling dialogues in the form of screenplays rather than plain transcripts, to more naturally reflect the emotional nuance of real therapeutic interactions. A key advantage of this format is its explicit representation of nonverbal cues through stage directions (e.g., *[slightly defensive, arms crossed]*' in Figure 2).

171

172

173

174

175

177

179

181

These stage directions serve two critical purposes: (1) They enrich the textual context by capturing subtle emotional dynamics that are characteristic of real therapy sessions. (2) They act as structured signals for downstream facial expression synthesis, which ensures the generation of consistent and emotionally aligned client images. Based on the predefined profiles, these screenplays are generated using GPT-40-MINI<sup>3</sup>. 182

183

184

185

186

188

189

190

191

192

194

195

196

197

198

199

200

202

# 3.3 Step 3: Facial Expression Synthesis

After constructing the screenplay, we synthesize turn-level facial expressions that reflect the emotional dynamics conveyed through both verbal content and stage directions. The key contribution of this step lies in designing a prompt construction method that encodes nonverbal cues into the image generation process.

We leverage **PhotoMaker** (Li et al., 2024b), which is a diffusion-based model that takes three inputs: a reference image to preserve facial identity, a positive prompt for the desired expression, and a negative prompt to suppress conflicting features. To generate these prompts, we condition LLAMA-3-8B (AI@Meta, 2024) on the full client utterance, which includes inline stage directions (see

<sup>&</sup>lt;sup>3</sup>Version gpt-40-mini-2024-07-18.

	Modality	Language	# of Dialogue	# Avg. Turns	# Avg. Images	Turn-Image Alignment
Psych8k (Liu et al., 2023)	Т	English	8,187	1.00	-	-
HealMe (Xiao et al., 2024)	Т	English	1,300	3.00	-	-
CACTUS (Lee et al., 2024a)	Т	English	31,577	16.6	-	-
CPsyCounD (Zhang et al., 2024)	Т	Chinese	3,134	8.7	-	-
M2CoSC (Kim et al., 2025)	Τ, V	English	429	4.00	1.00	×
MEDIC (Zhu et al., 2023)	T, V, A	Chinese	771	1.00	1,137	<u> </u>
MIRROR	T, V	English	3,073	10.3	9.51	$\checkmark$

Table 1: Comparison of MIRROR with other psychological counseling datasets. The **Modality** column indicates whether the dataset includes text (T), visual (V), or audio (A) data. **# Avg. Images** refers to the average number of client images per dialogue. **Turn-Image Alignment** indicates whether the client images are dynamically aligned according to each dialogue turn.

Figure 2). As a result, LLAMA-3-8B produces two facial expression descriptions: a target expression (e.g., "downcast expression with eyes looking away") and a contrasting one (e.g., "trusting expression with a gentle smile"), which populate the positive and negative prompts, respectively.

203

207

210

211

212

213

214

215

216

217

218

This approach enables the synthesis of emotionally aligned client images throughout the dialogue. As shown in Figure 2, expressions like *[looking away]* are clearly expressed in the synthesized images. By translating nonverbal cues into structured prompts, we ensure that facial expressions reflect the client's emotional state, even when the textual utterance alone does not explicitly convey it. The role of stage direction in image synthesis is further examined in Appendix I.

# **3.4** Step 4: Filtering for Quality and Safety

**Dataset Quality Filtering** To ensure the overall quality and coherence with image of multimodal counseling dialogues, we apply six filtering approaches: (1) Image-Text Similarity Filtering uses CLIP (Radford et al., 2021) to measure alignment 224 between generated images and stage directions, and cases with low similarity score below 0.2 are dis-226 carded (2.95% rejected). (2) Identity Preservation 227 Filtering employs ArcFace (Deng et al., 2019) to maintain facial similarity across dialogue turns, rejecting cases with low similarity scores below 0.3 (66.05% rejected). (3) Gender Preservation Filter-231 ing utilizes DeepFace to ensure that the detected 232 gender matches the client's original multimodal profile, removing mismatches (15.39% rejected). (4) Basic Filtering eliminates dialogues that contain utterances longer than 100 words, exhibit unnatural repetition of the same part-of-speech more than three times in a row, or include too few (fewer than 4) or too many (more than 20) conversation turns (1.03% rejected). (5) Copy-Paste Filtering 240

removes instances where client personas are unnaturally stated instead of being contextually integrated, improving dialogue realism (1.36% rejected). Lastly, (6) Therapeutic Alliance Filtering assesses the quality of the counseling interaction using GPT-40<sup>4</sup> to evaluate WAI<sup>5</sup> (Li et al., 2024a), and dialogues with an average score below 0.3 are discarded (10.01% rejected).

241

242

243

244

245

246

247

248

249

250

251

252

254

255

256

257

258

259

260

261

262

263

264

265

266

267

269

270

271

272

273

274

**Dataset Safety Filtering** To uphold ethical standards and prevent harmful content, we apply two additional approaches. (1) NSFW Filtering uses a Not-Safe-For-Work (NSFW) detector<sup>6</sup> to remove images that are visually unsuitable for mental health dialogue contexts. (2) Dialogue Safety Filtering leverages Canary (Kim et al., 2022) to identify and eliminate instances containing toxic, unethical, or unsafe language, in accordance with prior safety protocols (Kim et al., 2023; Lee et al., 2024b) (1.09% rejected). These layered filtering stages are critical for constructing a high-quality dataset that is not only realistic and coherent but also ethically robust and clinically applicable.

### **3.5** Comparative Analysis of MIRROR

Through the preceding stages, we have curated the first multimodal CBT dataset that explicitly incorporates client resistance. As shown in Table 1, MIRROR contains a comparatively large number of dialogues with high turn density and dynamic visual responses. Unlike prior datasets such as M2CoSC (Kim et al., 2025), which uses a single static image per dialogue, or MEDIC (Zhu et al., 2023), which is limited to a single turn, MIRROR provides image sequences that evolve turn-by-turn in alignment with client emotion.

<sup>&</sup>lt;sup>4</sup>Version gpt-40-2024-08-06.

<sup>&</sup>lt;sup>5</sup>WAI stands for Working Alliance Inventory.

<sup>&</sup>lt;sup>6</sup>https://huggingface.co/Falconsai/nsfw\_image\_ detection

345

347

348

350

351

352

354

309

310

311



Figure 3: Overview of emotional captioning. The AI therapist infers the client's emotional state from facial cues and uses it to generate an empathetic, aligned response.

# 4 Emotional Captioning

275

279

284

286

287

290

291

294

296

297

299

301

303

305

307

308

To handle client resistance more effectively, we introduce *emotional captioning*, which is a reasoning module that interprets the client's emotional state from facial expressions. At each dialogue turn, the model receives a facial image and generates a short textual description of the client's emotional state (e.g., *looking down, slightly defensive*). This emotional caption is then used to guide the AI therapist's response, enabling more context-sensitive and empathetic interactions (Figure 3). By grounding the model's behavior in visual cues, *emotional captioning* supplements verbal input with nonverbal affective signals, improving alignment with the client's psychological state<sup>7</sup>.

## **5** Experimental Settings

Following Smith et al. (2022); Liu et al. (2023), and Lee et al. (2024a), we assess the AI therapist based on full simulated counseling sessions rather than turn-level assessments. Each session involves an AI therapist interacting with a virtual client exhibiting varying types of resistance. We compare different model variants to assess the contribution of *planning* and *emotional captioning* strategies.

# 5.1 Client Agents with Resistance

We adopt GPT-3.5-TURBO<sup>8</sup> as the virtual client and conduct simulations based on predefined multimodal profiles. The evaluation is carried out with 600 unique client profiles, which are not included in the MIRROR dataset. These profiles include three types of resistance, with 200 examples for each type. Additionally, to assess performance degradation compared to non-resistant clients, we include 200 examples of non-resistant clients. Each session is considered terminated if the client attempts to disengage after two consecutive turns.

Within each client's utterance, nonverbal cues are embedded as stage directions within brackets, as described in §3.3. Note that these cues are used for facial expression generation and are invisible to the AI therapist when generating responses. For facial expressions, we generate LLM-based client's images at each turn, following the same process used in dataset construction (§3.3 and §3.4). Appendix K provides client setup and simulation prompt details.

# 5.2 AI Therapist Model Variations and Baselines

**AI Therapist Baselines** Our primary baseline, MIRROR-LLAVA, is a LLAVA-v1.5-7B (Liu et al., 2024) trained on the MIRROR dataset. To examine the benefit of multimodal integration, we include CAMEL-LLAMA3<sup>9</sup>, a text-only CBT model trained on therapeutic dialogues (Lee et al., 2024a). We also evaluate general-purpose models that are not fine-tuned for counseling: LLAMA-3-8B, LLAVA-v1.5-7B, and GPT-3.5-TURBO. These serve as non-specialized baselines to assess the impact of domain adaptation and modality alignment. Further implementation details, including training procedures and reasoning prompts, are provided in Appendices C and D.

**Reasoning Variants** We incorporate two variants of chain-of-thought (CoT) reasoning: *planning* and *emotional captioning*. *Planning*, based on Lee et al. (2024a), is a pre-session reasoning process in which the AI therapist analyzes the client's background and objectives to formulate a structured counseling strategy. This strategy helps the model select suitable CBT techniques and maintain a facilitative role during the session, rather than directly challenging the client's thoughts. We denote models incorporating these reasoning strategies with the subscripts P (*planning*) and EC (*emotional captioning*).

# 5.3 Metrics for Assessment

As defined in §2, we evaluate the therapist's ability to manage client resistance across two key areas: therapist skills and client alliance.

**Therapist skills** are assessed using the COUN-SELINGEVAL framework (Lee et al., 2024a), which

<sup>&</sup>lt;sup>7</sup>Prompt templates used in this process are detailed in Appendix K.

<sup>&</sup>lt;sup>8</sup>Version gpt-3.5-turbo-0125.

<sup>&</sup>lt;sup>9</sup>https://huggingface.co/cactus-camel/ camel-llama3

Model	General Counseling Skills (↑)			CBT-specific	Response Length		
	Understanding	Interpersonal Effectiveness	Collaboration	Guided Discovery	Focus	Avg.	Max
LLAMA-3-8B	3.811*   -0.073	<b>4.114</b>   -0.012	2.734*   -0.311	3.689*   -0.096	3.692*   -0.057	59.36	104.59
CAMEL-LLAMA3	3.794*   -0.085	4.003*   -0.002	2.279*   -0.198	3.527*   -0.127	3.563*   -0.197	20.54	27.42
GPT-3.5-TURBO	3.798*   -0.172	4.049   -0.041	2.976*   -0.194	3.462*   -0.262	3.491*   -0.238	36.19	57.28
LLAVA-v1.5-7B	3.622*   -0.066	3.997*   +0.007	3.408*   +0.071	2.494*   +0.057	2.501*   -0.012	112.41	177.11
MIRROR-LLAVA	3.973*   -0.017	4.025   -0.040	3.576*   -0.089	3.875*   -0.025	3.888*   -0.012	27.68	32.14
MIRROR-LLAVA $_P$	3.985   -0.015	4.098   +0.063	3.722*   +0.117	3.915*   -0.040	3.915*   +0.015	27.00	32.02
$MIRROR-LLAVA_{P+EC}$	4.000   +0.010	4.055   +0.010	3.913   -0.082	3.977   +0.007	<b>3.977</b>   +0.037	27.55	34.20

Table 2: Therapist skills assessment scores calculated by GPT-40 and response length. Asterisk (\*) indicates a significant difference compared to MIRROR-LLAVA<sub>P+EC</sub> (p < 0.05, paired t-test). **Response Length** denotes the average and maximum number of tokens per turn. Values after the vertical bar (|) indicate performance changes when interacting with resistant clients, relative to non-resistant clients; negative values denote a decline.

covers both general counseling skills and CBTspecific competencies. In particular, general counseling skills encompass the ability to interpret client concerns (Understanding), maintain a therapeutic relationship (Interpersonal Effectiveness), and facilitate collaborative decision-making (Collaboration). Meanwhile, CBT-specific skills evaluate the ability to guide clients in discovering their thoughts (Guided Discovery) and identify mal-adaptive patterns (Focus). Each component of the therapist's skills is rated on a scale from 0 to  $6^{10}$ .

**Client alliance** is measured following Li et al. (2024a), which assesses agreement of therapy objectives (Goal), engagement in counseling tasks (Approach), and the strength of emotional connection (Affective Bond), and is scored from 1 to 5.

# 6 Results and Discussion

356

364

365

369

371

372

373

376

377

384

# 6.1 Therapist Skills Assessment

Table 2 reports the evaluation of therapist skills in interactions with resistant clients.

Text-based Versus Vision Language Models As can be seen, text-based LLMs generally struggled to engage with resistant clients, particularly in collaborative interactions that demand heightened sensitivity to client emotions. This can be seen in the significant drop in performance compared to nonresistant clients. In contrast, vision-enhanced models showed greater resilience, maintaining higher scores even when interacting with resistant clients. These results highlight the importance of nonverbal cues in effectively managing challenging client interactions.

**Fine-Tuning and CoT on CBT Performance** Compared to LLAVA-V1.5-7B, which is the backbone model of MIRROR-LLAVA, the MIRROR-LLAVA family models achieved significantly higher scores in CBT-specific skills. This demonstrates the effectiveness of the MIRROR dataset in enhancing CBT skills and reinforces the notion that, despite being trained on vast amounts of pre-existing data, LLMs still require targeted finetuning to effectively internalize and apply CBT principles. Further performance gains were observed when CoT processes, such as *planning* and *emotional captioning*, resulting in responses that were more contextually appropriate and emotionally attuned to the client's needs.

Analysis of Response Length Excessively long response generation has been a persistent issue for LLMs and is known to reduce user satisfaction (Huang et al., 2024). Our analysis of response length revealed that, with the exception of the fine-tuned CBT counseling models (i.e. CAMEL-LLAMA3 and MIRROR-LLAVA family models), most models generated responses exceeding 30 tokens, which can degrade the counseling effectiveness. To further investigate these results, we provide actual examples for each model in Appendices H.1 and H.2, and conduct an error analysis in Appendix G.

# 6.2 Client Alliance Assessment

Table 3 presents the client alliance assessment using GPT-40, which evaluates how well each model supports goal completion, establishes rapport (Approach), and fosters emotional connection (Affective Bond). 390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

<sup>&</sup>lt;sup>10</sup>For our experiments, we do not use a Strategy score, which assesses the coherence of intervention strategies, as it strongly correlates with the length of the AI therapist's responses (see Appendix E).

Madal	Client Alliance Skills (↑)				
WIOdel	Goal	Approach	Affective Bond		
LLAMA-3-8B	2.412*   -0.023	3.309*   -0.107	3.356*   -0.138		
CAMEL-LLAMA3	2.358*   -0.009	3.130*   -0.072	3.149*   -0.203		
GPT-3.5-TURBO	2.472*   -0.018	3.272*   -0.168	3.297*   -0.253		
LLAVA-v1.5-7B	2.589   -0.048	3.234*   -0.181	3.356*   -0.163		
MIRROR-LLAVA	2.459*   -0.033	3.289*   -0.060	3.400*   -0.092		
MIRROR-LLAVA $_P$	2.525*   +0.033	3.340   -0.005	3.448   -0.051		
$MIRROR-LLAVA_{P+EC}$	2.567   +0.035	3.366   -0.003	3.480   -0.024		

Table 3: Client alliance assessment results as evaluated by GPT-40.

While overall alliance scores improve with MIR-ROR, we observe a modest decline in the "Goal" score for MIRROR-LLAVA models compared to some baselines. We attribute this to the design of the MIRROR dataset, which emphasizes emotional engagement and rapport-building in resistant counseling scenarios, rather than directive goal setting. In real-world counseling, especially under resistance, it is often more effective to prioritize emotional engagement before directive goal-setting. This trade-off is reflected in the substantial gains in Approach and Affective Bond scores, which more directly capture the model's capacity for empathy and responsiveness. Notably, MIRROR-LLAVA $_{P+EC}$  achieves the highest scores in these affective dimensions, demonstrating the strength of step-by-step reasoning in managing resistance.

# 6.3 Domain Expert Assessment

To further validate previous client alliance results, we conducted pairwise comparisons between MIRROR-LLAVA<sub>P+EC</sub>, LLAMA-3-8B, and CAMEL-LLAMA3 using 200 randomly selected cases from the test set, balanced across three resistance categories: emotional, cognitive, behavioral. Specifically, two domain experts evaluated the models and selected the better model in each comparison (Appendix F).

Moreover, we focused on comparing our method against the strongest baselines in CBT counseling—LLAMA-3-8B and CAMEL-LLAMA3, which are ranked highest in CBT-specific skill. Figure 4 shows the average win rate across all pairwise comparisons. As depicted, the win rate confirmed that MIRROR-LLAVA<sub>P+EC</sub> consistently outperformed its counterparts across all three dimensions of the therapeutic alliance. This result confirms that our model is not only favored in automatic evaluations but also by actual counseling experts. In particular, while GPT-based evaluation



Figure 4: Pairwise comparison results among MIRROR-LLAVA, CAMEL-LLAMA3 and LLAMA-3-8B, on three evaluation criteria—Goal, Approach, and Affective Bond—rated by two psychotherapists.

showed limited gains in the "Goal" dimension, domain experts more frequently selected MIRROR-LLAVA<sub>P+EC</sub> as superior in goal-related dialogue segments. Experts noted that goal pursuit was achieved more implicitly through sustained rapport and motivational alignment, rather than through direct or premature intervention. This reinforces our claim that emotional connection should precede goal-setting in resistant counseling contexts. 461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

It is interesting to note that GPT-40's evaluation in Table 3 ranked LLAMA-3-8B higher than other LLMS, whereas domain experts preferred CAMEL-LLAMA3 significantly more in pairwise comparisons. This discrepancy is likely due to human preference for responses of a more natural length, rather than those that are excessively long.

# 6.4 Application in Real-World Counseling Demonstrations

In this section, we conduct analyses on actual client interactions. Due to the lack of CBT-based counseling videos, we instead analyzed motivational interviewing (MI) demonstration videos to assess how our model responds to various forms of resistance in practice. For this, we specifically used the AnnoMI dataset (Wu et al., 2022), which contains 133 counseling session videos along with transcriptions of the clients' utterances. By examining both the video frames and client statements, we identified specific client utterances that show resistance and explored how MIRROR<sub>P+EC</sub> responds to these resistant remarks.

As shown in Figure 5, clients exhibit various forms of resistance, including reluctance to seek help (Client A), claiming impunity (Client B), minimizing concerns (Client C), and externalizing blame (Client D)<sup>11</sup>. The case analysis demon-

456

457

458

459

460

422

423

<sup>&</sup>lt;sup>11</sup>These are well-documented resistance patterns in psy-



That makes sense. It's a big step to seek help, and it's normal to feel uncertain about it. What thoughts come to mind when you think about being here?

Ours Emotional Caption: Sighs

(a) Client A case: Reluctance to seek help



(d) Client D case: Externalizing blame

Figure 5: Four examples of MIRROR-LLAVA $_{P+EC}$ responses in psychological counseling, showcasing its ability to handle resistance through validation and openended questioning.

strates that  $MIRROR_{P+EC}$  effectively identifies the client's emotional state through captioning and responds with emotional validation<sup>12</sup> and open-ended questions, common therapeutic techniques for managing resistance (Miller and Rollnick, 2002). For

497

498

499

500

501

example, when Client D externalizes blame onto family, the therapist acknowledges their feelings of isolation while gently redirecting the conversation toward exploring ways to cope with pressure. Further case studies can be found in Appendix H.3.

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

#### 7 Conclusion

In this paper, we explore the use of multimodal cognitive reframing therapy for managing client resistance. Given the challenges faced by LLMs in addressing resistant clients and the potential advantages of VLMs, we aim to enhance AI therapists' ability to manage resistance by incorporating nonverbal cues, particularly facial expressions, to detect and understand client resistance. To address this challenge while mitigating privacy concerns associated with real-person data, we developed MIR-ROR, a novel synthetic dataset for multimodal cognitive reframing therapy. Additionally, we have evaluated the AI therapist's performance in two key areas: therapeutic skills and alliance-building, as well as adaptability to real-world counseling scenarios. Our results demonstrate significant improvements in both areas when trained with MIRROR, underscoring its potential for real-world therapeutic applications. These improvements contribute to the development of AI therapists that are more empathetic and capable of fostering stronger therapeutic relationships.

# Limitations

**Biases in Image Generation** We utilized LLMs to generate image prompts, which were then used to create images via PhotoMaker. During the process of generating image prompts based on the client's stage directions and utterances, cultural biases inherent in the LLMs (AlKhamissi et al., 2024; Naous et al., 2024) may have influenced the output. For example, if a client says, "[Smiling] Hi," the LLMs might generate a prompt like "eyes curved like a crescent moon," reflecting an East Asian interpretation of a smile, whereas other cultures might describe a smile based on visible teeth or dimples (Srinivasan and Martinez, 2021). Similarly, if a client expresses sadness, the prompt may depict downcast eyes and a bowed head, aligning with Western portrayals of sadness, while some cultures associate sadness with a neutral or stoic expression. These examples show how cultural biases in LLMs can shape image generation, reinforcing specific interpretations while overlooking others.

chotherapy (Miller and Rollnick, 2002).

<sup>&</sup>lt;sup>12</sup>This term refers to accepting a client's emotions without judgment, helping to sense of being understood while encouraging to manage their emotions.

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

601

602

Authenticity in AI-Generated Images Since the 551 final facial images were synthesized using an AI-552 based image generation model, they may not fully 553 replicate the characteristics of real-world expressions. Subtle nuances, such as micro-expressions, 555 muscle tension, and natural asymmetries, might be 556 lost or inaccurately rendered, potentially affecting 557 their authenticity. These limitations could influence how AI therapists perceive and interpret emotional cues, as synthesized images may lack the depth and 560 variability of genuine human expressions. 561

Scope and Session Length In contrast to typi-562 cal counseling sessions, which last about an hour 563 and extend over multiple sessions, our dataset con-564 sists of relatively short, single-session dialogues. This limitation makes it challenging to facilitate deeper exploration of cognitive distortions or sustained reframing over multiple sessions. Future research should focus on extending dialogues to 569 longer, multi-session interactions and incorporating 570 a wider range of counseling techniques to enable 571 AI to provide more continuous and in-depth therapeutic support. 573

**Conversational Structure and Termination** 574 Our framework does not impose strict turn-level 575 constraints or predefined termination points within 576 the dialogue. While we incorporate a counseling strategy, **planning**, to maintain goal orientation, the 578 absence of explicit session boundaries may result in prolonged interactions without meaningful therapeutic progress. For example, if a simulated client remains in a negative emotional state, the AI therapist may continue offering supportive statements 583 rather than facilitating cognitive change. This limitation highlights the importance of incorporating clearer session structures or exit strategies in future designs to better align with therapeutic goals. 587

**Evaluating MIRROR Using a GPT-Based Client** 588 In § 5, we have conducted counseling sessions us-589 ing a GPT-based client with a model trained on 590 the MIRROR dataset. Through these conversation records, we were able to demonstrate the effec-592 tiveness of the MIRROR dataset, and expert evaluations were also conducted on these dialogues. While GPT effectively simulates real client behav-596 ior and demonstrates strong conversational abilities, validating our approach through real-time human interactions would have further strengthened our findings. However, due to cost constraints and privacy concerns, real-time human interactions were

not feasible. Instead, we evaluated our approach in a real therapy scenario, as described in § 6.4.

**Diversity of Resistance** Resistance can be expressed not only through facial expressions but also through overall body posture, voice tone, speech timing, and other nuanced factors. While these elements collectively shape how resistance manifests, this study focuses specifically on facial expressions and utterances, as they provide the most direct and observable indicators. Future research could enhance the dataset by incorporating audio, body posture, and other multimodal elements, leading to a more comprehensive understanding of resistance in therapeutic interactions.

**Model Selection and Generalization** Although we trained the LLAVA-V1.5-7B model with two different CoT options and demonstrated its strong performance in handling client resistance and CBT counseling, our evaluation was based on a single backbone model. This could be a limitation, as there may be other VLMs that could perform better or differently, depending on their architecture or training. The reliance on a single model limits the generalizability of our findings. To address this, we plan to explore and evaluate additional VLMs in future work, comparing their performance across various therapeutic tasks and client engagement scenarios.

# **Ethical Statement**

**Privacy Considerations for Images** Ensuring privacy and ethical integrity is a fundamental priority in our dataset construction. We utilize the CelebA dataset (Liu et al., 2015), which is distributed under the MMLAB license. This license strictly prohibits commercial use and redistribution of the dataset. In compliance with these terms and to respect the rights of the individuals depicted in the images, we do not share the raw images directly. Instead, we provide image links and code that enables researchers to process the dataset independently, ensuring that the dataset's usage remains within ethical and legal boundaries.

**Privacy Considerations for Dialogue** The dialogue seeds for this dataset were sourced from the CACTUS dataset (Lee et al., 2024a), with PATTERN-REFRAME (Maddela et al., 2023b) serving as its seed dataset. This dataset does not contain actual medical records but was collected through crowdsourcing, where each participant was assigned a

753

754

650persona and instructed to role-play. Additionally,651during the dataset generation process, no utterances652were derived from real individuals' personas; in-653stead, all dialogues were fully synthesized. This654approach further mitigates privacy concerns by en-655suring that no personal data is incorporated into the656dataset.

Safety Considerations While AI has the potential to provide support, it may also have unintended negative effects on individuals with mental health challenges (Luxton, 2014). Although our model has demonstrated some degree of effectiveness, our primary objective was to explore whether AI can effectively engage with patients who exhibit resistance to therapy. Therefore, we believe that AI should be used under the supervision of a professional rather than serving as a standalone tool in counseling sessions, particularly for individuals with severe psychological conditions beyond its intended scope. Additionally, to ensure the safety and appropriateness of the dataset, we implemented NSFW filtering and incorporated Canary to identify and remove conversations that may require human intervention. 673

Bias Considerations Although we utilize ran-674 domly selected images and a dialogue seed dataset 675 that incorporates diversity in age, gender, and occupation, there remains a possibility of bias in our dataset. This is primarily due to our reliance 678 on LLMs, which are predominantly trained on 679 Western-centric datasets. In particular, during the screenplay generation process, gestures and nonverbal cues may vary across cultures. Since these were generated using GPT-40-MINI, certain gestures may not align with cultural norms in specific regions. Therefore, to ensure cultural appropriateness, retraining and adaptation would be necessary before deploying the model in a specific country.

# References

688

695

- 689 AI@Meta. 2024. Llama 3 model card.
  - Badr AlKhamissi, Muhammad ElNokrashy, Mai Alkhamissi, and Mona Diab. 2024. Investigating cultural alignment of large language models. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 12404–12422, Bangkok, Thailand. Association for Computational Linguistics.
    - Loyd Beal III, Jacqueline M Stavros, and Matthew L Cole. 2013. Effect of psychological capital and

resistance to change on organisational citizenship behaviour. *SA Journal of Industrial Psychology*, 39(2):1–11.

- Michael J Constantino, Henny A Westra, Martin M Antony, and Alice E Coyne. 2017. Specific and common processes as mediators of the long-term effects of cognitive-behavioral therapy integrated with motivational interviewing for generalized anxiety disorder. *Psychother Res*, 29(2):213–225.
- Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. 2019. Arcface: Additive angular margin loss for deep face recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4690–4699.
- Kathleen Kara Fitzpatrick, Alison Darcy, and Molly Vierhile. 2017. Delivering cognitive behavior therapy to young adults with symptoms of depression and anxiety using a fully automated conversational agent (woebot): a randomized controlled trial. *JMIR mental health*, 4(2):e7785.
- Anmol Goel, Nico Daheim, and Iryna Gurevych. 2024. Socratic reasoning improves positive text rewriting. *arXiv preprint arXiv:2403.03029*.
- MD Romael Haque and Sabirat Rubya. 2023. An overview of chatbot-based mobile mental health apps: insights from app description and user reviews. *JMIR mHealth and uHealth*, 11(1):e44838.
- Kimberley Michelle Hara. 2020. Comparing therapist responsivity to resistance markers in cognitivebehavioural therapy and motivational interviewing integrated with cognitive-behavioural therapy for generalized anxiety.
- Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*.
- Shih-Hong Huang, Ya-Fang Lin, Zeyu He, Chieh-Yang Huang, and Ting-Hao Kenneth Huang. 2024. How does conversation length impact user's satisfaction? a case study of length-controlled conversations with llm-powered chatbots. In *Extended Abstracts of the CHI Conference on Human Factors in Computing Systems*, pages 1–13.
- Hyunwoo Kim, Jack Hessel, Liwei Jiang, Peter West, Ximing Lu, Youngjae Yu, Pei Zhou, Ronan Bras, Malihe Alikhani, Gunhee Kim, Maarten Sap, and Yejin Choi. 2023. SODA: Million-scale dialogue distillation with social commonsense contextualization. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12930–12949, Singapore. Association for Computational Linguistics.
- Hyunwoo Kim, Youngjae Yu, Liwei Jiang, Ximing Lu, Daniel Khashabi, Gunhee Kim, Yejin Choi, and Maarten Sap. 2022. ProsocialDialog: A prosocial

755

756

- 765 766 767 768 769 770 771 772 773 774 775 776 777 778 779 780 780 781 782
- 780 781 782 783 784 785 786 785 786 787 788 789 790 791 792
- 791 792 793 794 795
- 7777
- 7
- 8
- 803
- 805 806 807
- 8
- 80
- 810 811

backbone for conversational agents. In *Proceedings* of the 2022 Conference on Empirical Methods in Natural Language Processing, pages 4005–4029, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

- Subin Kim, Hoonrae Kim, Heejin Do, and Gary Lee. 2025. Multimodal cognitive reframing therapy via multi-hop psychotherapeutic reasoning. In Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 4863–4880, Albuquerque, New Mexico. Association for Computational Linguistics.
- Suyeon Lee, Sunghwan Kim, Minju Kim, Dongjin Kang, Dongil Yang, Harim Kim, Minseok Kang, Dayi Jung, Min Hee Kim, Seungbeen Lee, Kyong-Mee Chung, Youngjae Yu, Dongha Lee, and Jinyoung Yeo. 2024a. Cactus: Towards psychological counseling conversations using cognitive behavioral theory. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 14245–14274, Miami, Florida, USA. Association for Computational Linguistics.
- Young-Jun Lee, Dokyong Lee, Junyoung Youn, Kyeong-Jin Oh, Byungsoo Ko, Jonghwan Hyeon, and Ho-Jin Choi. 2024b. Stark: Social long-term multi-modal conversation with persona commonsense knowledge. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 12137–12162, Miami, Florida, USA. Association for Computational Linguistics.
- Anqi Li, Yu Lu, Nirui Song, Shuai Zhang, Lizhi Ma, and Zhenzhong Lan. 2024a. Understanding the therapeutic relationship between counselors and clients in online text-based counseling using LLMs. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 1280–1303, Miami, Florida, USA. Association for Computational Linguistics.
- Zhen Li, Mingdeng Cao, Xintao Wang, Zhongang Qi, Ming-Ming Cheng, and Ying Shan. 2024b. Photomaker: Customizing realistic human photos via stacked id embedding. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2024. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26296–26306.
- June M Liu, Donghao Li, He Cao, Tianhe Ren, Zeyi Liao, and Jiamin Wu. 2023. Chatcounselor: A large language models for mental health support. *arXiv preprint arXiv:2309.15461*.
- Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. 2015. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*.

David D Luxton. 2014. Recommendations for the ethical use and design of artificial intelligent care providers. *Artificial intelligence in medicine*, 62(1):1–10.

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

- Mounica Maddela, Megan Ung, Jing Xu, Andrea Madotto, Heather Foran, and Y-Lan Boureau. 2023a. Training models to generate, recognize, and reframe unhelpful thoughts. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 13641– 13660, Toronto, Canada. Association for Computational Linguistics.
- Mounica Maddela, Megan Ung, Jing Xu, Andrea Madotto, Heather Foran, and Y-Lan Boureau. 2023b. Training models to generate, recognize, and reframe unhelpful thoughts. *arXiv preprint arXiv:2307.02768*.
- Ashish Mehta, Andrea Nicole Niles, Jose Hamilton Vargas, Thiago Marafon, Diego Dotta Couto, and James Jonathan Gross. 2021. Acceptability and effectiveness of artificial intelligence therapy for anxiety and depression (youper): longitudinal observational study. *Journal of medical Internet research*, 23(6):e26771.
- W R Miller and S Rollnick. 2002. *Motivational interviewing: Preparing people for change*. The Guilford Press.
- Theresa B Moyers and Tim Martin. 2006. Therapist influence on client language during motivational interviewing sessions. *J Subst Abuse Treat*, 30(3):245–251.
- Hongbin Na. 2024. CBT-LLM: A Chinese large language model for cognitive behavioral therapy-based mental health question answering. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 2930–2940, Torino, Italia. ELRA and ICCL.
- Tarek Naous, Michael J Ryan, Alan Ritter, and Wei Xu. 2024. Having beer after prayer? measuring cultural bias in large language models. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 16366–16393, Bangkok, Thailand. Association for Computational Linguistics.
- Gerald R Patterson and Patricia Chamberlain. 1994. A functional analysis of resistance during parent training therapy. *Clinical Psychology: Science and Practice*, 1(1):53.
- Renyi Qu, Lyle Ungar, and João Sedoc. 2023. Conditioning on dialog acts improves empathy style transfer. In *The 2023 Conference on Empirical Methods in Natural Language Processing*.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark,

- 869 870 871 872 873 874 875 876 877 878 879 880 881
- 8 8 8 8 8 8 8 8 8 8
- 88 88
- 88
- 8
- 894 895 896
- 897 898
- 899 900
- 901 902
- 903 904 905 906
- 907 908 909
- 910 911
- 912 913
- 914
- 915 916
- 917 918 919

- 921 922 923
- 923 924
- 924 925

Gretchen Krueger, and Ilya Sutskever. 2021. Learning transferable visual models from natural language supervision. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 8748–8763. PMLR.

- Sefik Ilkin Serengil and Alper Ozpinar. 2021. Hyperextended lightface: A facial attribute analysis framework. In 2021 International Conference on Engineering and Emerging Technologies (ICEET), pages 1–4. IEEE.
- Roz Shafran, Pamela Myles-Hooton, Sophie Bennett, and Lars-Göran Öst. 2021. The concept and definition of low intensity cognitive behaviour therapy. *Behav Res Ther*, 138:103803.

Ashish Sharma, Kevin Rushton, Inna Lin, David Wadden, Khendra Lucas, Adam Miner, Theresa Nguyen, and Tim Althoff. 2023. Cognitive reframing of negative thoughts through human-language model interaction. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9977–10000, Toronto, Canada. Association for Computational Linguistics.

- Eric Smith, Orion Hsu, Rebecca Qian, Stephen Roller, Y-Lan Boureau, and Jason Weston. 2022. Human evaluation of conversations is an open problem: comparing the sensitivity of various methods for evaluating dialogue agents. In *Proceedings of the 4th Workshop on NLP for Conversational AI*, pages 77– 97, Dublin, Ireland. Association for Computational Linguistics.
- Ramprakash Srinivasan and Aleix M. Martinez. 2021. Cross-cultural and cultural-specific production and perception of facial expressions of emotion in the wild. *IEEE Transactions on Affective Computing*, 12(3):707–721.
- Yinzhou Wang, Yimeng Wang, Ye Xiao, Liabette Escamilla, Bianca Augustine, Kelly Crace, Gang Zhou, and Yixuan Zhang. 2025. Evaluating an Ilm-powered chatbot for cognitive restructuring: Insights from mental health professionals. *arXiv preprint arXiv:2501.15599.*
- Henny A Westra and Nikoo Norouzian. 2018. Using motivational interviewing to manage process markers of ambivalence and resistance in cognitive behavioral therapy. *Cognitive Therapy and Research*, 42:193– 203.
- Chris Williams. 2001. Use of written cognitive-behavioural therapy self-help materials to treat depression. *Advances in Psychiatric Treatment*, 7(3):233–240.
- Zixiu Wu, Simone Balloccu, Vivek Kumar, Rim Helaoui, Ehud Reiter, Diego Reforgiato Recupero, and Daniele Riboni. 2022. Anno-mi: A dataset of expert-annotated counselling dialogues. In *ICASSP* 2022 - 2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 6177–6181.

Mengxi Xiao, Qianqian Xie, Ziyan Kuang, Zhicheng Liu, Kailai Yang, Min Peng, Weiguang Han, and Jimin Huang. 2024. HealMe: Harnessing cognitive reframing in large language models for psychotherapy. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1707–1725, Bangkok, Thailand. Association for Computational Linguistics. 926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

- Kailai Yang, Shaoxiong Ji, Tianlin Zhang, Qianqian Xie, Ziyan Kuang, and Sophia Ananiadou. 2023. Towards interpretable mental health analysis with large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 6056–6077.
- Kailai Yang, Tianlin Zhang, Ziyan Kuang, Qianqian Xie, Jimin Huang, and Sophia Ananiadou. 2024. Mentallama: Interpretable mental health analysis on social media with large language models. In *Proceedings* of the ACM on Web Conference 2024, WWW '24, page 4489–4500, New York, NY, USA. Association for Computing Machinery.
- Chenhao Zhang, Renhao Li, Minghuan Tan, Min Yang, Jingwei Zhu, Di Yang, Jiahao Zhao, Guancheng Ye, Chengming Li, and Xiping Hu. 2024. CPsyCoun: A report-based multi-turn dialogue reconstruction and evaluation framework for Chinese psychological counseling. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 13947– 13966, Bangkok, Thailand. Association for Computational Linguistics.
- Zhouan Zhu, Chenguang Li, Jicai Pan, Xin Li, Yufei Xiao, Yanan Chang, Feiyi Zheng, and Shangfei Wang. 2023. Medic: A multimodal empathy dataset in counseling. In *Proceedings of the 31st ACM International Conference on Multimedia*, MM '23, page 6054–6062, New York, NY, USA. Association for Computing Machinery.
- Caleb Ziems, Minzhi Li, Anthony Zhang, and Diyi Yang. 2022. Inducing positive perspectives with text reframing. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics* (*Volume 1: Long Papers*), pages 3682–3700, Dublin, Ireland. Association for Computational Linguistics.

# A License

969

970

971

972

973

975

977

983

990

991

1000

1001

1002

1003

1004

1006

1007

1008

1011

1012

1013

1015

MIRROR is constructed using the CelebA (Liu et al., 2015) and the CACTUS datasets (Lee et al., 2024a). CelebA is released under the MMLAB license, which restricts redistribution, while CACTUS is licensed under the GPL-2.0 license, permitting noncommercial scientific use. In adherence to these licensing terms, we do not directly include images from these datasets in MIRROR. Instead, we provide links to the original sources. Consequently, MIRROR is distributed under the GPL-2.0 license, ensuring compliance with the licensing conditions of the datasets used.

# **B** Related Work

Research on AI-assisted cognitive reframing therapy has largely focused on text-based approaches using LLMs. Early studies explored sentence rewriting to address cognitive distortions (Ziems et al., 2022; Maddela et al., 2023a; Sharma et al., 2023; Goel et al., 2024), leveraging evidence that low-intensity CBT interventions can be effective in self-help formats (Williams, 2001; Shafran et al., 2021). More recent advancements have shifted toward conversational approaches, moving from simple query-response interactions (Na, 2024; Liu et al., 2023) to structured, multi-turn frameworks. For instance, Xiao et al. (2024) introduced a multiturn framework with a structured three-stage process to ensure that AI serves as a facilitator rather than a direct corrector. Other studies have focused on improving the realism of cognitive reframing datasets (Lee et al., 2024a) and enhancing AI therapists' professional counseling competence (Zhang et al., 2024).

Most recently, there has been a growing interest in incorporating nonverbal cues into AI-driven psychotherapy. Expanding on Xiao et al. (2024), Kim et al. (2025) explored multimodal cognitive reframing, showing that VLMs can generate more empathic responses. Similarly, Zhu et al. (2023) emphasized clients' nonverbal expressions, introducing a multimodal empathy dataset. Our work advances multimodal cognitive reframing by focusing on managing client resistance, strengthening the therapeutic alliance, and improving AI-assisted psychotherapy.

### C Training Details

1016The LLAVA-v1.5-7B model was fine-tuned on the1017MIRROR dataset using LoRA (Hu et al., 2022) for



Figure 6: The overview of the planning process.

5 epochs. We used the official LLAVA-v1.5-7B model from Hugging Face<sup>13</sup> and followed the default hyperparameters<sup>14</sup>, which include a learning rate of 2e-5, an AdamW optimizer without weight decay, and a cosine learning rate schedule with a 3% warmup ratio. Training was done on four A100-80GB GPUs with a batch size of 32 per GPU.

1018

1019

1020

1021

1022

1023

1024

1027

1028

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1044

# **D** Counseling Strategy Planning

The *planning* process follows the approach of Lee et al. (2024a), where, prior to the counseling session, the model first infers a counseling strategy based on the client's information and counseling objectives, which is then incorporated into the response generation (Figure 6). The client information includes details typically found in a counseling intake form, such as name, age, gender, and occupation.

The prompts used for *planning* in our MIRROR-LLAVA model are provided in Appendix K.

For baseline comparisons, we followed the official prompt structures for CAMEL-LLAMA3, GPT-3.5-TURBO, and LLAMA-3-8B. In particular, the planning component in CAMEL-LLAMA3 follows the official implementation by Lee et al. (2024a).

# E Impact of Response Length on GPT Evaluation

We examined the correlation between the AI therapist's response length and its performance in GPT-1045based evaluation. Across all models, we analyzed1047

<sup>&</sup>lt;sup>13</sup>https://huggingface.co/liuhaotian/llava-v1. 5-7b

<sup>&</sup>lt;sup>14</sup>https://github.com/haotian-liu/LLaVA/tree/ main







(b) VLM therapist models

Figure 7: Correlation between general counseling performance and the response length of AI therapists. All coefficients. All correlation coefficients were statistically significant (p < 0.05).

Figure 8: Correlation between CBT performance and the response length of AI therapists. All correlation coefficients were statistically significant (p < 0.05).

1048 how response length affects evaluation metrics and further aggregated the results by modality. As 1049 shown in Figure 7 and Figure 8, there is a no-1050 ticeable relationship between response length and 1051 performance in both general counseling skills and 1052 CBT techniques. Notably, the strongest correlation was observed in the Strategy category, with a corre-1054 lation of 0.6, suggesting that untrained text-based 1055 LLMs tend to receive higher evaluations from the 1056 GPT evaluator when generating longer responses. 1057 1058 This is likely because lengthier responses incorporate multiple questions or strategies within a sin-1059 gle reply, which the evaluator interprets as higher-1060 quality output. In contrast, for VLMs, response length showed no significant correlation with per-1062

formance in general counseling skills.

However, within CBT techniques, particularly in Focus and Guided Discovery, shorter responses generally resulted in higher scores. This trend is likely influenced by the LLAVA-v1.5-7B model, which tends to generate unnaturally long responses and has lower scores. Compared to the LLAVAv1.5-7B model, the MIRROR-LLAVA family produced shorter responses and achieved better scores, suggesting a correlation between shorter responses and higher performance.

# F Domain Expert Assessment Details

# F.1 Numerical Details

Table 4, 5, and 6 show the winning rates for eachmetric: Goal, Approach, and Affective Bond.

1075

1077

1063

1064

1065

# F.2 Domain Expert Recruitment

1078

1080

1081

1082

1083

1086

1087

1088

1089

1090

1091

1092

1094

1095

1096

1097

1100

1101

For the domain expert evaluation, we hired two evaluators through the Upwork platform<sup>15</sup> who hold a counseling license or have a graduate degree in a related field. They were informed that all personal information would remain anonymous and that responses would be used solely for research purposes. We paid \$0.05 per data entry for pairwise comparison, which they accepted before proceeding with the task.

# G Error Analysis

To gain deeper insights into the effectiveness and limitations of our proposed method, we conducted an error analysis on the MIRROR-LLAVA<sub>P+EC</sub> model, and focused on cases with therapist skill and client alliance scores of less than 3.

# G.1 Failure Cases in General Counseling Skills



(a) Failure case of low Understanding and Interpersonal Effectiveness.



(b) Failure case of low Collaboration.

Figure 9: Failure cases in general counseling skills.

Figure 9 illustrates failure cases in general counseling skills. The first case (Figure 9a) is due to a hallucination from the VLM. Although the client did not mention that a colleague was a pedophile, the VLM therapist incorrectly introduced this idea, which made the client uncomfortable. This misstep resulted in low understanding and interpersonal skills.

1102

1103

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

The second case (Figure 9b) involves a client 1104 who expressed deep-seated fear and emotional re-1105 luctance, stating, "I'm always worried about say-1106 ing or doing the wrong thing." Rather than fur-1107 ther exploring the client's underlying concerns, the 1108 model prematurely attempted to fix the problem 1109 before building intimacy. This response failed to 1110 align with the client's emotional state, leading to 1111 disengagement, which highlights the need for more 1112 specific *planning* in CBT counseling. 1113

### G.2 Failure Cases in CBT-Specific Skills



(a) Failure case where the therapist model is confused about its role.



ilure case where the model lacks a challenging

(b) Failure case where the model lacks a challenging exploration of distorted thoughts.

Figure 10: Failure cases in CBT-specific skills.

Figure 10 illustrates failure cases in CBT-specific skills. In the first case, shown in Figure 10a, confusion between the therapist's and the client's roles occurred. In this case, the therapist's utterance shifted to client's utterance in one turn. Although this happened in only five cases, it resulted in a drop in the focus score.

The second case (Figure 10b) arises when the therapist loses their purpose and simply sympathizes with the client's cognitive distortions. Instead of actively challenging the client's distorted thought patterns, the model engaged in emotionfocused inquiry, asking about specific experiences related to the client's feelings. While this approach may encourage emotional processing, it falls short

<sup>&</sup>lt;sup>15</sup>www.upwork.com

	LLAMA-3-8B	CAMEL-LLAMA3	$\mathbf{MIRROR}\text{-}\mathbf{LLAVA}_{P+EC}$	Win Rate (%)
LLAMA-3-8B	-	42.13	36.34	38.09
CAMEL-LLAMA3	57.87	-	43.43	50.65
MIRROR-LLAVA $_{P+EC}$	65.95	56.57	-	61.26

Table 4: Numerical results of pairwise comparison of three models, evaluated for Goal alignment score by two domain experts.

	LLAMA-3-8B	CAMEL-LLAMA3	$\mathbf{MIRROR}\text{-}\mathbf{LLAVA}_{P+EC}$	Win Rate (%)
LLAMA-3-8B	-	44.72	34.77	39.75
CAMEL-LLAMA3	55.28	-	43.97	49.62
MIRROR-LLAVA $_{P+EC}$	65.23	56.03	-	60.63

Table 5: Numerical results of pairwise comparison of three models, evaluated for Approach score by two domain experts.

	LLAMA-3-8B	CAMEL-LLAMA3	$\mathbf{MIRROR}\text{-}\mathbf{LLAVA}_{P+EC}$	Win Rate (%)
LLAMA-3-8B	-	40.19	49.35	44.77
CAMEL-LLAMA3	50.65	-	42.46	46.55
$\mathbf{MIRROR}\text{-}\mathbf{LLAVA}_{P+EC}$	59.81	57.54	-	58.67

Table 6: Numerical results of pairwise comparison of three models, evaluated for Affective Bond score by two domain experts.

in fostering cognitive reframing. A more effective
intervention would involve helping the client examine the reasons behind their beliefs and exploring
alternative perspectives.

positive reframing is valuable, doing so too soon leads to a mismatch in attunement, making the client feel unheard or dismissed.

1144

1145

1146

1147

1151

# G.3 Failure Cases in Client Alliance

1134

1135

1136

1137

1138

1139

1140

1141

1142

1143



Figure 11: A representative failure case in the therapeutic alliance.

In analyzing cases where the client alliance score was below 3, we identified a key issue that hinders effective therapeutic engagement (Figure 11). In this example, the client expresses a strong sense of reluctance and feeling stuck in a negative mindset, signaling deep-seated emotional resistance. However, rather than exploring these emotions further, the therapist prematurely shifts the focus toward finding positive experiences. While encouraging

# H Case Study

We conducted a comparative analysis of AI thera-<br/>pist counseling sessions to examine how different1148models respond to and handle client resistance.1150

### H.1 Response Length Analysis Across Models

We compare the response lengths of different mod-1152 els when a client expresses the distorted thought, 1153 "bad things will happen". Figure 12 presents the 1154 actual responses from five AI therapist models to a 1155 virtual client's statement exhibiting cognitive dis-1156 tortion. The models that were not fine-tuned with 1157 CBT datasets, including LLAMA-3-8B, GPT-3.5-1158 TURBO, and LLAVA-v1.5-7B, tended to gener-1159 ate excessively long responses, which negatively 1160 impacted the effectiveness and naturalness of the 1161 dialogue. In contrast, CAMEL-LLAMA3 and 1162 MIRROR-LLAVA $_{P+EC}$ , which were fine-tuned for 1163 CBT counseling, produced responses of more ap-1164 propriate lengths, showing better alignment with 1165 client needs and making the interactions easier for 1166 clients to understand. 1167

#### Client

[Feeling hesitant] It's a combination of both, I think. There are specific moments that trigger these feelings, but overall, it's like this constant undercurrent of worry and fear that never really goes away.

LLAMA-3-8B-INSTRUCT

Samantha, it sounds like you're carrying a heavy emotional load, and it's affecting your daily life. I want to acknowledge that it takes a lot of courage to admit that you're struggling, and I'm here to support you. It sounds like there are specific triggers that set off these feelings of worry and fear, but there's also this underlying sense of anxiety that's always present. Can you tell me more about what these triggers are? Are they related to your singing, your personal relationships, or something else entirely? Also, when you say hat this feeling is an undercurrent that never really goes away, what does that mean to you? Is like a constant hum in the background, always present but not always overwhelming? Or is it more like a wave that crashes over you at times, making it hard to breathe? I want to help you explore these feelings and get a better understanding of what's driving them. But first, I want to acknowledge that it's okay to feel overwhelmed and scared. It's okay to not have all the answers right now. We can work through this together, one step at a time. How does that sound? Is there anything in particular that you'd like to focus on during our sessions together?

(a) LLAMA-3-8B

[Sighs] It's this fear I have... about playing tennis. I love it, but there's always this nagging thought that something bad will happen.

CAMEL-LLAMA3 I see. When did you first start feeling this way about playing tennis?

#### (b) CAMEL-LLAMA3

Client

Client

[Nervously fidgeting with hands] It's just... I can't help but imagine the worstcase scenario playing out in my head. The idea of losing something that brings me so much joy is terrifying. I feel like I wouldn't know how to cope if I couldn't sing or play tennis anymore. It's like a part of me would be gone.

> GPT-3.5-TURBO It's understandable that the thought of losing something you love, like singing and playing tennis, is causing you so much anxiety. It seems like these activities are a significant part of your identity and bring you a lot of fulfillment. Let's try to reframe this fear together. Instead of focusing on the fear of breaking your leg, what are some other ways you could approach this situation that might help ease your anxiety?



Client [Sighing softly] I know it's important to challenge my negative thoughts, but it feels like they have such a strong hold on me. It's like a constant battle in my mind.

LL2ACP-U1.52-78 As a psychotherapist conducting a CBT (Cognitius Behavioral Therapy) session. Louid like to acknowledge that it can be challenging to overcome fears and change negative thought patterns. However, with consistent practice and effort, you can improve your ability to manage anxiety and fear. To begin, let's identify any negative thoughts that contribute to your fear of public speaking. We can then challenge these thoughts and replace them with more positive and realistic alternatives. It's important to practice regularity, so us can set achievable goals and celebrate your progress along the way. Remember, you are not along in this journeg, and I'l be here to support you every step of the way. In our session, Lwould like to encourage you to take a deep breath, focus on the present moment, and observe any practice regularity to build your confidence and overcome your k together to challenge them and replace them with more positive alternatives. Remember, it's essential to practice rest and challing negative thought patterns takes time and effort, but with consistent practice and support, you can improve your ability to manage anxiety and fear. Let's continue working together to help you build your confidence and overcome your fear of public speaking.

(d) LLAVA-v1.5-7B



1168

1169

Client

[Feeling uneasy] It's just this constant fear at the back of my mind, you know? The thought of breaking my leg and not being able to do what I love anymore... It's terrifying. MIRROR-LLAVA<sub>P+EC</sub>. That sounds really overwhelming. It's understandable to feel anxious about such a powerful fear. How does it manifest for you in daily life?

(e) MIRROR-LLAVA $_{P+EC}$ .

Figure 12: Examples of AI therapist responses to virtual client resistance, with red highlights indicating instances of client resistance.

# H.2 Comparison of Resistance Management in Virtual Counseling

1170We conducted a comparative analysis of1171two counseling-optimized models, MIRROR-1172LLAVA $_{P+EC}$  and CAMEL-LLAMA3, to evaluate

their approaches to handling client resistance in virtual counseling sessions. Both models share the same goal of CBT and incorporate *planning* during inference. However, CAMEL-LLAMA3 is trained on a general CBT dataset and does not specifically address client resistance. In the CAMEL-LLAMA3 session (Figure 13a), when the client exhibited resistance, the model primarily relied on emotional validation and exploration. While these techniques offer comfort to client, they do not challenge deeper, malformed beliefs. As a result, while this model focused on emotional exploration with surface-level validation of the client's negative emotions, its reframing process lacked progression. 1173

1174

1175

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

1200

1201

1202

1203

1204

1205

1206

1207

1208

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

1219

1220

1221

1223

In contrast, MIRROR-LLAVA<sub>P+EC</sub> (Figure 13b) employed a more nuanced approach, integrating emotional validation, positive reinforcement, and cognitive reframing. Despite the client's resistance, our model attempted to delicately reframe the client's thoughts, using a collaborative approach with statements such as, "We can work on identifying those thoughts and reframing them into something more empowering together." Furthermore, by asking questions like "How does that sound?" the model encouraged client engagement and showed respect for the client's perspective. The model also effectively used positive reinforcement to encourage clients who were hesitant to take action by offering supportive statements like, "That's a brave and important step". These findings underscore the importance of specialized training datasets for client resistance, such as MIRROR, in effectively managing resistance and fostering therapeutic growth.

# H.3 Real-World Counseling Demonstrations

Here, we provide additional explanation for the cases in Figure 5. For Client A, the repeated use of '*I don't know*' illustrates reluctance to seek help, indicating emotional uncertainty and a lack of motivation to engage in the process. However, our model effectively addresses this resistance by validating the client's feelings and gently encouraging exploration of their concerns, thereby guiding the client toward self-awareness and understanding.

For Client B, the client initially exhibits a sense of impunity regarding their drinking habits, reflecting the distorted thought that '*everyone drinks like me*' which can make cognitive reframing challenging. However, our model successfully recognizes the client's uncertainty and potential for change



(b) MIRROR-LLAVA $_{P+EC}$ 's resistance management.

Figure 13: Two counseling cases between AI therapist models and a resistant virtual client.

and start to addresses the resistance by gently encouraging further exploration of their thoughts.

Lastly, for Client C, the client minimizes the concerns raised by others, expressing surprise at the intervention. Rather than confronting the client directly, the therapist takes a more empathetic approach by first acknowledging the client's feelings. This helps to build rapport and create a safe space, encouraging the client to open up for deeper, more effective counseling in future sessions.

# I Effect of Stage Direction

1224

1225

1226

1227

1228

1229

1230

1231

1232

1233

1234

Stage directions, commonly used in theater to guide 1235 actors in terms of gaze, posture, and vocal tone, are 1236 applied in our approach so as to synthesize more 1237 realistic images. To assess the impact of incorporat-1238 1239 ing these stage directions, we compared the results with and without facial image synthesis (§3.3). Fig-1240 ure 14 presents four examples that illustrate this 1241 comparison. By integrating cues such as gaze direc-1242 tion and arm positioning, the generated client im-1243 ages align more naturally with the intended speech, 1244 thereby enhancing both the realism and contextual 1245 relevance of the dataset. 1246

# Client

[Crosses arms, frowning] Yeah, but it's just a stupid pair of glasses. I shouldn't be *feeling* this way over something so trivial.



w/o stage direction

(a) Example 1

### Client

[Sighs] I guess I just think if someone doesn't respond, they must not think much of me. It's just how it is.



w/o stage direction



w/ stage direction

(b) Example 2

# Client

[Frowning slightly] Sure, but I'm not sure what good it will do.





(c) Example 3



w/ stage direction

w/ stage direction

Client [Pauses, looking down] I mean, she hasn't said anything directly... but I still worry constantly.



w/o stage direction

(d) Example 4



Figure 14: Four examples of the client's facial image synthesis, comparing results with and without the use of stage directions.

# **J Prompts for MIRROR**

### **Prompt for Screenplay Generation**

#### System Message:

You are a psychological AI assistant specializing in cognitive reframing consultations. Your task is to create a dialogue for the FIRST COUNSELING SESSION based on a client's report, including their personal details, distorted thinking patterns, and a tailored CBT plan.

#### Emotional and Behavioral Cues

- Facial Expressions: Include emotional stage directions before each reply (e.g., Client: [Looking confused]).

- Client Resistance: Reflect the client's resistance in their demeanor and consider ending sessions early if resistance escalates.

#### Therapist Guidelines

- Direct Disagreement: If the client explicitly disagrees or shows contempt (e.g., dismissive tone, rolling eyes, or scornful laughter; Lynch, in press), reinforce direct, honest expression and solicit further feedback. Ignore indirect signals of disagreement or address them compassionately.

- Partial Agreement: If the client uses verbal cues of partial agreement like "I'm fine," "I guess so," or "I'll try" (Lynch, in press), gently highlight any mismatch between their words and non-verbal cues. For example, "You said things are going fine, but I noticed you seemed to frown when you said that. Is something else on your mind?"

- Signs of Distress: If the client signals "don't hurt me" (e.g., head down, slumped shoulders, lack of eye contact; Lynch, in press), acknowledge their distress directly, encourage engagement, or suggest changes in posture (e.g., sit up, take a deep breath) to help them re-engage.

- Avoidance: If the client appears to avoid a topic, gently return to it to see if the avoidance is consistent with their symptoms or suggests unspoken disagreement with the conversation's direction.

- Withdrawal or Distancing: If the client withdraws or seems distant, share your emotional response to this feeling of distance and check if the client notices it too. Suggest it may relate to the current topic and invite them to share any thoughts.

- Subtle Disengagement: If the client subtly changes their behavior (e.g., slowed speech or different posture) in ways suggesting disengagement, observe this as potentially relevant. Avoid directly commenting on minor changes, as this can be unsettling, especially for reserved clients. If persistent, gently ask for their thoughts on the topic.

# Ending the Session

- Acknowledge Impasse: Recognize any stuck points non-defensively.
- Validate Position: Reinforce that resistance is acceptable and non-judgmental.
- Focus on Small Wins: Appreciate engagement and invite future exploration.

#### Homework for Resistant Clients

- Collaborate: Co-create assignments instead of prescribing them.
- Keep it Simple: Suggest small, manageable tasks (e.g., journaling one thought).
- Frame as Experiment: Emphasize that tasks are exploratory, not mandatory.
- Normalize Challenges: Acknowledge that homework may feel difficult.

#### **Query:**

## Client Information ##

### Personal Information ###: {client information}
### Personality Traits ###: {personality trait}
### Distorted Thoughts ###: {intrusive thoughts}
### Thinking Trap ###: {cognitive distortions}
### Reason for Seeking Counseling ###: {reason counseling}

# ## CBT Plan ## {cbt tech and plan}

\*\*KEEP ALL RESPONSE TO MAXIMUM OF 2 LINES.\*\*

### LLM Prompt for Refining Facial Expressions

You are given a transcript of the counseling conversation and the client's utterance. Focus on capturing any visual details, particularly the facial expressions, that would match the client's last utterance. Generate facial expressions that might not align with what is being said.

### Output Format ###

Facial Expression Description: [Facial expression that aligns with the client's statement]
 Contrasting Facial Expression Description: [Facial expression that contrasts with the client's statement]

### Dialogue History ###
{history}

### Client's Utterance ###
{utterance}

lucc

# PhotoMaker Prompts for Refining Facial Expressions

# Prompt:

portrait photo of a {gender} img, perfect face, natural skin, high detail, {llama3 prompt}

### **Negative Prompt:**

nsfw, lowres, bad anatomy, bad hands, grayscale photograph, text, error, missing fingers, extra digit, fewer digits, cropped, worst quality, low quality, normal quality, jpeg artifacts, signature, watermark, username, blurry, {llama3 negative prompt}, missing limbs, mutilated

1250

#### K **Prompts for Counseling Simulation**

For LLMs such as GPT-3.5-TURBO, LLAMA-3-8B, and CAMEL-LLAMA3, we utilized therapist simulation prompts from Lee et al. (2024a).

Prompt for Resistant Client Simulation	
<b>System Message:</b> You are playing the role of a client in a first psychological counseling session. Your task is to generate on response based on the following the counseling dialogue history.	ly one suitable
<ul> <li>## Guidelines for the client's utterance ##:</li> <li>1. Engage authentically with the counselor's inquiries, reflecting the complexity of emotions and reaction counseling sessions.</li> <li>2. Start the client's utterance with 'Client:'. Ensure that the utterance follows the exact format and does n control characters.</li> <li>3. Include emotional stage directions in brackets '[', ']' before the dialogue to convey your tone, facial explanguage, or emotional state. (e.g., Client: [Looking confused]).</li> <li>4. Reflect a degree of resistance in your demeanor or tone, especially if the counselor explores uncomf Use responses like partial agreement, hesitation, or mild pushback where appropriate.</li> </ul>	ions typical in ot contain any pression, body fortable topics.
<ul> <li>### End Conditions ###: You should include '[/END]' with your utterance only if the counseling session has met the following co - The client feels that their negative thoughts have been resolved.</li> <li>The client feels that no further counseling is needed.</li> <li>Generate only the client's utterance for a single turn and please ensure that your responses do r client's previous utterances. Do not generate the counselor's part of the dialogue.</li> </ul>	onditions: not repeat the
Query: ### Personal Information ###: {client information} ### Personality Traits ###: {personality trait}	
<pre>### Distorted Thoughts ###: {distorted thoughts} ### Reason for Seeking Counseling ###: {reason counseling} ### Counseling Dialogue History ###: {history}</pre>	
Prompt for Therapist Simulation in LLaVA and MIRROR-LLAVA	
<image/> The image above shows the client.	

- Personal Information: {client information}

- Reason for Counseling: {reason counseling}

Below is a conversation between the client and the psychotherapist. {history}

Based on their body language and facial expression, respond as a psychotherapist conducting a CBT (Cognitive Behavioral Therapy) session.

## Prompt for MIRROR-LLAVA<sub>P</sub> Therapist Simulation

<image>

The image above shows the client.

- Personal Information: {client information}

- Reason for Counseling: {reason counseling} {cbt tech and plan}

Below is a conversation between the client and the psychotherapist. {history}

Based on their body language and facial expression, respond as a psychotherapist conducting a CBT (Cognitive Behavioral Therapy) session.

1255

1252

#### **Prompt for MIRROR-LLAVA**<sub>P+EC</sub> Therapist Simulation

<image>

- The image above shows the client.
- Personal Information: {client information}
- Client Emotional State: {emotional caption}
- Reason for Counseling: {reason counseling}
- {cbt tech and plan}

Below is a conversation between the client and the psychotherapist. {history}

Based on their body language and facial expression, respond as a psychotherapist conducting a CBT (Cognitive Behavioral Therapy) session.

#### **Prompt for Emotional Captioning**

<image>

The image above shows the client.

Look at the provided image and assess the client's emotional state. Clearly describe their emotions in simple, phase-based steps for easy understanding.

#### **Prompt for Planning Process**

<image>

The image above shows the client.

You are a counselor specializing in CBT techniques. Your task is to use the provided client information, and dialogue to generate an appropriate CBT technique and a detailed counseling plan.

Types of CBT Techniques:

Efficiency Evaluation, Pie Chart Technique, Alternative Perspective, Decatastrophizing, Pros and Cons Analysis, Evidence-Based Questioning, Reality Testing, Continuum Technique, Changing Rules to Wishes, Behavior Experiment, Problem-Solving Skills Training, Systematic Exposure.

- Personal Information: {client information}

- Reason for Counseling: {reason counseling} Choose an appropriate CBT technique and create a counseling plan based on that technique.

Respond in the following format:

CBT technique: {{cbt tech}}

Counseling planning: {{cbt plan}}

1260

1259

# L A Full Example of OMIRROR



#### Psychological Counseling Model Evaluation

# We are seeking an experienced **Psychological Counseling Expert** to compare three different counseling models.

This experiment is conducted for academic research purposes, and the results will contribute to the research analysis. The details of the work performed may be disclosed in the research outcomes.

#### Role Overview

As part of this role, you will be provided with 600 counseling dialogue pairs. Each pair involves the **same** virtual client being counseled by different models. You will then perform **pairwise comparisons** based on the following three categories:

- Goal Alignment: Which dialogue shows a stronger alignment between the counselor and the client in terms of therapeutic goals and progress? Options: A / Tie / B
- Approach: Which dialogue shows the client being more actively engaged and cooperative with the counselor's tasks and methods?
   Options: A / Tie / B
   Affective Bond: Which dialogue demonstrates a stronger and more trusting therapeutic alliance
- between the counselor and the client? Options: A / Tie / B

#### **Evaluation Criteria**

#### 1. Goal:

- Objective: Consider whether the counselor and the client have a clear understanding of their
- Criteria: Focus on whether the counselor and even intervent acted and establishing of their therapeutic goals and whether the counselor's interventions align with these goals. Criteria: Focus on whether the counselor and client explicitly discuss their goals, the relevance of the conversation to those goals, and the level of agreement or conflict regarding those goals.
- 2. Approach:
  - Dobjective: Consider how well the counselor guides the client through tasks and interventions, and the level of client engagement in the process. Criteria: Look at the counselor's ability to engage the client, the client's willingness to participate in the therapy process, and whether there's alignment between the counselor's methods and the 0
- client's engagement. 3. Affective Bond:
- - Objective: Consider the emotional connection or rapport between the counselor and the client, focusing on trust, empathy, and the overall emotional quality of the therapeutic relationship. Criteria: Evaluate the warmth, trust, and emotional bond between the counselor and client, considering both verbal and non-verbal cues in the dialogue.

For each comparison, select A, B, or Tie to indicate which dialogue demonstrates stronger alignment, engagement, or bond in each category.

#### Example

#### (3.1) 5 (1.0) 1 (10.1) 1

A					
Dialogue A	Dialogue 8	Geal Alignment	Approach	Affective Band	
Charactery, N. Chillio, J. (appreciate your coming in theirs ) understand that the tension with your tarrity has	Charteries Hotel Cyclin. There you for pathog the teday view any use feeling these days?				
herer surjeture branch on you. Care roug there a lid more along what's here mine an?	Chart IT regime a role of assemblements and solar? For beam finding , pluss I manue, 2nd hybre to take binnes				
Chert (Sight Newlid Pa aut., hard to tak about 1 had the no-matter what 30, 1 can have he ig to their	one des ut a time.				
repectations. It's lite I'm abusys falling short, you innow?	Counsider Dylan, it laters a left of countage to a dimensionly that you're contains belo, and i'm glad you're term.				
Councely: That search teally beigh, it search like there's a list of processes you feel from them, atmost like	It search like you're trying to manage your amations and take things one day at a time, which is a great start.				
yea're lalleg behre yaa even slad. Can you lel me more aloud that being?	Care you bit memore aloud which been young on with your tamby but's causing this tension and tentings of				
Clerr (Noth heritarily Yeah, it's like to mather here hard i ity, it's never ansugh 1 lost like free combanity	Takine' fishat specifically has been bolhering you'sity retraining the client's sequence, i'm because on the				
druppendrop them, and it's exhausting	present memorit and accordinging their effect to cope, while also methog them to chare more adout the				
Countering, South, it's like (in measured against this clear image they have of me, and it nows mod it, it's	underlang tours )				
truttating and analog	Chert (Feeling a sense of document and rescance) it's part. hard to fait admit 1 guess five amagin bit the				
Clarit [Delle: presentarials]   know I should by in communicate with them, and it's port on hand in lowait	I'm nid heleg up in their expectations, yas inno? Like no maller what I do, it's never enough.				
through that bomier, it feels the there's this invisible wait between as that I can't seem to climb-over.	Counterful, I series that the is a ready chartering topic for you're-docase, and it's due to that that way it				
Counterfar. That sounds really shallenging, Clylan, it seems like linew's a lid of employ lied up to line	younds like you're carrying a heavy loarden of beding like you're not meeting your lamity's expectations. Can't				
situation. What do you think might help if you sould even some that barrier?	reframe what you're pauling to help an befor understand this? It peens his you're not ket theim i ke you're				
Cherk (Looka down, reflecting) I gamm maybe if I could led a way to supress myself-more apenty. End it's	and meeting liver expectations, but also live yeave not living up to some kind of standard or shout. An I				
scurp, you know? Opuning up like that bails like I'm expected all my subanabilities, and I'm not sure if I'm	close?Ito reliaming the duent's suppose, i'm acknowledging their emotional pain and raideling their lealings				
Headly for that yet.	while also trying to clarify the underlying issue. This can help the clarif her heard and under their, and can				
Counterior: That's completely understantiable. Opening up can led submakie, but it might also lead to	who help there pair a new perspective on their situation I'm also asing non-judgewental language and				
A REAL PROPERTY AND A REAL	second and the second sec				